

# Analog and Digital Electronics

Daniel Adam Steck

*Department of Physics, University of Oregon*



# **Analog and Digital Electronics**

Daniel Adam Steck

*Department of Physics, University of Oregon*

Copyright © 2015, by Daniel Adam Steck. All rights reserved.

This material may be distributed only subject to the terms and conditions set forth in the Open Publication License, v1.0 or later (the latest version is available at <http://www.opencontent.org/openpub/>). Distribution of substantively modified versions of this document is prohibited without the explicit permission of the copyright holder. Distribution of the work or derivative of the work in any standard (paper) book form is prohibited unless prior permission is obtained from the copyright holder.

Original revision posted 28 January 2015.

This is revision 1.0.9, 7 March 2023.

Cite this document as:

Daniel A. Steck, *Analog and Digital Electronics*, available online at <http://steck.us/teaching> (revision 1.0.9, 7 March 2023).

Author contact information:

Daniel Steck

Department of Physics

1274 University of Oregon

Eugene, Oregon 97403-1274

[dsteck@uoregon.edu](mailto:dsteck@uoregon.edu)



# Acknowledgements

For comments and corrections, thanks to:

- Harrison Allen–Sutter (U. Oregon)
- Patrick Baus (TU Darmstadt)
- Gary Bentley (<http://www.daltonbentley.com>)
- Cameron Dennis (U. Oregon)
- Kyle Eichenberger (U. Oregon)
- Austin Ferrie (U. Oregon)
- Cody Jarrett (U. Oregon)
- Eric Torrence (U. Oregon)



# Contents

<b>Acknowledgements</b>	<b>5</b>
<b>Contents</b>	<b>7</b>
<b>Using This Book</b>	<b>17</b>
<b>I Analog Electronics</b>	<b>21</b>
<b>1 Resistors</b>	<b>23</b>
1.1 Basic Definitions . . . . .	23
1.2 Ohm's Law . . . . .	23
1.2.1 Resistors . . . . .	24
1.3 Networks and Kirchoff's Laws . . . . .	24
1.3.1 Series Resistors . . . . .	25
1.3.2 Parallel Resistors . . . . .	26
1.3.3 Voltage Divider . . . . .	26
1.4 Thévenin's Theorem . . . . .	27
1.4.1 Voltage Divider . . . . .	28
1.4.2 Connected Circuits and Power Transfer . . . . .	29
1.5 Matrix Solution of Resistor Networks . . . . .	30
1.5.1 Review of Linear Algebra . . . . .	30
1.5.2 Matrix Form of the Resistance Network and Example . . . . .	31
1.5.3 Solution for the Effective Resistance . . . . .	32
1.5.4 Proof of Thévenin's Theorem . . . . .	33
1.5.4.1 Example Network . . . . .	35
1.5.4.2 Extension to Current Sources . . . . .	36
1.5.4.3 Norton's Theorem . . . . .	37
1.5.4.4 Superposition of Sources . . . . .	38
1.6 Circuit Practice . . . . .	39
1.6.1 Reflection-Symmetric Network . . . . .	39
1.6.2 Series and Parallel Light Bulbs . . . . .	39
1.6.3 Thévenin Circuit . . . . .	40
1.7 Exercises . . . . .	41
<b>2 Capacitors and Inductors</b>	<b>53</b>
2.1 Capacitor Basics . . . . .	53
2.2 Simple R–C Circuits . . . . .	54
2.2.1 Integrator . . . . .	54
2.2.1.1 Solution by Integrating Factor . . . . .	55
2.2.1.2 Constant Input: Exponential Charging . . . . .	55
2.2.1.3 Integration . . . . .	56

2.2.2	Differentiator . . . . .	56
2.3	AC Signals and Complex Notation . . . . .	57
2.3.1	Complex Phase . . . . .	57
2.3.2	Capacitive Reactance . . . . .	58
2.3.3	Inductive Reactance . . . . .	59
2.3.4	Impedance . . . . .	59
2.3.5	Low-Pass Filter . . . . .	59
2.3.6	Example Problem: Alternate Scaling . . . . .	61
2.3.7	Example Problem: High-Pass Filter . . . . .	61
2.4	Phase . . . . .	62
2.4.1	Example: Low-Pass Filter . . . . .	63
2.5	Power . . . . .	64
2.6	Resonant Circuits . . . . .	65
2.6.1	$Q$ Factor . . . . .	66
2.6.1.1	Fundamental Definition . . . . .	67
2.7	Coupled, Resonant Circuits . . . . .	68
2.8	Circuit Practice . . . . .	72
2.8.1	Tesla Coil . . . . .	72
2.9	Exercises . . . . .	73
<b>3</b>	<b>Diodes</b> . . . . .	<b>77</b>
3.1	Ideal Diode . . . . .	77
3.2	Vacuum Diodes . . . . .	77
3.3	Semiconductor Diodes . . . . .	78
3.4	Current–Voltage Characteristics . . . . .	80
3.4.1	Diode Law . . . . .	81
3.5	Zener Diodes . . . . .	82
3.6	Rectifier Circuits . . . . .	83
3.6.1	Half-Wave Rectifier . . . . .	83
3.6.2	Full-Wave Rectifier . . . . .	84
3.7	Schottky Diodes . . . . .	85
3.8	Circuit Practice . . . . .	86
3.8.1	Cockroft–Walton Multiplier . . . . .	86
3.9	Exercises . . . . .	88
<b>4</b>	<b>Bipolar Junction Transistors</b> . . . . .	<b>91</b>
4.1	Overview . . . . .	91
4.2	Usage . . . . .	92
4.3	Mechanism . . . . .	93
4.4	Packaging . . . . .	94
4.5	Transistor Switch . . . . .	95
4.5.1	Saturation Mode . . . . .	95
4.5.2	Forward-Active Mode . . . . .	96
4.5.3	Summary . . . . .	96
4.6	Emitter Follower . . . . .	96
4.6.1	Input and Output Impedance . . . . .	97
4.7	Transistor Current Source . . . . .	99
4.7.1	Compliance . . . . .	100
4.7.2	Bias Network . . . . .	100
4.8	Common-Emitter Amplifier . . . . .	102
4.9	Bias Network (AC Coupling) . . . . .	103
4.10	Transistor Differential Amplifier . . . . .	104
4.10.1	Differential-Only Input . . . . .	106

4.10.2	Common-Mode-Only Input . . . . .	106
4.10.3	General Input and Common-Mode Rejection . . . . .	107
4.10.4	Improving the Differential Amplifier . . . . .	107
4.11	Ebers–Moll Model . . . . .	108
4.11.1	Magnitudes . . . . .	109
4.11.2	Relation to $\beta$ . . . . .	109
4.11.3	Intrinsic Emitter Resistance . . . . .	110
4.11.4	Current Mirror . . . . .	110
4.11.4.1	Application to the Differential Amplifier . . . . .	111
4.11.5	Other Refinements to the Transistor Model . . . . .	112
4.11.5.1	Temperature Dependence of the Base–Emitter Voltage . . . . .	112
4.11.5.2	Early Effect . . . . .	113
4.11.5.3	Miller Effect . . . . .	113
4.11.5.4	Variation of $\beta$ . . . . .	114
4.12	Common-Emitter Amplifier Revisited: Emitter Degeneration . . . . .	114
4.13	Biasing the Push–Pull Pair . . . . .	116
4.14	Mathematical Modeling of DC BJT Behavior . . . . .	119
4.14.1	Ebers–Moll Model . . . . .	119
4.14.2	Ebers–Moll–Early Model . . . . .	121
4.14.3	Gummel–Poon Model . . . . .	121
4.15	Little-“ $h$ ” Notation . . . . .	123
4.16	Circuit Practice . . . . .	124
4.16.1	Transistor Switching an Inductive Load . . . . .	124
4.16.2	Joule Thief . . . . .	125
4.16.3	Solid-State Tesla Coil . . . . .	126
4.16.4	Eric Clapton Signature Stratocaster Preamplifier . . . . .	127
4.16.5	Op-Amp Internals . . . . .	127
4.17	Exercises . . . . .	130
<b>5</b>	<b>Field-Effect Transistors and Semiconductor Switching Devices</b>	<b>139</b>
5.1	JFET (Depletion-Mode FET) . . . . .	139
5.2	MOSFET (Enhancement-Mode FET) . . . . .	141
5.3	Quantitative FET Behavior . . . . .	142
5.3.1	Visualization . . . . .	143
5.4	Basic FET Circuits . . . . .	146
5.4.1	JFET Current Source . . . . .	146
5.4.2	JFET Source Follower . . . . .	146
5.4.3	JFET Common-Source Amplifier . . . . .	147
5.4.4	JFET Differential Amplifier . . . . .	149
5.4.5	MOSFET Analog Switch . . . . .	150
5.5	Thyristors . . . . .	151
5.5.1	SCRs . . . . .	151
5.5.1.1	Example: Latching Switch with Power-Supply Fault Protection . . . . .	153
5.5.2	DIACs and TRIACs . . . . .	153
5.5.2.1	Light Dimmer . . . . .	154
5.6	IGBTs: Switching Very Large Voltages and Currents . . . . .	155
5.6.1	Driver Circuitry . . . . .	157
5.6.2	Inverter Circuits . . . . .	158
5.6.2.1	Inverter-Based Tesla Coil . . . . .	160
5.7	Circuit Practice . . . . .	162
5.7.1	Touch Switch . . . . .	162
5.7.2	JFET Op-Amp Internals . . . . .	163
5.8	Exercises . . . . .	165

<b>6</b>	<b>Vacuum Tubes</b>	<b>167</b>
6.1	Vacuum Diodes	167
6.1.1	Child–Langmuir Law	168
6.1.2	Vacuum Full-Wave Rectifier	168
6.2	Vacuum Triodes	169
6.2.1	Triode Voltage Amplifier (Common-Cathode Amplifier)	170
6.2.1.1	DC Bias	171
6.2.1.2	Naïve AC Analysis	172
6.2.1.3	Proper AC Analysis: Plate Resistance	173
6.2.2	Design Example: 12AX7 Preamplifier	175
6.2.3	Phenomenological Triode Model	175
6.3	Vacuum Tetrodes	177
6.4	Vacuum Pentodes and Beam Power Tubes	178
6.4.1	Design Example: 6V6 Power Output Stage	180
6.4.2	Triode Connection	183
6.5	Circuit Practice: Simple Tube Amplifier	184
6.5.1	First Preamplifier Stage	185
6.5.2	Second Preamplifier Stage	186
6.5.3	Power Amplifier Stage	186
6.5.4	Feedback Loop	186
6.5.5	Power Supply	187
<b>7</b>	<b>Operational Amplifiers</b>	<b>189</b>
7.1	Op-Amp Basics	189
7.1.1	Usage: Open-Loop	189
7.1.2	Usage: Closed-Loop	190
7.2	Op-Amp “Golden Rules”	190
7.3	Basic Op-Amp Circuits	190
7.3.1	Unity-Gain Buffer/Follower	190
7.3.2	Inverting Amplifier	191
7.3.2.1	Stability	192
7.3.3	Noninverting Amplifier	193
7.3.4	Summing (Inverting) Amplifier	193
7.3.5	Circuit Practice: Differential Amplifier	194
7.4	Op-Amp Filters	194
7.4.1	Op-Amp Differentiator	195
7.4.2	Op-Amp Integrator	195
7.4.3	Differentiator Issues	196
7.4.4	Integrator Issues	197
7.4.5	Sources of Integrator Error	198
7.4.5.1	Input Bias Current	198
7.4.5.2	Input Offset Voltage	200
7.4.6	Integrator Applications	200
7.5	Instrumentation Amplifiers	201
7.5.1	“Classic” Instrumentation Amplifier	202
7.5.2	Instrumentation-Amplifier Applications	203
7.5.2.1	Thermocouple Amplifier	203
7.5.2.2	Differential Transmission for Noise Rejection	204
7.5.2.3	AC-Coupled Inputs with High Impedance	205
7.6	Practical Considerations	206
7.6.1	Input-Bias Currents and Precision Amplifiers	206
7.6.1.1	Inverting Amplifier	206
7.6.1.2	Balanced Input-Impedances: Inverting Amplifier	207

7.6.1.3	Balanced Input-Impedances: Noninverting Amplifier . . . . .	207
7.6.1.4	Input Offset Currents . . . . .	208
7.6.1.5	Common-Mode Rejection Ratio . . . . .	208
7.6.2	Power Supplies . . . . .	208
7.6.2.1	Power-Supply Rejection . . . . .	209
7.6.2.2	Power-Supply Bypass Capacitors . . . . .	209
7.7	Finite-Gain Analysis . . . . .	211
7.7.1	Noninverting Amplifier . . . . .	211
7.7.1.1	Gain Limits and Error . . . . .	212
7.7.1.2	Insensitivity to Gain Variation . . . . .	213
7.7.2	Feedback and Input Impedance . . . . .	213
7.7.3	Feedback and Output Impedance . . . . .	214
7.7.4	Circuit Practice: Finite Gain in the Inverting Amplifier . . . . .	215
7.8	Bandwidth . . . . .	216
7.8.1	Slew Rate . . . . .	217
7.8.1.1	Slew Rate and Power-Boosted Op-Amps . . . . .	217
7.8.2	Stability and Compensation . . . . .	219
7.8.2.1	Op-Amp Output and Capacitive Loads . . . . .	220
7.9	Comparators . . . . .	221
7.9.1	Schmitt Trigger . . . . .	222
7.10	Positive Feedback and Oscillator Circuits . . . . .	223
7.10.1	Relaxation Oscillator . . . . .	223
7.10.2	Buffered Phase-Shift Oscillator . . . . .	224
7.11	Amplifier Noise Analysis . . . . .	225
7.11.1	Sources of Noise . . . . .	226
7.11.1.1	Johnson–Nyquist Noise . . . . .	226
7.11.1.2	“ $kT/C$ ” Noise . . . . .	228
7.11.1.3	Shot Noise . . . . .	228
7.11.1.4	$1/f$ Noise . . . . .	228
7.11.2	Estimation of Amplifier Noise . . . . .	229
7.11.2.1	Example: Precision Current Source . . . . .	231
7.12	Circuit Practice . . . . .	235
7.12.1	Analog Computers . . . . .	235
7.12.1.1	Proportional–Integral Amplifier . . . . .	235
7.12.1.2	Damped Harmonic Oscillator . . . . .	235
7.12.2	Gyrator . . . . .	236
7.12.3	Guitar Preamp with Midrange Boost/Cut . . . . .	238
7.12.4	Active Rectifiers . . . . .	241
7.12.5	Pulse-Area Stabilizer . . . . .	242
7.13	Exercises . . . . .	244
<b>8</b>	<b>PID Control</b> . . . . .	<b>261</b>
8.1	Basics of Linear Control . . . . .	261
8.2	Example: First-Order Plant, Proportional Control . . . . .	262
8.2.1	General Result: Closed-Loop Transfer Function . . . . .	262
8.2.2	Frequency-Domain Solution . . . . .	263
8.2.3	Time-Domain Solution . . . . .	263
8.2.4	Constant Input and Proportional Droop . . . . .	263
8.3	Integral Control . . . . .	264
8.3.1	Example: First-Order Plant, Integral Control . . . . .	264
8.3.2	Frequency Domain . . . . .	264
8.4	Proportional–Integral (PI) Control . . . . .	265
8.5	Proportional–Integral–Derivative (PID) Control . . . . .	265



<b>II</b>	<b>Digital Electronics</b>	<b>267</b>
<b>9</b>	<b>Binary Logic and Logic Gates</b>	<b>269</b>
9.1	Binary Logic . . . . .	269
9.2	Binary Arithmetic . . . . .	269
9.2.1	Unsigned Integers . . . . .	269
9.2.1.1	Binary-Coded Decimal . . . . .	270
9.2.1.2	Hexadecimal . . . . .	270
9.2.2	Negative Values and Sign Conventions . . . . .	270
9.2.2.1	Sign-Magnitude Convention . . . . .	270
9.2.2.2	Two's Complement . . . . .	270
9.3	Logic Gates . . . . .	271
9.3.1	One-Input Gates . . . . .	271
9.3.2	Two-Input Gates . . . . .	272
9.3.2.1	AND and NAND . . . . .	272
9.3.2.2	OR and NOR . . . . .	272
9.3.2.3	Universal Gates . . . . .	272
9.3.2.4	XOR and XNOR . . . . .	272
9.3.3	More Complex Gates . . . . .	273
9.4	Circuit Practice . . . . .	273
9.5	Exercises . . . . .	274
<b>10</b>	<b>Boolean Algebra</b>	<b>277</b>
10.1	Algebras and Boolean Algebra . . . . .	277
10.2	Boolean-Algebraic Theorems and Manipulations . . . . .	278
10.2.1	De Morgan's Theorems . . . . .	278
10.2.2	Absorption Theorems . . . . .	278
10.2.3	Another Theorem . . . . .	278
10.2.4	Example: XOR Gate . . . . .	278
10.2.4.1	NAND-Gate Realization . . . . .	279
10.2.5	Example: Algebraic Simplification . . . . .	279
10.3	Karnaugh Maps . . . . .	280
10.3.1	Three-Input Example . . . . .	280
10.3.2	Four-Input Example . . . . .	282
10.3.3	XOR Example . . . . .	282
10.3.4	Race Hazards . . . . .	282
10.4	Circuit Practice . . . . .	284
10.4.1	Boolean-Algebra Theorems . . . . .	284
10.4.2	Karnaugh Map . . . . .	284
10.5	Exercises . . . . .	286
<b>11</b>	<b>Physical Implementation of Logic Gates</b>	<b>291</b>
11.1	Simple Mechanical Switches . . . . .	291
11.2	Diode Logic (DL) . . . . .	292
11.2.1	Diode Review . . . . .	292
11.2.2	DL AND Gate . . . . .	293
11.2.3	DL OR Gate . . . . .	293
11.3	Resistor-Transistor Logic (RTL) . . . . .	294
11.3.1	BJT Review . . . . .	294
11.3.2	RTL NOT Gate . . . . .	294
11.3.3	RTL NOR Gate . . . . .	295
11.4	The Real Thing: Transistor-Transistor Logic (TTL) . . . . .	295
11.4.1	TTL Nomenclature . . . . .	297

11.5	The Modern Thing: CMOS Logic . . . . .	297
11.5.1	MOSFET Review . . . . .	297
11.5.2	NMOS . . . . .	298
11.5.3	PMOS . . . . .	299
11.5.4	CMOS . . . . .	300
11.5.5	CMOS NAND and NOR . . . . .	302
11.5.6	CMOS Nomenclature . . . . .	303
11.6	Circuit Practice . . . . .	303
11.6.1	Mystery RTL/DL Gate . . . . .	303
11.6.2	CMOS Level-Shifting Buffers . . . . .	304
11.7	Exercises . . . . .	306
<b>12</b>	<b>Multiplexers and Demultiplexers</b>	<b>309</b>
12.1	Multiplexers . . . . .	309
12.1.1	Example: 74151 . . . . .	309
12.2	Demultiplexers . . . . .	310
12.2.1	Example: 74138 . . . . .	310
12.3	Making a MUX . . . . .	310
12.4	Expanding a MUX (or DEMUX) . . . . .	311
12.5	Analog MUX/DEMUX . . . . .	312
12.6	Circuit Practice: Multiplexed Thermocouple Monitor . . . . .	312
12.7	Exercises . . . . .	316
<b>13</b>	<b>Flip Flops</b>	<b>319</b>
13.1	Flip-Flop Construction: SR Flip Flop . . . . .	319
13.1.1	Application: Debounced Switch . . . . .	320
13.2	Clocked Flip-Flops . . . . .	321
13.2.1	D-Type Flip-Flop . . . . .	321
13.2.2	Edge-Triggered, D-Type Flip-Flop . . . . .	322
13.2.3	JK Flip-Flop (Edge-Triggered) . . . . .	322
13.3	Counters . . . . .	323
13.3.1	Asynchronous (Ripple) Counter . . . . .	324
13.4	Memory and Registers . . . . .	324
13.4.1	Register . . . . .	324
13.4.2	Shift Register . . . . .	325
13.5	Sequential Logic and the State Machine . . . . .	325
13.5.1	Example: Synchronous, Divide-by-3 Counter . . . . .	326
13.5.2	State Diagrams . . . . .	327
13.6	Memory . . . . .	328
13.6.1	Example: 8×1-bit RAM . . . . .	329
13.6.2	Example: 6116 SRAM . . . . .	330
13.6.3	Other Memory Types . . . . .	331
13.7	State Machines with Memory . . . . .	331
13.7.1	Example: Divide-by-3-With-Hold Counter . . . . .	332
13.7.2	General Considerations: Towards a Microprocessor . . . . .	333
13.7.3	Programmable ROM as Logic . . . . .	334
13.7.4	Programmable Logic Devices . . . . .	335
13.8	Circuit Practice . . . . .	336
13.8.1	Basic Flip-Flops . . . . .	336
13.8.2	Pulse-Area Stabilizer . . . . .	337
13.8.3	Circuit Practice: Backwards Divide-by-3 Counter . . . . .	339
13.8.4	Memory: RAM vs. ROM . . . . .	339
13.8.5	Circuit Practice: Divide-by-2-or-3 Counter . . . . .	339

13.9 Exercises . . . . .	341
<b>14 Comparators</b>	<b>353</b>
14.1 Overview and Review . . . . .	353
14.1.1 Example: TL3016 . . . . .	353
14.2 Open-Collector Output . . . . .	354
14.3 Schmitt Trigger . . . . .	355
14.3.1 Example: Analog-to-Digital Clock-Signal Conversion . . . . .	357
14.4 Circuit Practice . . . . .	358
14.5 Exercises . . . . .	360
<b>15 Pulse and Waveform Generation</b>	<b>361</b>
15.1 The Classic 555 Timer . . . . .	361
15.1.1 Equivalent Circuit . . . . .	361
15.1.2 Astable Multivibrator . . . . .	362
15.1.2.1 Frequency Modulation . . . . .	364
15.1.2.2 Pulse-Width Modulation: LED Dimmer . . . . .	364
15.2 Monostable Multivibrators . . . . .	365
15.2.1 555 as a One-Shot . . . . .	365
15.2.1.1 The 74121 . . . . .	366
15.2.1.2 Combining One-Shots: Pulse Delay . . . . .	367
15.3 Circuit Practice . . . . .	367
15.3.1 Duty-Cycle Control . . . . .	367
15.3.2 Astable Multivibrator . . . . .	368
15.4 Exercises . . . . .	370
<b>16 Digital–Analog Interfaces</b>	<b>373</b>
16.1 Digital-to-Analog Conversion . . . . .	373
16.1.1 Resolution . . . . .	373
16.1.2 DAC Circuitry . . . . .	374
16.1.3 R–2R Ladder . . . . .	375
16.2 Analog-to-Digital Conversion . . . . .	376
16.2.1 Flash ADC . . . . .	376
16.2.2 Successive Approximation . . . . .	377
16.2.3 Single/Dual-Slope ADC . . . . .	379
16.3 Circuit Practice . . . . .	380
16.3.1 Computer-Interface DAC Controller . . . . .	380
16.3.2 3-Bit ADC . . . . .	380
16.4 Exercises . . . . .	381
<b>17 Phase-Locked Loops</b>	<b>385</b>
17.1 Frequency Multiplier . . . . .	385
17.1.1 Feedback Loop . . . . .	388
17.2 Example PLL . . . . .	388
17.3 Other Applications . . . . .	389
17.3.1 FM Demodulation . . . . .	389
17.3.2 Direct Digital Synthesis . . . . .	389
17.4 Dynamical Model . . . . .	389
17.4.1 Equation of Motion . . . . .	390
17.4.2 Damping . . . . .	390
17.5 Exercises . . . . .	393

<b>Appendices</b>	<b>395</b>
<b>Appendix A Homemade Printed Circuit Boards</b>	<b>397</b>
A.1 LED Blinker Circuit . . . . .	397
A.2 PCB Fabrication . . . . .	398
A.2.1 LED-Blinker PCB . . . . .	398
A.2.2 Presensitized PCBs . . . . .	400
A.2.3 Cutting to Size . . . . .	400
A.2.4 Exposure . . . . .	400
A.2.5 Developing . . . . .	401
A.2.6 Etching . . . . .	402
A.2.7 Drilling . . . . .	403
A.2.8 Cleaning Up the Board Dimensions . . . . .	403
A.2.9 Preparing to Solder . . . . .	403
A.2.10 Soldering . . . . .	404
A.2.11 Stuffing in Stages . . . . .	404
A.2.12 Protection . . . . .	405
A.3 Characterizing the Circuit . . . . .	405
<b>Appendix B State-Machine Emulation of the 74LS502/3 Successive-Approximation Register</b>	<b>407</b>
B.1 74502 . . . . .	407
B.1.1 General Emulation Notes . . . . .	408
B.1.2 22V10 Emulation . . . . .	409
B.2 74503 . . . . .	412
B.2.1 ATF750C Emulation . . . . .	413
B.3 Testing the State Machines . . . . .	416
<b>Appendix C Gallery of Characteristic Curves</b>	<b>421</b>
C.1 Diodes . . . . .	422
C.1.1 Small-Signal Diodes . . . . .	422
C.1.2 Medium-Current Rectifiers . . . . .	423
C.1.3 High-Current Rectifiers . . . . .	424
C.2 Bipolar Junction Transistors . . . . .	425
C.2.1 Small-Signal BJTs . . . . .	425
C.2.2 Power BJTs . . . . .	427
C.3 MOSFETs . . . . .	429
C.3.1 Small-Signal MOSFETs . . . . .	429
<b>Index</b>	<b>431</b>



# Using This Book

## Assumptions and Background

Electronics, as far as an intermediate-level area of study for physics students, is unlike other traditional physics topics in that the level of mathematical sophistication involved is relatively low—things like  $V = IR$  are pretty standard fare. However, being *good* at electronics requires the ability to think intuitively, breaking complex problems down into manageable bits. Studying electronics is also a great way to develop this way of thinking. Of course, it's also a useful addition to your bag of tricks if you need to get stuff done in the lab.

Now more than ever, the modern state of electronics is a fast-moving field, with classic (not to mention *useful*) parts continually being “end-of-lifed.” The main focus here is on covering the fundamentals, rather than to try to keep up with the latest fashion. Studying classic circuits and components is useful in analyzing an enormous library of existing circuits and instruments. And once you master the basics, incorporating more modern components is a matter of relating them to the ones you already know. (You'll *always* have to do this, no matter how “current” your training is.)

This book contains more than enough material for two ten-week quarters, one on analog electronics, and one on digital electronics; the level is aimed at senior physics undergraduate students, who have had some exposure to resistor, capacitors, and inductors in a freshman-level sequence in introductory physics (preferably calculus-based). Some experience with differential equations is handy, too.

## Essential Topics: Analog Electronics

Among the topics covered here, there is a subset that is more or less essential for a course—leftover topics can be incorporated in the time left. On the analog side, the essential topics are below. (A non-essential topic, for example, would be vacuum tubes: These days you can get pretty far in life not knowing how a vacuum tube works. But they're still a lot of fun.)

- **Chapter 1: Resistors and Networks.** Sections 1.1–1.4.2 (pp. 23–30).
- **Chapter 2: Capacitors and Inductors.** Sections 2.1–2.6.1.1 (pp. 53–68).
- **Chapter 3: Diodes.** Sections 3.1–3.6.2 (pp. 77–85).
- **Chapter 4: Bipolar Transistors.** Sections 4.1–4.11.4.1 (pp. 91–112).
- **Chapter 5: Field-Effect Transistors.** Sections 5–5.4.5 (pp. 139–151).
- **Chapter 7: Op-Amps.** Sections 7.1–7.10.2 (pp. 189–225).

In the list above, some of the topics included can even be considered only marginally essential. For example, JFETs could be covered in somewhat less depth (though MOSFETs are quite handy).

## Essential Topics: Digital Electronics

On the digital side, the emphasis on old-school components and techniques is more evident. These days, much of the electronic “heavy lifting” is done on processors and microcontrollers, but programming these is

more a matter of programming than electronics (though they are well worth studying). The main philosophy here is that it's good to get some experience with the various basic types of chips. Anything more modern or complicated is then not hard to tackle, once armed with this background and the data sheet. In this sense, digital electronics is relatively simple compared to analog electronics. In analog circuit design, there is a *lot* to worry about: impedance, loading, bandwidth, distortion, noise, etc. Digital electronics is often fairly close to the ideal of 1's and 0's. The compensation for this, of course, is that digital circuits are often way more complicated, involving many parallel signals.

- **Chapter 9: Binary Logic and Logic Gates.** The whole chapter (pp. 269–273).
- **Chapter 10: Boolean Algebra.** The whole chapter (pp. 277–283).
- **Chapter 11: Physical Implementation of Logic Gates.** The whole chapter (pp. 291–303).
- **Chapter 12: Multiplexers and Demultiplexers.** The whole chapter (pp. 309–312).
- **Chapter 13: Flip Flops.** The whole chapter (pp. 319–336).
- **Chapter 14: Comparators.** The whole chapter (pp. 353–358).
- **Chapter 15: Pulse and Waveform Generation.** The whole chapter (pp. 361–367).
- **Chapter 16: Digital–Analog Interfaces.** The whole chapter (pp. 373–380).

Some of the essential topics here also straddle the line between analog and digital electronics, and in this sense it's generally better to study analog electronics *first*, and *then* move on to digital. Digital circuits are fundamentally made up of analog components, of course, and so some analog considerations are necessary particularly in understanding the limits of digital circuits (as well as other more obvious areas, such as the important case of interfacing analog and digital signals).

## Circuit Practice

Each chapter contains a “Circuit Practice” section near the end. These are exercises that are useful to cement the material soon after going through it, and before tackling more difficult problems. The exercises range from useful basic problems (slight variations on problems considered in the main text) to more complex, real-world circuits that serve as “reading exercises” (to, for example, recognize how some of the building blocks from the main text are used in real circuits, but also to get some experience in analyzing circuits that you might encounter inside or outside the laboratory). In some cases, where the practice problem is particularly urgent, the circuit-practice problem may appear earlier in the chapter. In these cases you should complete the problem *before* proceeding with the following material.

In the classroom setting, these are particularly good as end-of-class exercises for the students to work through as review. The instructor can interact with the students individually as they work through the problems, allowing for valuable feedback on how well the students are absorbing the material.

## Exercises

It can be difficult to find electronics exercises suitable for the advanced physics undergraduate. There is a multitude of electronic-engineering texts, with a multitude of exercises in tow, but many of the problems are thinly disguised copies of the same simple circuits with different numbers. Advanced physics undergrads are generally pretty capable of plugging in different numbers in a familiar problems, so the emphasis here is on more substantial and sophisticated variations on the basics covered in the main text. The problems range from fairly straightforward variations to some relatively involved problems. On the analog side, the problems are mostly restricted to circuit analysis, which is difficult enough; on the digital side, where life is more straightforward, there are more circuit-design problems. Many problems here are quite “practical,” in that they involve analyzing or designing useful circuit examples; but some are also useful in developing ways of thinking like a physicist. This follows on the philosophy of some of the material in the main text—the



proof of Thévenin's theorem is not particularly useful in circuit design, for example, but it's a useful example of the mathematical structure of resistor networks and of the analysis of constrained static systems.

Hopefully this book serves as a reasonably good reservoir of “physics-style” electronics exercises.

## Further Reading

This book is far from complete, and there are many options for augmenting this material and filling in the gaps. A couple of standouts are worth mentioning, however.

- Paul Horowitz and Winfield Hill, *The Art of Electronics*, 3rd ed. (Cambridge, 2015) (<https://artofelectronics.net>) (ISBN: 9780521809269).
- Dennis Barnaal, *Analog Electronics for Scientific Application* (Breton, 1982) (ISBN: 0534010156); also Dennis Barnaal, *Digital Electronics for Scientific Application* (Waveland, 1982) (ISBN: 0881334219).

Horowitz and Hill is the towering classic of electronics, particularly among physicists. If you're going to work with or design circuits, or just want to be serious about electronics, this book needs to be at the ready on your bookshelf. It covers in great detail all of the things you need to worry about when designing a circuit, and it also acts as a nice cookbook for useful, practical circuit ideas. The authors have made efforts to make it a “beginner-friendly” book, but because it contains so much good information, and it moves quickly and intuitively through many topics, it can be something of an intimidating firehose on the first pass. (It can be better with a skilled teacher telling you what to read, and more importantly what *not* to read.) It's now in the 3rd edition, but older 2nd editions are worth having too for the collections of circuit ideas, and more importantly the “bad circuits” collections, which serve as a good barometer of the current state of your electronics expertise (for the 3rd edition they have been relegated to the web site, but so far it's just not the same).

The Barnaal books, though now somewhat dated, are well-written and work well as introductory books. The level of the discourse and especially exercises is more appropriate for sophomore-level physics students, so it's somewhat easier than what you'll find right here. Being out of print for a while, these books can usually be had at a bargain, and are well worth picking up.

## Other Hyperlinks and Navigating this Document

To make it easier to access information within and beyond this document, there are many hyperlinks throughout. To keep the document “pretty,” the hyperlinks are not highlighted in colors or boxes by default. Some of the more obvious ones are the spelled-out URLs, like <http://steck.us>, are clickable, as are the QR codes mentioned above. Some of the less obvious ones are:

- DOI (document object identifier) codes in article citations (for locating articles online).
- ISBN codes in book citations (these resolve to pages on [amazon.com](https://www.amazon.com)).
- Section titles in the Contents.
- Page numbers in the Index.
- Chapter, section, page, and equation numbers throughout the text.



**Part I**

# **Analog Electronics**



# Chapter 1

## Resistors

### 1.1 Basic Definitions

Here, we’re going to breeze through a few fundamental notions in electromagnetism. At the most basic level, electronics studies the flow of “stuff,” or more specifically, **charge** (measured in **Coulombs**). The *flow* of charge is **current** (*not* “amperage”), defined by

$$I = \frac{dQ}{dt}. \quad (1.1)$$

(current)

What causes charge to move around and form currents? It’s the **potential** associated with an **electric field**. To move a charge between two points, say  $A$  and  $B$ , this requires some work (energy)  $W$  done against the force due to the field. Then the **potential difference** or **voltage difference** is

$$V_{AB} = \frac{W}{Q}. \quad (1.2)$$

That is, the work is proportional to the charge and to the difference in potential between the two points:

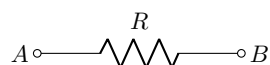
$$V_{AB} := V_A - V_B. \quad (1.3)$$

For a *static* electric field, it turns out that  $V_{AB}$  is independent of the path that the charge takes between the points, so we can represent this as a simple difference between the endpoint potentials. It’s important to note that only *differences* in potential matter: if we raise both  $V_A$  and  $V_B$  by the same amount, the work to transport the charge isn’t affected. Finally, note that voltage/potential is measured in **volts (V)**, which is the same as joules per coulomb (J/C), as we can see from the work relation (1.2).

An **electromotive force (EMF)** is a special name for a voltage difference due to an energy source (say, a battery, or a power supply).

### 1.2 Ohm’s Law

Since electric fields exert forces on charges, you might think that a constant electric field makes a charge move *ballistically*, or like a mass moves under constant gravity. But charges (electrons) moving through a material (metal, semiconductor), due to interactions with the material, quickly settle into a **terminal velocity**, like a particle falling through air under gravity. Under these conditions, and for small voltages, the velocity of the charges is proportional to the applied voltage, so the *current* is proportional to the voltage. This is the content of **Ohm’s law**. Consider a *resistor* (essentially any conducting material, say a wire, where the material “resists” the flow of charge), which we represent by the following schematic symbol:



Here the **resistance** (measured in **ohms** or  $\Omega$ ) is  $R$ , and the resistor connects points  $A$  and  $B$ . Then Ohm's law states

$$V_{AB} = IR. \quad (1.4)$$

(Ohm's law)

That is, for a fixed voltage, the current is inversely proportional to the resistance, which is sensible.

The voltage (and hence, electric field) does work on the charges. The **power** is the rate of work, or

$$P = \frac{dW}{dt}. \quad (1.5)$$

From Eq. (1.2),  $W = VQ$  (dropping subscripts on  $V$ ), and at constant voltage,

$$P = V \frac{dQ}{dt} = VI. \quad (1.6)$$

(electrical power)

Of course, with Ohm's law, we can also write

$$P = VI = I^2 R = \frac{V^2}{R} \quad (1.7)$$

(electrical power)

for a few useful alternate forms of the electrical power

### 1.2.1 Resistors

Essentially anything short of a superconductor has resistance. Wires that carry current have resistance, but usually it's desirable to keep their resistance small. But in virtually all electronic circuits, it's useful to introduce controlled quantities of resistance, and these are the electrical components we call resistors. A few basic types are:

1. **wirewound resistors**: are just wires wrapped around a form. These are usually expensive, but can be precise (using thin wire to make a large resistor) or able to handle high power (using thick wire embedded in ceramic).
2. **carbon film**: are a thin layer of carbon deposited on some insulating form (usually a small cylinder). They're cheap, but not particularly accurate in value. In the type with axial leads, these are usually recognizable by their tan color.
3. **metal film**: are a thin layer of metal deposited on some insulating form (again, usually a small cylinder). They're more expensive than carbon, but more accurate. In the type with axial leads, these are usually recognizable by a blue or blue/green color.

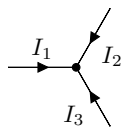
## 1.3 Networks and Kirchoff's Laws

A simple **circuit** is any network of resistors, batteries, and wires (later to include more stuff!). There are two basic laws that govern the circuit if all we have is batteries and resistors, and these are called **Kirchoff's laws**, one for current, and one for voltage.

1. The **current law** states that at any junction, the current going in to the junction must exactly balance the current going out, for charge not to accumulate there:

$$\sum_j I_{\text{in},j} = \sum_j I_{\text{out},j}. \quad (1.8)$$

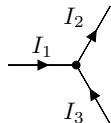
We can also keep track of the sense of “in” and “out” by keeping track of the *sign* of the current (always important to do!). Thus, for example, in the following junction,



we should write

$$I_1 + I_2 + I_3 = 0. \quad (1.9)$$

On the other hand, if we draw the currents like *this*,



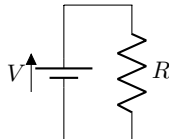
then we should write

$$I_1 - I_2 + I_3 = 0. \quad (1.10)$$

2. The **voltage law** states that around a *closed circuit* or *loop* in a circuit, the EMFs must balance the voltage drops, or

$$\sum_j \mathcal{E}_j = \sum_j V_j, \quad (1.11)$$

where the  $\mathcal{E}_j$  are the EMFs, and the  $V_j$  are the voltage drops. For example, in the circuit below, the EMF is  $V$  due to the battery, so the voltage drop across the resistor must also be  $V$ .



### 1.3.1 Series Resistors

Two resistors in series are the same as a single resistor, with an effective resistance that is the sum of the individual resistors.

$$\text{---} \underbrace{\text{zigzag}}_{R_1} \text{---} \underbrace{\text{zigzag}}_{R_2} \text{---} = \text{---} \underbrace{\text{zigzag}}_{R_{\text{eff}} = R_1 + R_2} \text{---}$$

Why? (Try to work this out on your own!)

The idea is that any current  $I$  that flows through one must flow through the other. So the voltages across the resistors are  $V_1 = IR_1$  and  $V_2 = IR_2$ . Then the total drop across both resistors is

$$V = V_1 + V_2 = I(R_1 + R_2) =: IR_{\text{eff}}. \quad (1.12)$$

Thus,

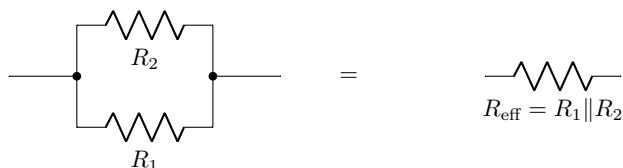
$$R_{\text{eff}} = R_1 + R_2 \quad (1.13)$$

is the effective resistance of the series pair. Of course, this generalizes to multiple resistors.



### 1.3.2 Parallel Resistors

Two resistors in parallel also behave as a single resistor, as shown below.



The shorthand notation here is

$$\frac{1}{R_1 \parallel R_2} := \frac{1}{R_1} + \frac{1}{R_2}, \quad (1.14)$$

so that more parallelism in resistors *decreases* resistance.

Why is this? In this case, the *voltage*  $V$  is common to the two resistors. The currents are  $I_1 = V/R_1$  and  $I_2 = V/R_2$ . But the total current through the pair must be  $I = I_1 + I_2$ , which satisfies  $I = V/R_{\text{eff}}$ . This means

$$\frac{V}{R_{\text{eff}}} = \frac{V}{R_1} + \frac{V}{R_2}, \quad (1.15)$$

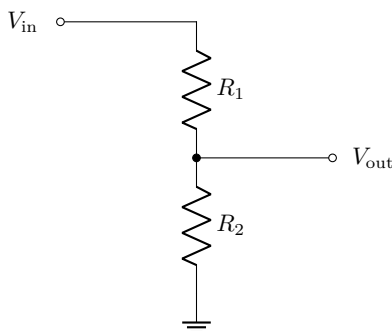
and canceling the voltage gives

$$\frac{1}{R_{\text{eff}}} = \frac{1}{R_1} + \frac{1}{R_2}, \quad (1.16)$$

which agrees with the shorthand above.

### 1.3.3 Voltage Divider

This is a useful combination of resistors that occurs *all the time* in circuits.



Here at the bottom of the circuit diagram, we are drawing the **ground** symbol, which means we are declaring this point to be a *fixed* voltage (say, zero), and all other voltages are *differences* with respect to ground. (In the “ground” or “earth” pin on an ac power receptacle, this is *literally* the ground outside. Often the case of electrical devices is connected to ground for safety, and in cars the entire chassis is ground.)

Now due to the input voltage  $V_{\text{in}}$ , some current flows in from the input. This must satisfy

$$I = \frac{V_{\text{in}}}{R_1 + R_2}. \quad (1.17)$$

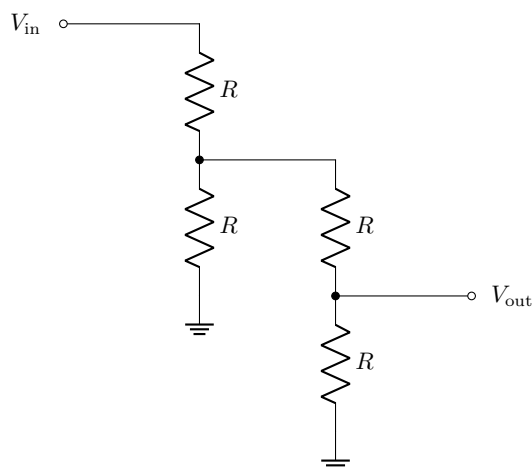
The same current flows through  $R_2$ , so

$$V_{\text{out}} = IR_2 = \left( \frac{R_2}{R_1 + R_2} \right) V_{\text{in}}. \quad (1.18)$$

(voltage divider)

The output voltage is reduced from the input by the ratio of  $R_2$  to the total resistance. **This is important; you should memorize this.** Especially if this is made from an adjustable resistor (**potentiometer**), this can be used to make an adjustable voltage source.

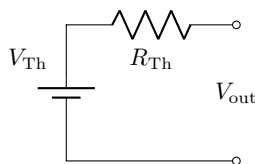
However, suppose we *chain* two voltage dividers, as follows.



Is the output voltage just divided twice, or  $1/4$  the input voltage? Why *not*?

## 1.4 Thévenin's Theorem

Thévenin's theorem<sup>1</sup> is very useful in analyzing passive networks. It says that **any** network of resistors and EMFs—if we interact with it only at two points—can be replaced by an equivalent circuit of a series EMF and resistor, as shown.



The EMF and resistance are called the **Thévenin equivalent voltage** and **Thévenin equivalent resistance**, respectively. This equivalent circuit is a direct consequence of the *linearity* of a passive network. We will defer the proof to Section 1.5.4 because it is fairly technical and not necessary for understanding it; for now we'll concentrate on making good use of it.

How do we find the Thévenin equivalent component values? A couple of observations:

1. If nothing is connected to the output, then  $V_{\text{out}} = V_{\text{Th}}$ .
2. If the output is short-circuited,<sup>2</sup> then  $V_{\text{out}} = 0$  and  $V_{\text{Th}} = R_{\text{Th}} I_{\text{short}}$ .

These allow you to infer the Thévenin values. The second rule is useful experimentally, provided shorting the output does not destroy the circuit!

**An alternate rule to find the Thévenin resistance**, especially in analyzing a circuit on paper, is as follows.

2. (alternate rule) Replace all EMFs by short circuits, and compute the equivalent resistance at the output.

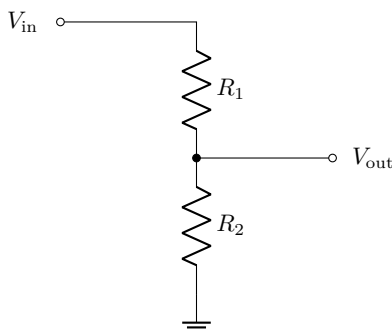
This is, in fact, typically the more useful version of rule 2.

<sup>1</sup>Named for L. Thévenin, “Extension de la loi d’Ohm aux circuits électromoteurs complexes,” *Annales Télégraphiques (Troisième série)* **10**, 222 (1883); L. Thévenin, “Sur un nouveau théorème d’électricité dynamique,” *Comptes Rendus Hebdomadaires des Séances de l’Académie des Sciences* **97**, 159 (1883). For a historical overview as well as an English translation of part of Thévenin’s papers, see Don H. Johnson, “Origins of the Equivalent Circuit Concept: The Voltage-Source Equivalent,” *Proceedings of the IEEE* **91**, 636 (2003) (doi: 10.1109/JPROC.2003.811716) <https://pdfs.semanticscholar.org/b55d/3f11daa8943d76f6e645b13020d4b1602648.pdf>. The proof as mentioned by Johnson is a **very** brief version of the one we present in Section 1.5.4.

<sup>2</sup>A “short circuit” or “short” is a direct, low-resistance path between two points in a circuit, like a piece of wire. It acts like a “shortcut” for current between the two points, compared to going through the regular, higher-resistance paths in the circuit.

### 1.4.1 Voltage Divider

Back to the voltage divider of Section 1.3.3.



1. We already found the Thévenin voltage as the unloaded-output voltage:

$$V_{\text{Th}} = \left( \frac{R_2}{R_1 + R_2} \right) V_{\text{in}}. \quad (1.19)$$

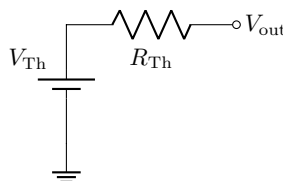
(voltage divider, Thévenin voltage)

2. If we short the output (to ground), the current only flows through  $R_1$ , so  $I_{\text{short}} = V_{\text{in}}/R_1$ . Then using  $R_{\text{Th}} = V_{\text{Th}}/I_{\text{short}}$ , we get

$$R_{\text{Th}} = \frac{R_1 R_2}{R_1 + R_2} = R_1 \parallel R_2. \quad (1.20)$$

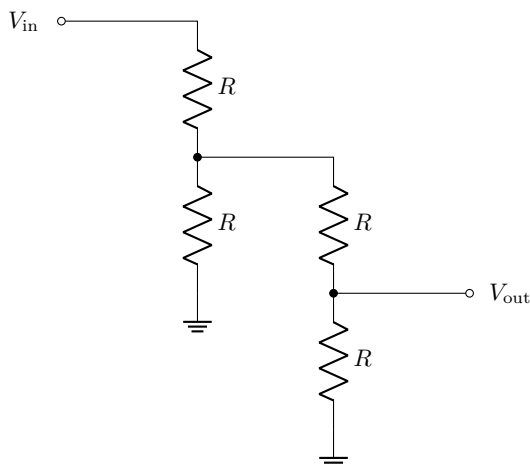
(voltage divider, Thévenin resistance)

That is, the equivalent circuit is as below

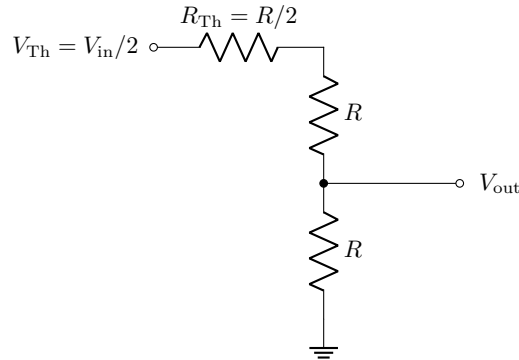


The equivalent voltage is the divided voltage, and the equivalent resistance is the **parallel resistance** of the two resistors. **This is important; you should memorize this.** This result also follows directly from the alternative rule: replacing  $V_{\text{in}}$  with a short to ground just puts  $R_1$  and  $R_2$  in parallel from the point of view of the output.

Now back to the example of two cascaded voltage dividers.



We can replace the first divider by the Thévenin equivalent:



Now we are back to a simple voltage divider, so

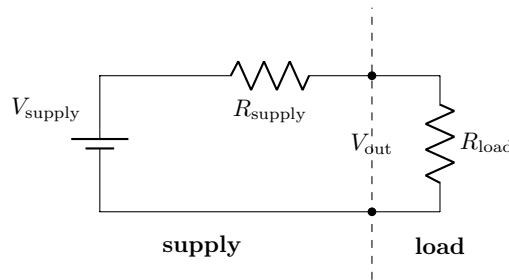
$$V_{\text{out}} = \left(\frac{2}{5}\right) \frac{V_{\text{in}}}{2} = \frac{V_{\text{in}}}{5}. \quad (1.21)$$

### 1.4.2 Connected Circuits and Power Transfer

The idea of Thévenin equivalence is also useful when analyzing what happens when you connect two fairly arbitrary circuits together. Suppose we take the common situation of one circuit “powering” another. That is, two circuits interact,

1. The “**supply**” is some circuit **with** EMFs.
2. The “**load**” is some circuit **without** EMFs.

Using the Thévenin-equivalent circuits for both, we can represent the connection thusly:



This is just a voltage divider, so we can write the output as

$$V_{\text{out}} = \left(\frac{R_{\text{load}}}{R_{\text{load}} + R_{\text{supply}}}\right) V_{\text{supply}}. \quad (1.22)$$

Now a few remarks are in order.

1. If  $R_{\text{load}} \gg R_{\text{supply}}$ , then  $V_{\text{out}} \approx V_{\text{supply}}$ , and the supply acts like an ideal voltage source. This applies to the “unloaded” (large  $R_{\text{load}}$ ) and “good supply” (small  $R_{\text{supply}}$ ) regimes.
2. If the supply is a battery,  $R_{\text{supply}}$  is the “internal resistance” of the battery. The internal resistance of a battery is larger for smaller or “used-up” batteries. The symptom is that the voltage sags when you try to draw current.
3.  $R_{\text{supply}}$  is called the **output impedance** of the supply.  $R_{\text{load}}$  is called the **input impedance** of the load.

4. The **impedance-matching condition** answer, under what conditions is maximum power transferred from source to load? The power in the load is

$$P_{\text{load}} = \frac{V_{\text{out}}^2}{R_{\text{load}}} = \frac{R_{\text{load}}}{(R_{\text{load}} + R_{\text{supply}})^2} V_{\text{supply}}^2. \quad (1.23)$$

Maximizing this via

$$\frac{d}{da} \left( \frac{a}{(a+b)^2} \right) = \frac{1}{(a+b)^2} - \frac{2a}{(a+b)^3} = 0, \quad (1.24)$$

which leads to  $a = b$ , we have the matching condition

$$R_{\text{load}} = R_{\text{supply}}. \quad (1.25)$$

(impedance-matching condition)

This is saying, for a *fixed* source impedance, the most power we can get out of the source and into the load is if the load impedance matches the supply impedance. In older tube amplifiers, this was an important consideration. For efficient matching to different speaker loads, amplifier output transformers would often have different “taps” for 4 Ω, 8 Ω, 16 Ω, etc. speakers.

## 1.5 Matrix Solution of Resistor Networks

Have a look at the XKCD comic “Circuit Diagram,”<sup>3</sup> and enjoy (you’ll recognize more stuff here as you learn more about electronics). Make sure to hover the cursor over the comic so you see the last joke.

Now look at part of the circuit labelled “Oh, so you think you’re such a whiz at EE 201?” (If you can’t access the circuit for whatever reason, it is a rat’s nest of resistors, and the idea is to find the equivalent resistance.) Randall Munroe was joking, but we’ll develop a systematic way to handle this kind of problem, which you can use to tackle that mess without much difficulty.

### 1.5.1 Review of Linear Algebra

First, we’re going to use some linear algebra (in practice, we will want the help of a computer), so let’s review the notation. A **matrix** is a group of numbers indexed by two numbers. For example, we can write down a  $2 \times 2$  matrix as

$$\begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}. \quad (1.26)$$

We can refer to the whole matrix as **A**. We can also refer an **element** (one of the entries) of the matrix as  $A_{ij}$ . Note that  $i$  refers to the element’s *row*, while  $j$  refers to the element’s *column*. We can write a **system of linear equations** as

$$\begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}. \quad (1.27)$$

This is just another way to write down the pair of equations

$$\begin{aligned} A_{11}x_1 + A_{12}x_2 &= b_1 \\ A_{21}x_1 + A_{22}x_2 &= b_2 \end{aligned} \quad (1.28)$$

and the shorthand notation for the matrix form is

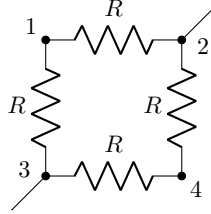
$$\mathbf{Ax} = \mathbf{b}, \quad (1.29)$$

where **x** and **b** are **vectors** (i.e.,  $n \times 1$  matrices, or specifically here,  $2 \times 1$  matrices). Under certain conditions, it is possible to solve for the  $x_j$  in terms of the  $b_i$  and the  $A_{ij}$  (we’ll let a computer help here). Make sure you understand the pattern of the matrix-vector multiplication in the equations above *before* you continue.

<sup>3</sup><http://xkcd.com/730/>

### 1.5.2 Matrix Form of the Resistance Network and Example

To develop the method to solve general networks of resistors, we'll apply it to the simple circuit below, to make the method more intuitive.



The goal is to find the effective resistance between points 2 and 3.

Instead of working with resistance, we'll work with the **conductance**  $G = 1/R$ , so that Ohm's law is

$$I = GV. \quad (1.30)$$

Why? Well, conductance is well-defined between disconnected points ( $G = 0$ ), and it will turn out to add up in the correct way for this formalism. Now between any two nodes in the network, we can write Ohm's law as

$$I_{ij} = G_{ij}(V_i - V_j), \quad (1.31)$$

where  $V_j$  is the voltage at node  $j$ ,  $G_{ij}$  is the conductance of the connection between the two points (with  $G_{jj} = 0$ , so there is no "self-conductance", and  $G_{ij} = G_{ji}$ ), and  $I_{ij}$  is the current flowing from node  $i$  to node  $j$ . Specifically, in this example problem, we have  $G_{12} = G_{24} = G_{43} = G_{31} = G$ , while  $G_{23} = G_{14} = 0$ .

Now Kirchhoff's current law states that the sum of currents flowing out of a node must be zero (note that we aren't yet considering any currents going *in* to a node, the way we defined  $I_{ij}$ ):

$$\sum_{j \neq i} I_{ij} = 0. \quad (1.32)$$

Note that we explicitly excluded the  $I_{ii} = 0$  case. This means

$$\sum_{j \neq i} G_{ij}(V_i - V_j) = 0. \quad (1.33)$$

Separating out terms we can rewrite this as

$$\left( \sum_{j \neq i} G_{ij} \right) V_i + \sum_{j \neq i} (-G_{ij}) V_j = 0. \quad (1.34)$$

We can interpret this as a linear system of equations as follows.

- The coefficient of  $V_i$  represents the diagonal elements of a matrix, call it  $\mathbf{A}$ . Then  $A_{ii}$  is the sum of all conductances connected to node  $i$ .
- The sum in the second term gives all the off-diagonal elements of the matrix. That is,  $A_{ij} = -G_{ij}$ , or the conductance between nodes  $i$  and  $j$ , with a minus sign.

To apply this to the example above, we will have a  $4 \times 4$  matrix  $\mathbf{A}$ , and the system (1.34) of equations becomes

$$R^{-1} \begin{bmatrix} 2 & -1 & -1 & 0 \\ -1 & 2 & 0 & -1 \\ -1 & 0 & 2 & -1 \\ 0 & -1 & -1 & 2 \end{bmatrix} \begin{bmatrix} V_1 \\ V_2 \\ V_3 \\ V_4 \end{bmatrix} = 0, \quad (1.35)$$

or  $\mathbf{A}\mathbf{v} = 0$  for short. Note that  $R^{-1}$  multiplies every element in the matrix  $\mathbf{A}$ . Note also that the matrix  $A_{ij}$  is **symmetric**:  $A_{ij} = A_{ji}$ . **This is a very good sanity check for this matrix**, especially when you set up the matrix for part (c). Go through the elements of this matrix to make sure they make sense. The diagonal elements are all 2 because every node is connected to two identical resistors. There is a  $-1$  representing every connection (via an identical resistor) between two nodes.

### 1.5.3 Solution for the Effective Resistance

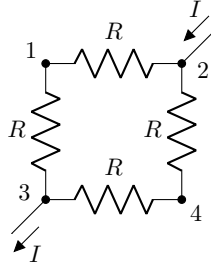
Unfortunately, the above matrix equation is not useful as is, because it is satisfied by  $V_i = V$  for any constant  $V$ . Physically, without any applied voltage, this circuit really doesn't *do* much. But first, we'll take care of *another* problem: the matrix  $\mathbf{A}$  is *singular*, meaning one of the equations in the linear system is redundant. Physically, this is because the absolute voltage of the circuit is not defined (the equations, by construction, only determine voltage *differences*). We can fix this by explicitly tying one of the nodes to zero voltage. Since node 3 is one of the nodes of interest, let's set

$$V_3 = 0, \quad (1.36)$$

which means we replace one of the equations in the linear system by this one. In matrix form, we will modify the third row of the matrix as follows:

$$R^{-1} \begin{bmatrix} 2 & -1 & -1 & 0 \\ -1 & 2 & 0 & -1 \\ 0 & 0 & 1 & 0 \\ 0 & -1 & -1 & 2 \end{bmatrix} \begin{bmatrix} V_1 \\ V_2 \\ V_3 \\ V_4 \end{bmatrix} = 0. \quad (1.37)$$

To solve the other problem (i.e., to make the circuit do something), let's introduce a current  $I$ , which flows into node 2. The same current  $I$  must flow *out* of node 3.



To handle this, we will modify Eq. (1.34) to say that the sum of currents flowing out of a node is equal to the currents flowing into a node (counting the external currents  $I$  as inputs, with minus signs to properly reflect their direction):

$$\left( \sum_{j \neq i} G_{ij} \right) V_i + \sum_{j \neq i} (-G_{ij}) V_j = I_{\text{in},i}. \quad (1.38)$$

Here  $I_{\text{in},i}$  is the current flowing in to node  $i$ . In our example problem,  $I_{\text{in},2} = I$ , while  $I_{\text{in},3} = -I$ . Now writing the linear system including these currents, we have

$$\begin{bmatrix} 2 & -1 & -1 & 0 \\ -1 & 2 & 0 & -1 \\ 0 & 0 & 1 & 0 \\ 0 & -1 & -1 & 2 \end{bmatrix} \begin{bmatrix} V_1 \\ V_2 \\ V_3 \\ V_4 \end{bmatrix} = \begin{bmatrix} 0 \\ IR \\ 0 \\ 0 \end{bmatrix}. \quad (1.39)$$

Note that we kept the third row of the matrix to reflect  $V_3 = 0$ , and we only introduced the current  $I$  at node 2. This is a well-defined system of equations, which we can solve to obtain  $V_2$ . The effective resistance of the network is defined by

$$R_{\text{eff}} = \frac{V_2 - V_3}{I} = \frac{V_2}{I}. \quad (1.40)$$

When you solve the matrix equation, as you might expect, the result should be  $R_{\text{eff}} = R$ . We will leave the details and the application to the XKCD problem to an exercise (Problem 1.20).

### 1.5.4 Proof of Thévenin's Theorem

A slightly more general version of the above matrix formalism for resistor networks allows us to prove Thévenin's theorem. Recall that the theorem dealt with networks of resistors and EMFs; first let's handle the resistors and then put the EMFs in later.

But before proceeding it is useful to recast the above matrix formalism as an optimization problem. To do this, we can begin by noting that the resistor connecting nodes  $i$  and  $j$  dissipates a power

$$P_{ij} = I_{ij}(V_i - V_j) = G_{ij}(V_i - V_j)^2, \quad (1.41)$$

where we used Eq. (1.31) in the last step. Summing over all resistors, we can then write the total power as

$$P = \sum_{i,j < i} P_{ij} = \sum_{i,j < i} G_{ij}(V_i - V_j)^2 = \sum_{i,j \leq i} G_{ij}(V_i - V_j)^2, \quad (1.42)$$

where we are taking advantage of the definition  $G_{jj} = 0$  for any  $j$ . So far, we are being careful to count any resistor only once, but we are free to double-count provided we divide the result by two. Thus,

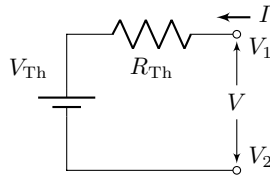
$$P(\mathbf{v}) = \frac{1}{2} \sum_{i,j} G_{ij}(V_i - V_j)^2, \quad (1.43)$$

where we are explicitly regarding the power to be a function of the voltage “coordinate vector”  $\mathbf{v} = (V_1, V_2, \dots, V_N)$ . The solution corresponds to a minimum in  $P$  in the sense that the vanishing-gradient condition  $\nabla P(\mathbf{v}) = 0$  gives

$$\frac{\partial P}{\partial V_i} = \sum_j G_{ij}(V_i - V_j) = 0. \quad (1.44)$$

This is just the Kirchoff current law in the form (1.33), which we used to set up the matrix solution of a general resistor network above. Intuitively, the circuit voltages (and currents) will arrange themselves in such a way as to *minimize* the power dissipated in the network, subject to any constraints on the network, such as applied voltages. (Technically, we have only showed that the power is *stationary* with respect to voltage changes, and thus could correspond to a saddle point or even a maximum, though minimization of power is nicely intuitive.) It is the constraints that we will now have to deal with explicitly to tackle Thévenin's theorem.

Now we need to put some EMFs in the circuit. First, recall that the Thévenin circuit is as below; one way to characterize its behavior is to impose a voltage  $V$  across the output, and see what current  $I$  flows.



Of course, these should be related by

$$V = IR_{Th} + V_{Th}. \quad (1.45)$$

To impose this behavior, we can add in an input current as in the previous section, but instead we will take a somewhat different (but ultimately equivalent) approach. To impose the voltage  $V$ , suppose we impose the constraint that  $V_1 - V_2 = V$  (i.e.,  $V_1$  and  $V_2$  are the “output terminals” of the circuit). Since  $V_1 - V_2 - V = 0$  anyway, it doesn't change anything if we multiply it by an some constant and then add the (still-zero) result to the power (1.43):

$$P(\mathbf{v}) = \frac{1}{2} \sum_{i,j} G_{ij}(V_i - V_j)^2 - I(V_1 - V_2 - V). \quad (1.46)$$

To make the dimensions work out, the constant multiplying the “zero” has to be a current, and we called the current  $-I$  in anticipation of the upcoming solution (it will indeed turn out to be the current flowing



via the output terminals). However, technically at this point it is an undetermined constant, and the constraint equation  $V_1 - V_2 - V = 0$  along with the  $\nabla P(\mathbf{v}) = 0$  equations will allow us to determine  $I$ . In optimization-theory parlance,  $I$  is called a **Lagrange multiplier**.

To add in the other EMFs in the resistor network, the idea proceeds by adding constraints in the same way. Suppose we add in a bunch of different EMFs  $\mathcal{E}_k$  connecting various pairs of network nodes. These constraints can be represented by vanishing constraint functions

$$f_k(\mathbf{v}) = V_{\sigma(k)} - V_{\tau(k)} - \mathcal{E}_k = 0, \quad (1.47)$$

which imposes the EMF  $\mathcal{E}_k$  between the  $\sigma(k)$ th and  $\tau(k)$ th nodes, where  $\sigma(k)$  and  $\tau(k)$  are indexing functions that simply track the node placement of the  $k$ th EMF. At most, with  $N$  total voltage nodes, there can be  $N-2$  more constraint EMFs (in addition to the one we already imposed) without overconstraining the circuit. We otherwise won't worry about inconsistent constraints, and we can note that in the case where the output nodes should be voltage-constrained by an EMF that is part of the circuit, the resulting Thévenin circuit is trivial, so we won't worry about that case here. Introducing extra Lagrange multipliers  $\mathcal{I}_k$ , Eq. (1.46) becomes

$$P(\mathbf{v}) = \frac{1}{2} \sum_{i,j=1}^N G_{ij}(V_i - V_j)^2 - I(V_1 - V_2 - V) - \sum_{k=1}^{N_C \leq N-2} \mathcal{I}_k(V_{\sigma(k)} - V_{\tau(k)} - \mathcal{E}_k). \quad (1.48)$$

after implementing all the EMFs as constraints. Here  $N_C \leq N-2$  is the number of constraints (EMFs).

Now differentiating Eq. (1.48) with respect to  $V_i$  leads to the system of equations

$$\sum_{j=1}^N G_{ij}(V_i - V_j) - I(\delta_{i,1} - \delta_{i,2}) - \sum_{k=1}^{N_C \leq N-2} \mathcal{I}_k(\delta_{i,\sigma(k)} - \delta_{i,\tau(k)}) = 0, \quad (1.49)$$

where  $\delta_{nm}$  is the Kronecker delta, satisfying  $\delta_{nn} = 1$  and  $\delta_{nm} = 0$  whenever  $n \neq m$ . If we rearrange the first sum as in Eq. (1.34), this becomes

$$\left( \sum_{j=1}^N G_{ij} \right) V_i + \sum_{j=1}^N (-G_{ij}) V_j - \sum_{k=1}^{N_C \leq N-2} \mathcal{I}_k(\delta_{i,\sigma(k)} - \delta_{i,\tau(k)}) = I(\delta_{i,1} - \delta_{i,2}), \quad (1.50)$$

and comparison with Eq. (1.38) shows that the Lagrange multipliers indeed have the interpretation as extra currents into the respective nodes. In particular the  $\mathcal{I}_k$  are the currents supplied by the corresponding EMFs  $\mathcal{E}_k$  in order to maintain their imposed voltages within the network.

The  $N$  equations (1.50), along with the constraint equations (1.47) in the form

$$V_{\sigma(k)} - V_{\tau(k)} = \mathcal{E}_k, \quad (1.51)$$

with  $V_1 - V_2 = V$ , give  $N + N_C + 1$  equations for the  $N + N_C + 2$  variables  $(V, V_1, \dots, V_N, I, \mathcal{I}_1, \dots, \mathcal{I}_{N_C})$ . However, recall that the system represented by Eq. (1.50) is singular in the sense that one is redundant, corresponding physically to the ambiguity of the absolute voltage. We can cure this by replacing one of the equations, say  $i = 2$ , by the condition  $V_2 = 0$ . Assuming this has been done, we can simply implement  $V_1 - V_2 = V$  by replacing  $V_1$  by  $V$ , and then we can write the above  $N + N_C$  remaining equations as

$$\mathcal{G}\mathbf{v} = I\hat{e}_1 + \boldsymbol{\epsilon}, \quad (1.52)$$

where  $\mathbf{v}$  has been extended to include the Lagrange-multiplier currents as  $\mathbf{v} := (V, V_2, \dots, V_N, \mathcal{I}_1, \dots, \mathcal{I}_{N_C})$ ;  $\mathcal{G}$  is a square matrix encoding the left-hand sides of the Kirchhoff and EMF-constraint equations;  $\hat{e}_1$  is the unit vector  $(1, 0, 0, \dots, 0)$ ; and  $\boldsymbol{\epsilon} := (0, 0, \dots, 0, \mathcal{E}_1, \dots, \mathcal{E}_{N_C})$ . This system has the solution

$$\mathbf{v} = (\mathcal{G}^{-1}\hat{e}_1)I + \mathcal{G}^{-1}\boldsymbol{\epsilon}, \quad (1.53)$$

Recall that the first equation is the equation determining  $V$ . Thus, the most relevant solution for our purposes is

$$V = \hat{e}_1 \cdot (\mathcal{G}^{-1}\hat{e}_1)I + \hat{e}_1 \cdot (\mathcal{G}^{-1}\boldsymbol{\epsilon}), \quad (1.54)$$

(Thévenin's theorem)

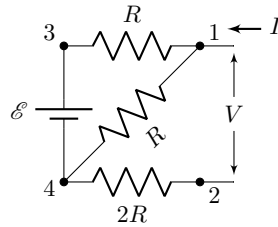
although the other equations are of course useful in giving the voltages and currents throughout the rest of the circuit. Notice that this has precisely the form of Eq. (1.45), provided we identify

$$R_{\text{Th}} = \hat{e}_1 \cdot (\mathcal{G}^{-1} \hat{e}_1), \quad V_{\text{Th}} = \hat{e}_1 \cdot (\mathcal{G}^{-1} \epsilon). \quad (1.55)$$

Note that this is sufficient to establish Thévenin's theorem. We can find the Thévenin voltage by setting  $V = V_{\text{Th}}$ , so that  $I = 0$  (i.e., find the Thévenin voltage by finding the voltage of the unloaded circuit); find the Thévenin resistance by shorting the EMFs ( $\epsilon = 0$ ), and then find the equivalent resistance (i.e., find the  $I$  response to an input  $V$ , and take the ratio to find  $R_{\text{Th}}$ ).

#### 1.5.4.1 Example Network

This proof gets a bit abstract at the end, so perhaps a more specific example will make the above algebra more clear. A simple example circuit with only one EMF is below, and again the idea is to use the above technique to find the Thévenin equivalent circuit, assuming connections at nodes 1 and 2 (where the input voltage  $V$  and current  $I$  will determine the Thévenin-equivalent action of the circuit).



For this circuit, it's not hard to find the Thévenin values  $V_{\text{Th}} = \mathcal{E}/2$  and  $R_{\text{Th}} = 5R/2$ . A direct translation of Eq. (1.50) into a matrix equation could read as below, recalling that the procedure for setting up the voltage matrix is the same as in Eq. (1.34):

$$R^{-1} \begin{bmatrix} 2 & 0 & -1 & -1 \\ 0 & 1/2 & 0 & -1/2 \\ -1 & 0 & 1 & 0 \\ -1 & -1/2 & 0 & 3/2 \end{bmatrix} \begin{bmatrix} V_1 \\ V_2 \\ V_3 \\ V_4 \end{bmatrix} = \begin{bmatrix} I \\ -I \\ \mathcal{I} \\ -\mathcal{I} \end{bmatrix}. \quad (1.56)$$

However, the idea here is a little different: while we can still treat (the external)  $I$  as an “input,” we should treat (the internal)  $\mathcal{I}$  as an unknown to solve for just like the  $V_i$ . To keep the large matrix square, we can also throw in the constraint equation  $V_3 - V_4 = \mathcal{E}$ . The resulting system is

$$R^{-1} \begin{bmatrix} 2 & 0 & -1 & -1 & 0 \\ 0 & 1/2 & 0 & -1/2 & 0 \\ -1 & 0 & 1 & 0 & 1 \\ -1 & -1/2 & 0 & 3/2 & -1 \\ 0 & 0 & 1 & -1 & 0 \end{bmatrix} \begin{bmatrix} V_1 \\ V_2 \\ V_3 \\ V_4 \\ \mathcal{I}R \end{bmatrix} = \begin{bmatrix} I \\ -I \\ 0 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \mathcal{E}/R \end{bmatrix}. \quad (1.57)$$

Before solving this, a couple of last modifications. First, we can cure the singular nature of the matrix by replacing the second row by  $V_2 = 0$ . Second, we can identify  $V_1 = V$ . With these modifications, the system becomes

$$\begin{bmatrix} 2 & 0 & -1 & -1 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ -1 & 0 & 1 & 0 & 1 \\ -1 & -1/2 & 0 & 3/2 & -1 \\ 0 & 0 & 1 & -1 & 0 \end{bmatrix} \begin{bmatrix} V \\ V_2 \\ V_3 \\ V_4 \\ \mathcal{I}R \end{bmatrix} = \begin{bmatrix} IR \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \mathcal{E} \end{bmatrix}, \quad (1.58)$$

and it's ready to solve. Inverting the matrix  $\mathcal{G}$  gives

$$\mathcal{G}^{-1} = \begin{bmatrix} 5/2 & 1 & 2 & 2 & 1/2 \\ 0 & 1 & 0 & 0 & 0 \\ 2 & 1 & 2 & 2 & 1 \\ 2 & 1 & 2 & 2 & 0 \\ 1/2 & 0 & 1 & 0 & -1/2 \end{bmatrix}, \quad (1.59)$$

and thus the solution

$$\begin{bmatrix} V \\ V_2 \\ V_3 \\ V_4 \\ \mathcal{I}R \end{bmatrix} = \mathcal{G}^{-1} \begin{bmatrix} IR \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} + \mathcal{G}^{-1} \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \mathcal{E} \end{bmatrix} = I \begin{bmatrix} (5/2)R \\ 0 \\ 2R \\ 2R \\ (1/2)R \end{bmatrix} + \begin{bmatrix} 1/2 \\ 0 \\ 1 \\ 0 \\ -1/2 \end{bmatrix} \mathcal{E}. \quad (1.60)$$

Reading off the top equation gives the Thévenin parameters,

$$V = I \left( \frac{5R}{2} \right) + \frac{\mathcal{E}}{2}, \quad (1.61)$$

so that  $R_{\text{Th}} = 5R/2$  and  $V_{\text{Th}} = \mathcal{E}/2$ , as expected. The other rows in the linear system give more information; for example the last equation gives  $\mathcal{I} = IR/2 - \mathcal{E}/2$  for the current through the EMF. The other rows give  $V_4 = 2IR$  and  $V_3 = \mathcal{E} + 2IR$ , which are both readily apparent from looking at the circuit diagram.

#### 1.5.4.2 Extension to Current Sources

In the proof above, we only considered a network of EMFs (ideal voltage sources) and resistors. However, Thévenin's theorem still works if you throw in a few current sources. We have already seen how to do this in Eq. (1.38), where an external “input” current flows into one node and out another. In terms of the power equation (1.48), we accomplished this same effect by inserting a constraint term  $-I(V_1 - V_2 - V)$ , which constrained  $V_{12} = V$  by imposing whatever input current  $I$  achieves this. The same idea can work to put in other current sources. Adding a term of the form  $-I_s(V_i - V_j - V_{ij})$  does this if we view it correctly: if we regard  $V_{ij}$  as an unknown, it will simply take on the value  $V_i - V_j$  while introducing the source current  $I_s$  between nodes  $i$  and  $j$ . The upshot is that Eq. (1.50) becomes

$$\left( \sum_{j=1}^N G_{ij} \right) V_i + \sum_{j=1}^N (-G_{ij}) V_j - \sum_{k=1}^{N_C \leq N-2} \mathcal{I}_k (\delta_{i,\sigma(k)} - \delta_{i,\tau(k)}) = I(\delta_{i,1} - \delta_{i,2}) + \sum_{k=1}^{N_I} I_{sk} (\delta_{i,\sigma'(k)} - \delta_{i,\tau'(k)}), \quad (1.62)$$

where there are  $N_I$  current sources, with the  $k$ th source between nodes  $\sigma'(k)$  and  $\tau'(k)$ . Note that these current sources are not real constraints; they simply redirect current through the circuit, and thus there is in principle no limit to how many current sources we can have.

These extra current sources can be collected in the same way as  $I$ , so that Eq. (1.52) becomes

$$\mathcal{G}\mathbf{v} = I\hat{\mathbf{e}}_1 + \boldsymbol{\epsilon} + \boldsymbol{\iota}, \quad (1.63)$$

where the  $i$ th entry in  $\boldsymbol{\iota}$  is the sum of all source currents into node  $i$ , minus the sum of all source currents out of the same node. The matrix system has the more general solution

$$\mathbf{v} = (\mathcal{G}^{-1}\hat{\mathbf{e}}_1)I + \mathcal{G}^{-1}\boldsymbol{\epsilon} + \mathcal{G}^{-1}\boldsymbol{\iota}, \quad (1.64)$$

and the first equation of this linear system is

$$V = \hat{\mathbf{e}}_1 \cdot (\mathcal{G}^{-1}\hat{\mathbf{e}}_1)I + \hat{\mathbf{e}}_1 \cdot (\mathcal{G}^{-1}\boldsymbol{\epsilon}) + \hat{\mathbf{e}}_1 \cdot (\mathcal{G}^{-1}\boldsymbol{\iota}), \quad (\text{Thévenin's theorem with current sources}) \quad (1.65)$$

Thus, we have the same theorem, but with a modified Thévenin voltage:

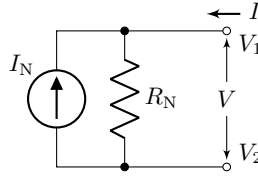
$$R_{\text{Th}} = \hat{e}_1 \cdot (\mathcal{G}^{-1} \hat{e}_1), \quad V_{\text{Th}} = \hat{e}_1 \cdot [\mathcal{G}^{-1}(\epsilon + \iota)]. \quad (1.66)$$

Note that it looks like we are conflating units, especially in the Thévenin voltage, but note that by construction the  $\epsilon$  and  $\iota$  vectors cannot have nonzero entries in the same position.

The procedure for handling current sources in applying Thévenin's theorem is the same: To find the Thévenin voltage, find the unloaded output voltage; to find the Thévenin resistance, find the resistance with all EMFs replaced by short circuits and all current sources replaced by open circuits (corresponding to a zero currents).

### 1.5.4.3 Norton's Theorem

Incidentally, there is another possible equivalent circuit that serves as an alternative to the Thévenin circuit. Consider the following circuit, analogous to the Thévenin circuit (p. 33).



Here, both the source current  $I_N$  and the external current  $I$  flow through  $R_N$ , and thus

$$V = IR_N + I_N R_N. \quad (1.67)$$

Note that Eq. (1.65) has this form and thus the proof for Thévenin's theorem also shows that this current-source circuit is also an equivalent representation, with the identifications

$$R_N = \hat{e}_1 \cdot (\mathcal{G}^{-1} \hat{e}_1), \quad I_N = \frac{\hat{e}_1 \cdot [\mathcal{G}^{-1}(\epsilon + \iota)]}{\hat{e}_1 \cdot (\mathcal{G}^{-1} \hat{e}_1)}. \quad (1.68)$$

This is called **Norton's theorem**,<sup>4</sup> and the procedure for finding the elements of the equivalent circuit is analogous to the Thévenin procedure:

1. With the output short-circuited, calculate the output current  $I_{\text{short}}$ , which is just  $I_N$ .
2. With the output unloaded, compute the output voltage, which gives  $R_N$  via  $V_{\text{out}} = V = I_N R_N$ .
2. (alternate rule) Replace voltage sources by short circuits and current sources by open circuits, then calculate the equivalent resistance across the output terminals to get  $R_N$ .

In fact, comparing Eq. (1.67) to Eq. (1.45), it is easy to convert between a Thévenin and a Norton equivalent circuit by identifying

$$R_{\text{Th}} = R_N, \quad V_{\text{Th}} = I_N R_N. \quad (1.69)$$

(Thévenin–Norton conversion)

In practice, the Thévenin circuit seems to be the more useful representation, but it's still handy to know about the Norton equivalent.

<sup>4</sup>Named after E. L. Norton, "Design of finite networks for uniform frequency characteristic, Technical Report TM2670?1860, Bell Laboratories, 1926. This theorem is also called the Mayer–Norton theorem, after H. F. Mayer, "Über das Ersatzschema der Verstärkerröhre," *Telegraphen- und Fernsprech-Technik* **15**, 335 (1926). For a nice history, see Don H. Johnson, "Origins of the Equivalent Circuit Concept: The Current-Source Equivalent," *Proceedings of the IEEE* **91**, 817 (2003) (doi: 10.1109/JPROC.2003.811795) <https://pdfs.semanticscholar.org/8b43/8720d7d7773d2a2b599932648878c84d40e2.pdf>.

#### 1.5.4.4 Superposition of Sources

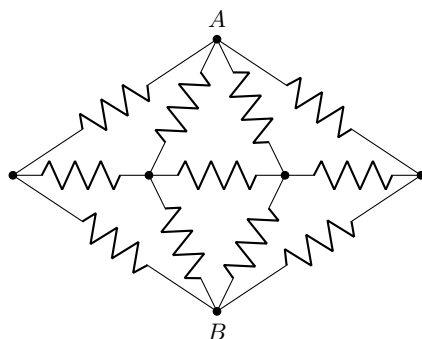
One final observation regarding Eq. (1.65) is in order here. The source vectors  $\epsilon$  and  $\iota$  for EMFs and current sources, respectively, can be decomposed into sums of vectors, each corresponding to a *single* source. The total effect on  $V$  is simply a sum over terms involving each source (including  $I$ ). The interpretation of this is, that if we know the value of  $V$  when we “turn on” only one source at a time, we can simply sum up all these voltages to find the value of  $V$  when all the sources are present together. The same conclusion holds for the current between two nodes, assuming they are not connected by a voltage source. This is sometimes called the **superposition theorem**, and is sometimes useful in analyzing linear circuits. (It is sometimes also used as the justification for Thévenin’s theorem.)

## 1.6 Circuit Practice

Here are a few example circuits to analyze with solutions; try to work these out and test your understanding so far. (Try *before* looking at the solutions!)

### 1.6.1 Reflection-Symmetric Network

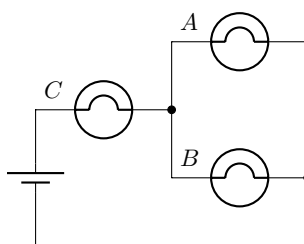
Given that all resistors here are equal and of resistance  $R$ , what is the equivalent resistance between  $A$  and  $B$ ?



**Solution.** Notice that by symmetry, the voltage across the center resistance is zero, so the current flowing through it is zero, no matter what the voltage  $V_{AB}$ . Thus, we can remove it from the circuit. The same goes for the two other “equatorial” resistors. Thus, we have 4 parallel resistances of  $2R$ , for a total resistance  $R/2$ .

### 1.6.2 Series and Parallel Light Bulbs

Suppose the three bulbs in the circuit below are identical.



What is the relative brightness of the bulbs? Do bulbs  $B$  and  $C$  get brighter or dimmer when you remove  $A$ ?

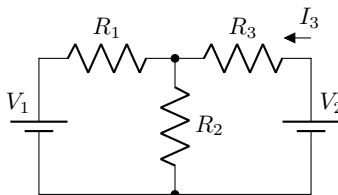
**Solution.** First, note that *half* the  $C$  current flows through  $A$  and  $B$ , while  $A$  and  $B$  drop half the voltage of  $C$ . So  $A$  and  $B$  are  $1/4$  as bright as  $C$ .

When  $A$  is removed, then  $B$  and  $C$  are equally bright. Less current flows overall, because the total resistance is larger.  $C$  drops less voltage than it used to,  $B$  drops more than it used to. Clearly  $C$  should be dimmer, but  $B$  should be brighter, because it now has more voltage and more current (the drop in overall current is outweighed by the factor of 2 this bulb gets by its neighbor disappearing).

To be careful, note in the original case that the voltage across  $C$  is  $(2/3)V$ , and across  $A$  and  $B$  it is  $(1/3)V$ , where  $V$  is the battery voltage. The current through  $C$  is  $(2/3)V/R$ , and half that flows through  $A$  and  $B$ . In the new case, each bulb drops  $V/2$ , and current  $(1/2)V/R$  flows through each. Before, the  $C$  power was the product of voltage and current, or  $(4/9)V^2/R$ , and now it is  $(1/2)V^2/R$ , so this is dimmer. Before, the  $B$  power was  $(1/9)V^2/R$ , so this one gets brighter.

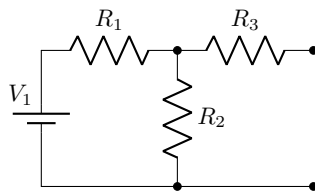
### 1.6.3 Thévenin Circuit

Compute  $I_3$  in the circuit below, by using the Thévenin equivalent to the circuit, thinking of  $V_2$  as the “load” for the rest of the circuit.



Use values  $R_1 = 1\ \Omega$ ,  $R_2 = 2\ \Omega$ ,  $R_3 = 4\ \Omega$ ,  $V_1 = 1\ \text{V}$ , and  $V_2 = 2\ \text{V}$ .

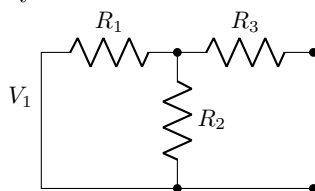
**Solution.** To work out the Thévenin equivalent, first look at the “unloaded” circuit:



The output voltage (and thus Thévenin-equivalent voltage) here is just the voltage-divider result

$$V_{\text{Th}} = \left( \frac{R_2}{R_1 + R_2} \right) V_1 = \frac{2}{3} \text{ V}. \quad (1.70)$$

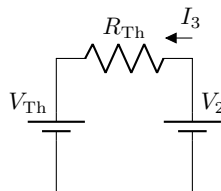
To get the resistance, replace the EMF by a short:



The Thévenin resistance is the resistance that appears at the output terminals, or

$$R_{\text{Th}} = R_3 + R_1 \parallel R_2 = R_3 + \frac{R_1 R_2}{R_1 + R_2} = \frac{14}{3} \text{ V}. \quad (1.71)$$

Now going back to the original circuit, we have the equivalent



This is pretty easy to solve, just use Ohm’s law:

$$I_3 = \frac{V_2 - V_{\text{Th}}}{R_{\text{Th}}} = \frac{(4/3) \text{ V}}{(14/3) \Omega} = \frac{2}{7} \text{ A}. \quad (1.72)$$

## 1.7 Exercises

### Problem 1.1

Ohm's law for the **conductance**  $G$  is

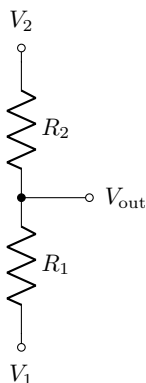
$$I = GV, \quad (1.73)$$

where the conductance is related to the resistance by  $G = 1/R$ , and  $G$  is measured in *mhos*, or  $\mathcal{U}$ . (An equivalent but less entertaining unit is the *siemens*, or  $S$ .)

Derive expressions for the conductance  $G$  of two conductors of conductances  $G_1$  and  $G_2$  in series. Do the same for two conductors in parallel. Use only the form of Ohm's law above to start; do *not* use the analogous results for parallel and series resistances.

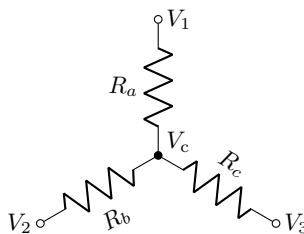
### Problem 1.2

Compute  $V_{\text{out}}$ . This is a voltage divider, but the “bottom” end is not grounded. *Put the result into some form you can remember and then memorize it.*



### Problem 1.3

In this circuit,



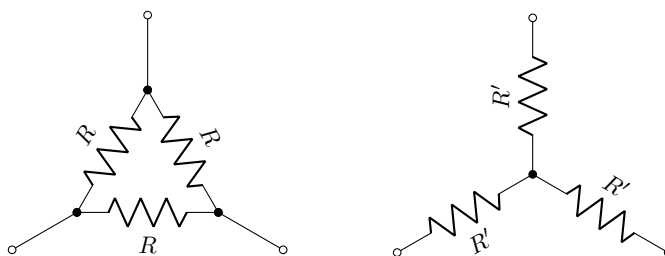
show that the “center voltage”  $V_c$  is given by

$$V_c = \frac{R_b R_c V_1 + R_c R_a V_2 + R_a R_b V_3}{R_a R_b + R_b R_c + R_c R_a}. \quad (1.74)$$

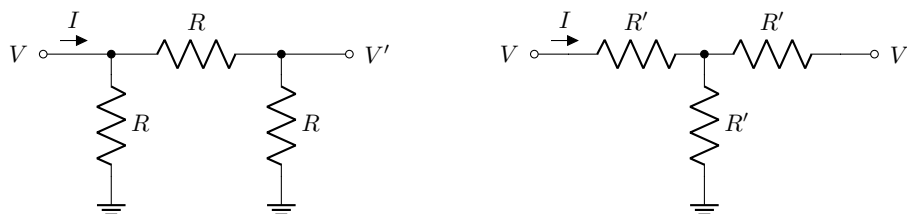
### Problem 1.4

Find the relation between  $R$  and  $R'$  that makes the two circuits below equivalent, given arbitrary connections at the three inputs.



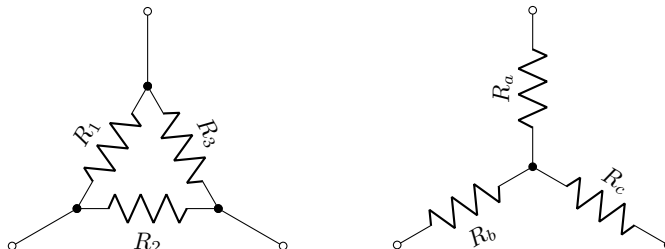


One way to do this is to redraw these circuits as the equivalent “pi” and “tee” networks below, and consider the current  $I$  into one port. However, you should *justify* that your treatment is sufficient to establish that the networks are *fully* equivalent.

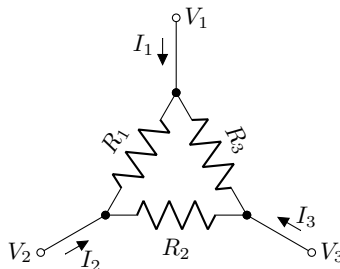


### Problem 1.5

Find relations between resistances  $R_1$ ,  $R_2$ , and  $R_3$  in the left-hand network, and resistances  $R_a$ ,  $R_b$ , and  $R_c$  in the right-hand network, that make the two networks equivalent.



What does it mean for two resistor networks to be equivalent? Label the voltages and currents, for example, as follows.



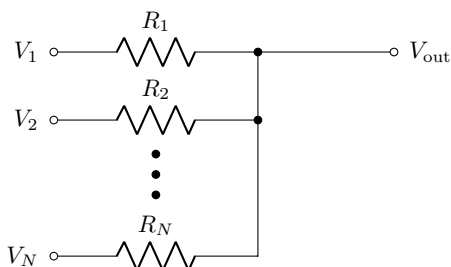
Now it's useful to think of some of the voltages and currents as “inputs,” and others as “outputs,” or “responses.” For example, if the three voltages are inputs, then they cause the corresponding currents to flow. You can also think of one or two of the currents as being inputs instead of their respective voltages. (But *not* all three currents; why?) The circuits are equivalent if every set of inputs produces the equivalent set of outputs. For the purposes of this problem, it is sufficient to show that the two networks produce equivalent relations between currents and voltages.

To do this:

- Derive expressions for the currents  $I_1$ ,  $I_2$ , and  $I_3$  in terms of the voltages for the left-hand (“Delta”) circuit.
- Derive expressions for the currents  $I_1$ ,  $I_2$ , and  $I_3$  in terms of the voltages for the right-hand (“star”) circuit.
- Set the corresponding currents in the two circuits equal to each other, and derive expressions for  $R_a$ ,  $R_b$ , and  $R_c$  in terms of  $R_1$ ,  $R_2$ , and  $R_3$ . (You may think of the voltages as “inputs,” and thus as independent parameters.)

### Problem 1.6

Consider the network below with  $N$  input voltages at each of  $N$  resistors.



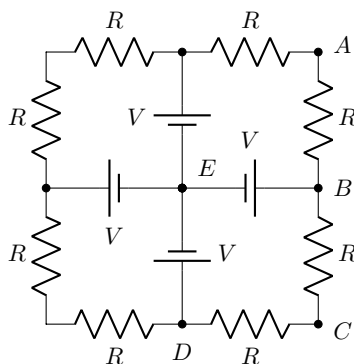
Prove **Millman’s theorem**, which gives

$$V_{\text{out}} = \frac{\sum_{j=1}^N G_j V_j}{\sum_{j=1}^N G_j} \quad (1.75)$$

for the (unloaded) output voltage, where  $G_j := 1/R_j$  is the conductance of the  $j$ th resistor.

### Problem 1.7

Consider the circuit below.

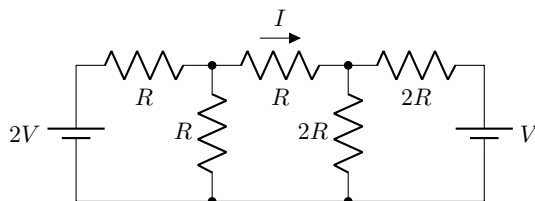


Give the Thévenin-equivalent circuit for the circuit in the diagram, as seen between

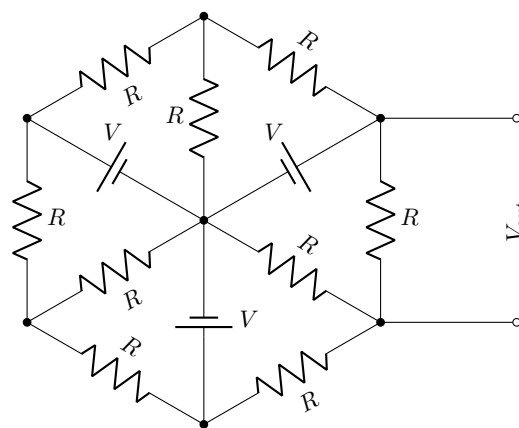
- points  $B$ – $D$ .
- points  $A$ – $B$ .
- points  $B$ – $E$ .
- points  $A$ – $E$ .

**Problem 1.8**

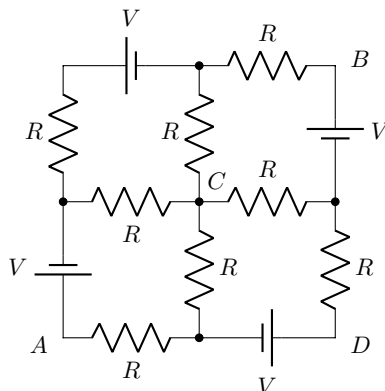
In the circuit below, compute the current  $I$  through the top-middle resistor.

**Problem 1.9**

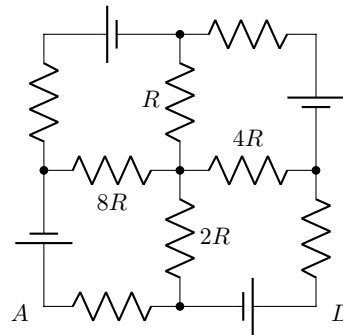
Give the Thévenin voltage  $V_{Th}$  and resistance  $R_{Th}$  for the Thévenin equivalent circuit of the circuit shown below.

**Problem 1.10**

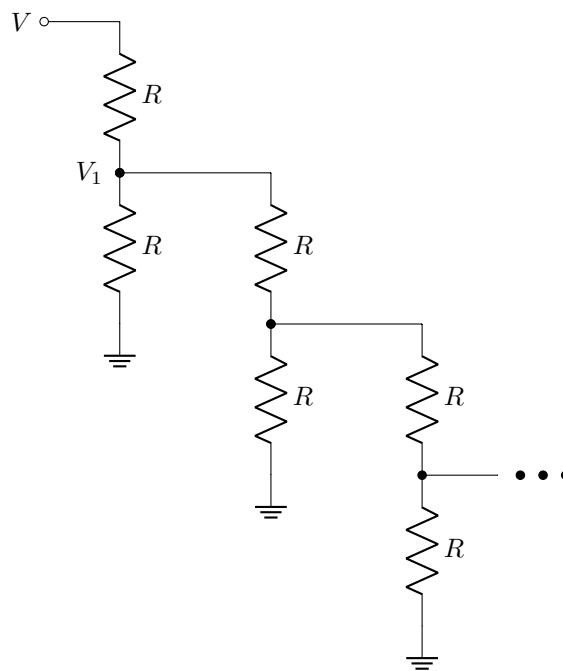
Consider the circuit below.



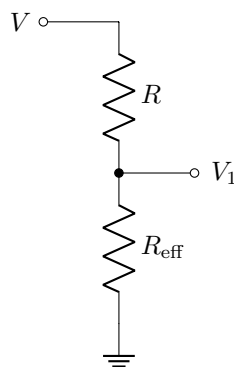
- Find the Thévenin equivalent of this circuit, as seen between points  $A$  and  $B$ .
- Do the same between points  $B$  and  $C$ . (The result of Problem 1.4 may help.)
- Suppose that we change the inner resistors as shown below, with the rest of the components the same. Calculate the voltage between points  $D$  and  $A$ .

**Problem 1.11**

Consider the following circuit, consisting of an infinite cascade of voltage dividers.



By how much does the rest of the divider chain load down the first divider? That is, in the equivalent circuit

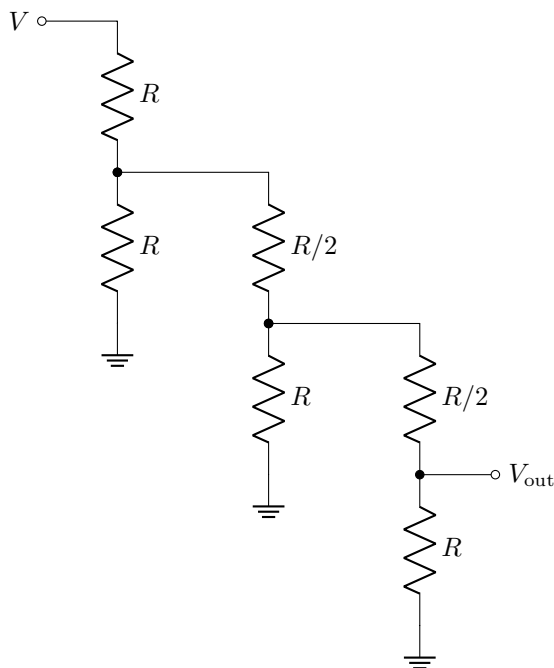


what is the value of  $R_{\text{eff}}$ ?

*Hint:* In the original divider chain, consider the *second* divider, along with the rest of the chain. *This* part of the circuit is *also* equivalent to the same effective circuit. Now you have two different, equivalent circuits in terms of  $R$  and  $R_{\text{eff}}$ ; compute the resistance from the point  $V$  to ground in both circuits, equate, and solve for  $R_{\text{eff}}$ .

### Problem 1.12

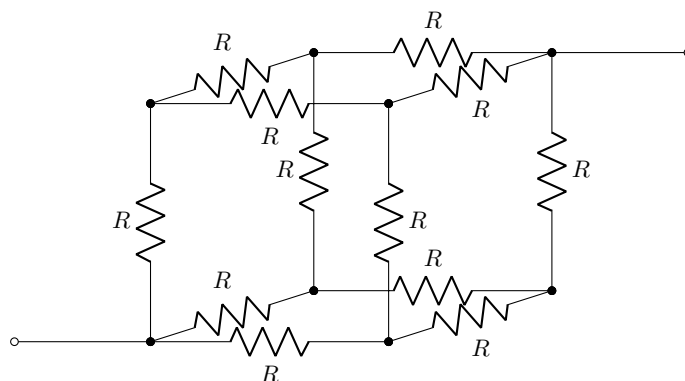
Consider the circuit below, with 3 cascaded voltage dividers (not all the same).



- Compute  $V_{\text{out}}$ .
- Compute the current in *each one* of the resistors in the circuit, assuming no load connected to  $V_{\text{out}}$ .

### Problem 1.13

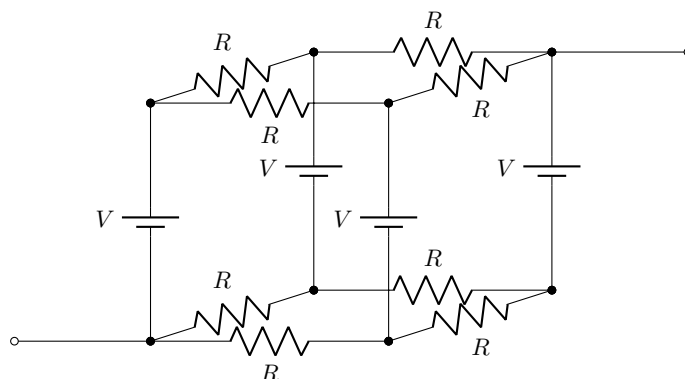
Consider a network of 12 identical resistors of resistance  $R$  in the shape of a cube, as shown in the circuit below. Compute the equivalent resistance between the two terminals shown at opposite corners of the cube.



*Hint:* use symmetry, and analyze a current flowing between the two terminals. It also may (or may not) help to identify which junctions are at equivalent potentials.

#### Problem 1.14

Consider the network below of 8 identical resistors and 4 identical voltage sources. With the two indicated terminals as the output, give the Thévenin equivalent circuit. *Explain* your result.



#### Problem 1.15

A “schmesistor” is a device that obeys “Schmohm’s law,”

$$V = I^2 S, \quad (1.76)$$

where  $S$  is the “schmesistance.”

(a) Show that two schmesistors  $S_1$  and  $S_2$  in **series** behave like one schmesistor with schmesistance  $S_{\text{eff}} = S_1 + S_2$ .

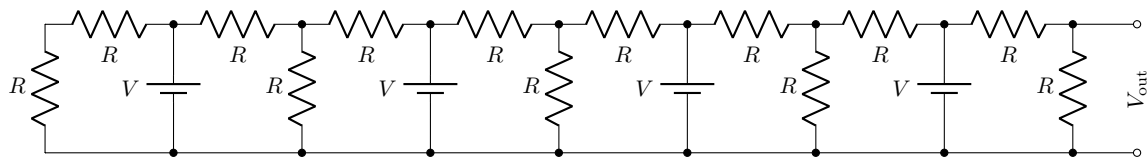
(b) Show that two schmesistors  $S_1$  and  $S_2$  in **parallel** behave like one schmesistor with schmesistance

$$S_{\text{eff}} = \left( \frac{1}{\sqrt{S_1}} + \frac{1}{\sqrt{S_2}} \right)^{-2}. \quad (1.77)$$

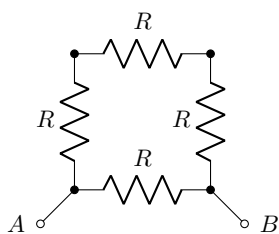
(c) Describe the difficulty in using complex notation to analyze a circuit that includes schmesistors.

**Problem 1.16**

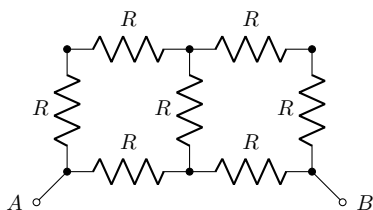
Give the Thévenin equivalent circuit for the circuit below, as “seen” by the output terminals marked by  $V_{\text{out}}$ .

**Problem 1.17**

(a) Compute the effective resistance between points  $A$  and  $B$ .

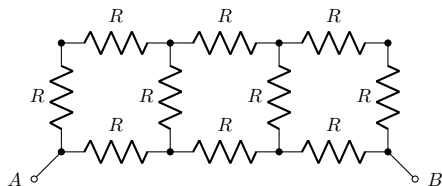


(b) Compute the effective resistance between points  $A$  and  $B$ .

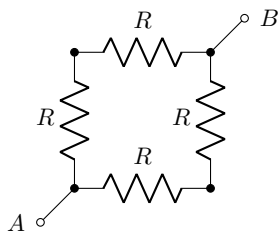


*Hint: symmetry.*

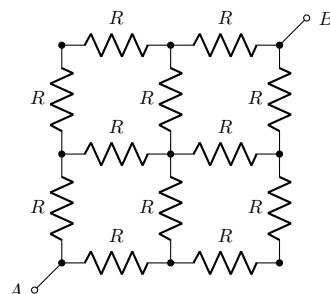
(c) Compute the effective resistance between points  $A$  and  $B$ .

**Problem 1.18**

(a) Compute the effective resistance between points  $A$  and  $B$ .

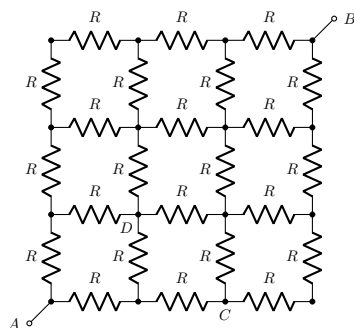


- (b) Compute the effective resistance between points  $A$  and  $B$ .



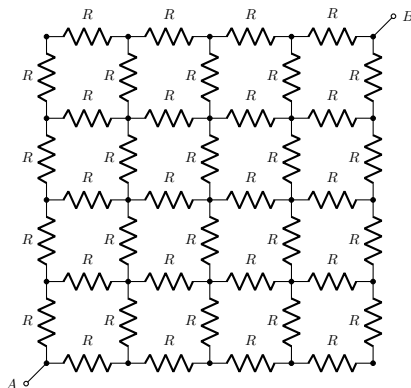
*Hint: symmetry.* There are (at least) two symmetries you should be taking advantage of here.

- (c) Compute the effective resistance between points  $A$  and  $B$ .



*Hint: symmetry,* but note that it is *not* true in general that  $V_C = V_D$ .

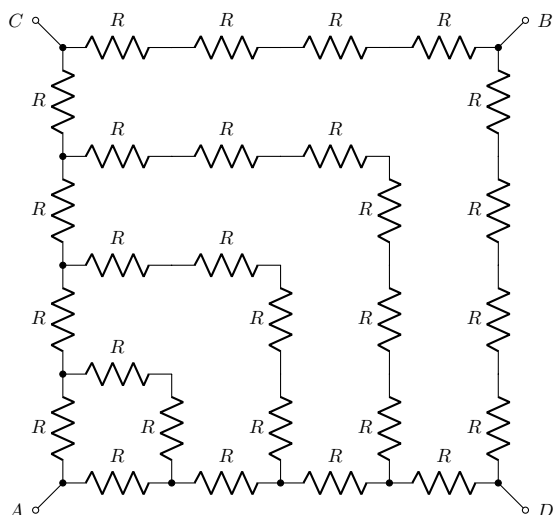
- (d) Generalize (a)–(c) to a resistor network forming a  $4 \times 4$  grid(!) of squares.



### Problem 1.19

- (a) Compute the equivalent resistance between points  $A$  and  $B$  in the circuit below.  
 (b) Compute the equivalent resistance between points  $C$  and  $D$  in the circuit below.





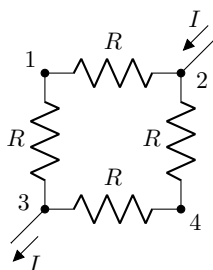
### Problem 1.20

Now back to the XKCD comic “Circuit Diagram.”<sup>5</sup> The point of this problem is to apply the matrix method of Section 1.5 to solve this problem.

(a) Solve the example linear system in Eq. (1.39), using a computer, to obtain a value for  $V_2$ , and thus for  $R_{\text{eff}}$ . I will describe how you can do this in *Mathematica* as an example (which is obnoxiously expensive but handy in that it will give *exact* results using this method), but a similar procedure will work in some other software packages as well, such as the open-source *Octave*.<sup>6</sup>

First in a *Mathematica* window, type “A=”, and then under the menu **Insert > Table/Matrix > New...**, create a  $4 \times 4$  matrix. Then fill in the elements as above. Do the same thing to create a vector ( $4 \times 1$  matrix) for the currents (call it “II”, and use “IR” to represent  $IR$  in the vector; *Mathematica* reserves the symbol “I” for  $\sqrt{-1}$ ). Finally, use `LinearSolve[A, II]` to obtain the solution vector **v**, and use it to deduce the effective resistance.

(b) Modify the calculation to compute the effective resistance between nodes 3 and 4 in the same example circuit (reproduced below).

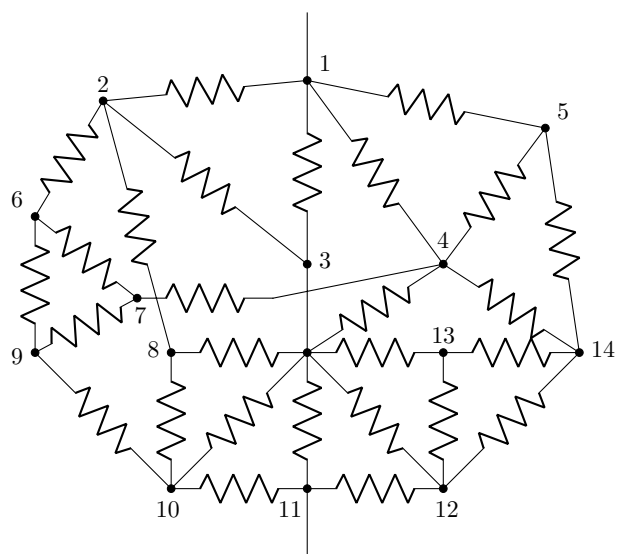


(c) Finally, modify the calculation to handle the XKCD mess.

The network is reproduced below, with nodes labeled (all  $1\text{-}\Omega$  resistors). You should set  $V_{11} = 0$  and let a current  $I$  flow into node 1 (and out node 11).

<sup>5</sup><http://xkcd.com/730/>

<sup>6</sup><http://octave.org>



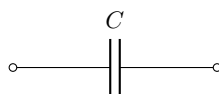


## Chapter 2

# Capacitors and Inductors

### 2.1 Capacitor Basics

A **capacitor** (an older, equivalent term is **condensor**) is, at minimum, a pair of conductors separated by vacuum or dielectric. The symbol itself depicts two parallel electrodes.



Generally, useful capacitors have a relatively large area to achieve a reasonably useful capacitance, so they are something like a pair of planar conductors, with little separation. They may have stacks of many planar conductors to increase the area even further, and the stacks may be rolled up to save space and stored in a can or dipped in epoxy for robustness.

Capacitors act as devices to store *charge*, and in doing so they also store *energy*. The charge stored on the plates generates an electric field and thus a potential difference between the capacitor plates. The potential and charge are related by the capacitor law,

$$Q = CV, \tag{2.1}$$

(capacitor law)

where  $C$  is the **capacitance**, measured in **Farads (F)**, which characterizes the capacitor. Larger conductor areas, “stronger” dielectrics, and smaller electrode spacing all result in larger capacitance. (For example, a few twists to intertwine two pieces of insulated hookup wire makes a capacitor of a few pF.) Differentiating this relation at constant capacitance, and using  $I = dQ/dt$ , we find

$$I = C \frac{dV}{dt}, \tag{2.2}$$

(capacitor charging)

which means that a current *charges* or *discharges* a capacitor, allowing it to build up or bleed off charge, thus changing the voltage.

There are many types of capacitors, and usually the different types are named according to their dielectrics. A few of the most important are:

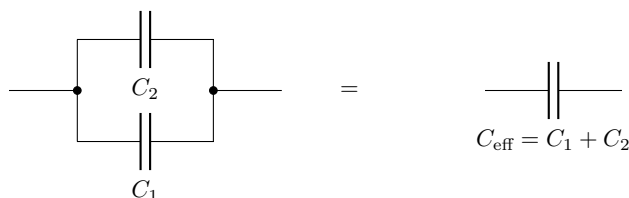
1. **ceramic/monolithic**: these tend to be cheap, and work reasonably well at high frequencies. The capacitances are fairly small in the overall spectrum, ranging from a  $\sim$ pF to about  $\sim$ 0.1  $\mu$ F.
2. **electrolytic**: generally this refers to **aluminum electrolytic capacitors**. In these capacitors, one electrode is aluminum foil, and the other is a liquid electrolyte. Under normal operation, an oxide layer grows on the aluminum and acts as a thin dielectric layer. The foil can be coiled to give large surface area, and the dielectric is very thin, so capacitances can be large, from  $\sim$ 1  $\mu$ F to  $\sim$ 1 F or more,

though these tend not to work as well at high frequencies. These are **polarized**, meaning they can only sustain voltage applied in one direction; the wrong voltage polarity can cause the insulating layer to break down, leading to failure of the capacitor.

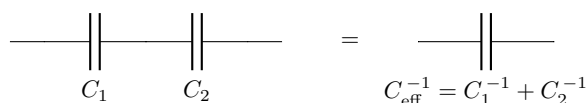
3. **tantalum**: these are also electrolytic capacitors, with a sintered titanium pellet as one electrode, and a solid electrolyte as the other. These typically have intermediate capacitances in the range of  $\sim 0.1 \mu\text{F}$  to  $\sim 10 \mu\text{F}$ , and are fairly compact compared to aluminum electrolytics.

There are many other kinds of capacitors. Some examples include paper, mica, mylar, polystyrene, polypropylene, oil, and niobium electrolytic.

Like resistors, networked capacitors can combine to form equivalent single capacitors. Two capacitors in *parallel* simply add their capacitances,



as resistors in series add. Two *series* capacitors add less straightforwardly,

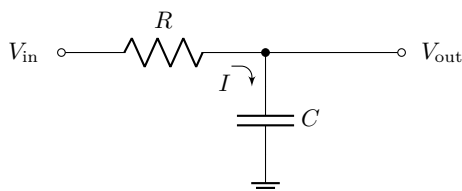


like two *parallel* resistors add.

## 2.2 Simple R-C Circuits

### 2.2.1 Integrator

Consider the simple R-C circuit below. It looks something like a voltage divider, with the lower resistor replaced by a capacitor.



To analyze this (unloaded) circuit, note that all of the current  $I$  going through the resistor must also pass through the capacitor. Applying Ohm's law to determine the current,

$$V_{\text{in}} - V_{\text{out}} = IR, \quad (2.3)$$

and then using the capacitor-charging law (2.2) to relate the current to the output voltage,

$$I = C \frac{dV_{\text{out}}}{dt}. \quad (2.4)$$

Solving for  $dV_{\text{out}}/dt$  and eliminating  $I$ ,

$$\frac{dV_{\text{out}}(t)}{dt} = -\frac{V_{\text{out}}(t)}{RC} + \frac{V_{\text{in}}(t)}{RC}. \quad (2.5) \quad \text{(differential equation for integrator)}$$

Note that by comparing left- and right-hand sides here, we can see that  $RC$  must have the units of time for the units to work out.

### 2.2.1.1 Solution by Integrating Factor

So how do we solve this differential equation? Note that the input  $V_{\text{in}}(t)$  is arbitrary and unknown, so the best we can do is to obtain the solution in terms of  $V_{\text{in}}(t)$ . For this type of **ordinary differential equation (ODE)**, there is a nice trick to simplify it into something more manageable. The idea is to first set  $V_{\text{in}} = 0$ , and look at the equation.

$$\frac{dV_{\text{out}}(t)}{dt} = -\frac{V_{\text{out}}(t)}{RC}. \quad (2.6)$$

This ODE is easy to solve: it's just exponential decay,

$$V_{\text{out}}(t) = V_{\text{out}}(0) e^{-t/RC}. \quad (2.7)$$

Now the idea is that even in the *presence* of  $V_{\text{in}}$ , we should expect more or less the *same* (exponential-decay) behavior. So let's "build this in" to the solution by assuming a solution of the form

$$V_{\text{out}}(t) = \tilde{V}(t) e^{-t/RC}, \quad (2.8)$$

where  $\tilde{V}(t)$  is the "deviation" from the simple solution  $e^{-t/RC}$ , which is called an **integrating factor**. (We aren't losing anything in this assumption, because  $\tilde{V}$  could be anything.) Then solving for  $\tilde{V}$ ,

$$\tilde{V} = V_{\text{out}} e^{t/RC}, \quad (2.9)$$

and differentiating,

$$\frac{d\tilde{V}}{dt} = \frac{dV_{\text{out}}}{dt} e^{t/RC} + \frac{1}{RC} V_{\text{out}} e^{t/RC}. \quad (2.10)$$

Now if we multiply Eq. (2.5) by  $e^{t/RC}$ , and bring both  $V_{\text{out}}$  terms to the left,

$$\frac{dV_{\text{out}}}{dt} e^{t/RC} + \frac{V_{\text{out}}}{RC} e^{t/RC} = \frac{V_{\text{in}}}{RC} e^{t/RC}, \quad (2.11)$$

then notice the left-hand side is the right-hand side of Eq. (2.10). Thus, we have

$$\frac{d\tilde{V}}{dt} = \frac{V_{\text{in}}}{RC} e^{t/RC}. \quad (2.12)$$

Integrating both sides from 0 to  $t$ ,

$$\tilde{V}(t) - \tilde{V}(0) = \frac{1}{RC} \int_0^t V_{\text{in}}(t') e^{t'/RC} dt'. \quad (2.13)$$

Using Eq. (2.9) to get rid of  $\tilde{V}$ , and multiplying through by  $e^{-t/RC}$ ,

$$V_{\text{out}}(t) = V_{\text{out}}(0) e^{-t/RC} + \frac{1}{RC} \int_0^t V_{\text{in}}(t') e^{(t'-t)/RC} dt'. \quad (2.14)$$

(integrator solution)

Thus, we have a solution as an integral to Eq. (2.5) in terms of an arbitrary input voltage.

### 2.2.1.2 Constant Input: Exponential Charging

Now let's take the simple case where the capacitor is initially uncharged [ $V_{\text{out}}(0) = 0$ ], and some **constant** voltage  $V_{\text{in}}$  appears at the input. Then it comes out of the integral, and we have

$$\begin{aligned} V_{\text{out}}(t) &= \frac{V_{\text{in}}}{RC} e^{-t/RC} \int_0^t e^{t'/RC} dt' \\ &= V_{\text{in}} e^{-t/RC} e^{t'/RC} \Big|_0^t, \end{aligned} \quad (2.15)$$

with final solution

$$V_{\text{out}}(t) = V_{\text{in}} \left(1 - e^{-t/RC}\right). \quad (\text{integrator solution, constant input}) \quad (2.16)$$

This is an exponential rise from 0 to  $V_{\text{in}}$ , with **time constant**  $RC$ . That is, looking at  $e^{-t/RC}$  when  $t = RC$ , this falls from unity to  $1/e$ , or about 37%. So  $1 - e^{-t/RC}$ , when  $t = RC$ , rises from 0 to  $1 - 1/e$ , or about 63%.

One important thing to note about this is that the capacitor tends to **smooth** the input. Here, we can regard the problem as an input of a sudden voltage step at  $t = 0$  [which is consistent with  $V_{\text{out}}(0) = 0$ ], and the output is now smoothed over a time scale  $RC$  due to the exponential action of the R-C circuit.

### 2.2.1.3 Integration

This R-C circuit is called a **passive integrator** or **integrator**. The reason is as follows. Going back to Eq. (2.5), if  $V_{\text{out}} \ll V_{\text{in}}$ , then we can ignore the  $V_{\text{out}}$  term:

$$\frac{dV_{\text{out}}(t)}{dt} \approx \frac{V_{\text{in}}(t)}{RC}. \quad (2.17)$$

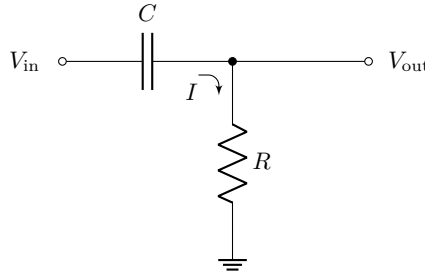
In this case, we can just integrate this equation directly from 0 to  $t$ :

$$V_{\text{out}}(t) - V_{\text{out}}(0) \approx \frac{1}{RC} \int_0^t V_{\text{in}}(t') dt'. \quad (2.18)$$

Thus, the output is the simple integral of the input signal, up to an offset and a factor of  $1/RC$ . In this regime, the capacitor simply stores up the charge that comes in from the resistor. If the capacitor charge, and hence voltage, becomes *too* large, then this reduces the voltage drop across the resistor, and the simple integral approximation breaks down.

### 2.2.2 Differentiator

We can get the other simple R-C circuit by interchanging the resistor and capacitor in the integrator, as shown below.



The capacitor charging equation here gives

$$I = C \frac{d}{dt} (V_{\text{in}} - V_{\text{out}}), \quad (2.19)$$

while Ohm's law gives

$$I = \frac{V_{\text{out}}}{R}. \quad (2.20)$$

Eliminating  $I$ , we get

$$\frac{V_{\text{out}}}{R} = C \frac{d}{dt} (V_{\text{in}} - V_{\text{out}}). \quad (\text{differentiator ODE}) \quad (2.21)$$

which is the ODE for this circuit. This circuit is called a **differentiator**, for reasons analogous to what we saw for the integrator. If  $V_{\text{out}} \ll V_{\text{in}}$ , then the ODE reduces to

$$\frac{V_{\text{out}}}{R} \approx C \frac{dV_{\text{in}}}{dt}, \quad (2.22)$$

and the output is approximately the derivative of the input.

It is possible to write down a general solution to Eq. (2.21) for the differentiator output in terms of an arbitrary input:

$$V_{\text{out}}(t) = V_{\text{in}}(t) + [V_{\text{out}}(0) - V_{\text{in}}(0)]e^{-t/RC} - \frac{e^{-t/RC}}{RC} \int_0^t V_{\text{in}}(t') e^{t'/RC} dt'. \quad (\text{differentiator solution}) \quad (2.23)$$

The idea is fairly similar to the integrator (the same integrating-factor trick applies); we will leave this as an exercise.

## 2.3 AC Signals and Complex Notation

Suppose we have an **ac signal** (alternating-current signal), oscillating at a single frequency  $\omega$ , described by voltage:

$$V(t) = V_0 \cos \omega t. \quad (2.24)$$

Remember to distinguish *angular frequencies* from “regular frequencies.” Angular frequencies are usually represented by  $\omega$ , and are measured in rad/s. “Regular” frequencies are usually represented by  $f$  or  $\nu$ , and are measured in Hz or cycles/s. We can write

$$\omega = 2\pi f \quad (\text{angular frequency}) \quad (2.25)$$

to relate the two. In physics, using the angular frequency saves writing a bunch of factors of  $2\pi$ , but when quoting a physical value, it’s best to stick to regular frequencies. That, is you could quote an angular frequency by saying “ $\omega/2\pi = 100$  Hz” instead of quoting the direct value of  $\omega$  in rad/s.

As an example, let’s go back to the capacitor charging law,

$$I = C \frac{dV}{dt}. \quad (2.26)$$

Then with the above ac signal, we have

$$I(t) = -\omega C V_0 \sin \omega t = \omega C V_0 \cos(\omega t + \pi/2), \quad (2.27)$$

where in the last step, we changed the negative sine to a phase-shifted cosine. This makes the current easier to compare to the original voltage. In fact, we can see that the current **leads** the voltage in phase, because the current’s phase is larger by  $\pi/2$  than the voltage’s phase. The *shape* of the current signal is otherwise the same as the voltage, except for amplitude and phase.

### 2.3.1 Complex Phase

There is a nicer way to handle this monochromatic time dependence, and that is to introduce a complex-number notation. The idea is to represent the time dependence at frequency  $\omega$  by a factor  $e^{-i\omega t}$ .<sup>1</sup> Then define a **complex voltage**

$$\tilde{V} = V_0 e^{-i\omega t}. \quad (2.28)$$

<sup>1</sup>Note that this is a common convention in *physics*; engineers usually use  $e^{j\omega t}$ , where  $j = -i$ . The physics notation comes from the time dependence of a right-going plane wave,  $e^{i(kx - \omega t)}$ .



The *real* (physical) voltage is just the real part of this, or

$$V(t) = \operatorname{Re}[\tilde{V}] = \operatorname{Re}[V_0 e^{-i\omega t}] = V_0 \cos \omega t, \quad (2.29)$$

if  $V_0$  is real. The imaginary part,  $iV_0 \sin \omega t$ , is “carried along” for mathematical convenience, and should be dropped at the end of the calculation to get physical results. Thus, this works for *linear* circuits (in a *nonlinear* circuit, a single frequency would be converted into other frequencies, so this analysis is best for a capacitor or resistor, but less so for a diode).

Why is this representation convenient? First, *other* phases are easy to represent as the phase of the *complex* coefficient  $V_0$ . For example, we can write a phase  $\phi$  as

$$\tilde{V} = V_0 e^{-i\omega t - i\phi}, \quad (2.30)$$

and the real part of this is

$$V(t) = V_0 \cos(\omega t + \phi), \quad (2.31)$$

as we expect for this phase. However, if we absorb the phase  $\phi$  into the *complex* voltage amplitude  $V_0$ , then  $\tilde{V}$  just looks like  $V_0 e^{-i\omega t}$ .

The other nice thing about this complex notation is that derivatives, integrals, and their associated phases are very easy to handle. We are assuming all time dependence is of the form  $e^{-i\omega t}$ . Then a time derivative acts on the phase of any object like

$$\frac{d}{dt} e^{-i\omega t - i\phi} = -i\omega e^{-i\omega t - i\phi}. \quad (2.32)$$

Since this always happens, we can formally identify

$$\frac{d}{dt} \equiv -i\omega, \quad (\text{derivative for monochromatic signals}) \quad (2.33)$$

provided we are discussing monochromatic signals. Then the capacitor-charging rule (2.26) becomes

$$\tilde{I} = -i\omega C \tilde{V} \quad (2.34)$$

in complex notation.

### 2.3.2 Capacitive Reactance

We can interpret the complex form of the capacitor-charging rule in a powerful way. Solving for  $\tilde{V}$ , we find

$$\tilde{V} = \tilde{I} \frac{i}{\omega C}. \quad (2.35)$$

In fact, this looks a lot like Ohm’s law, if we lump everything next to the  $\tilde{I}$  into an effective “resistance,”

$$X_C := \frac{i}{\omega C}. \quad (2.36)$$

(capacitive reactance)

This is called the **capacitive reactance**, and has the same units as resistance. It functions as something like a resistance in the “Ohm’s law for capacitors,”

$$\tilde{V} = \tilde{I} X_C, \quad (2.37)$$

(Ohm’s law for capacitors)

which is the same as Eqs. (2.35) and (2.26). However, because the reactance represents a derivative, it depends on frequency. In fact, what this is saying is that capacitors have very high “resistance” at small frequencies (a capacitor is basically a broken wire, after all, and no current flows once the capacitor is charged), but the capacitor acts like a short circuit at high frequencies, as we will see.

### 2.3.3 Inductive Reactance

The same idea applies to inductors, which satisfy the “inductive-kick law”

$$V = L \frac{dI}{dt}, \quad (2.38)$$

in terms of the **inductance**  $L$ . Switching to complex notation and replacing the derivative, we have

$$\tilde{V} = -i\omega L \tilde{I}, \quad (2.39)$$

in which case we can define an **inductive reactance**

$$X_L := -i\omega L, \quad (2.40)$$

(inductive reactance)

so the inductor law becomes

$$\tilde{V} = X_L \tilde{I}, \quad (2.41)$$

(inductive reactance)

or just Ohm’s law with  $X_L$  standing in for a resistance.

### 2.3.4 Impedance

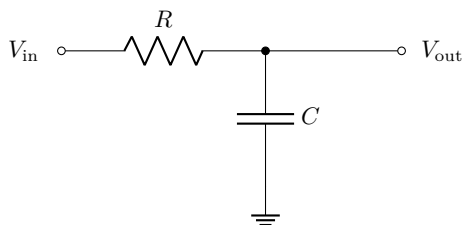
Since reactances and resistances “look” the same once they’re stuck into Ohm’s law, we can use the more general notion of **impedance** to represent any of these. To summarize the important points,

- **Resistances**  $R$  are always *real*.
- **Reactances**  $X_C$ ,  $X_L$  are always *purely imaginary*, and they depend on frequency as well.
- An **impedance** can be any combination of resistances and reactances, and can be any complex value, not necessarily purely real or purely imaginary.
- Resistances and reactances are *special cases* of impedances.

The point of all this: for capacitors, inductors, and resistors in ac circuits, *everything* we did for resistive networks carries over to the ac case, in terms of impedances. This includes all the parallel, series, and Thévenin stuff.

### 2.3.5 Low-Pass Filter

As an example, let’s return to the integrator circuit from Section 2.2.1.



If we think of the capacitor as being a resistor with “resistance”  $X_C$ , then this is just a voltage divider. Using the voltage-divider formula (1.18),

$$\tilde{V}_{out} = \frac{X_C}{R + X_C} \tilde{V}_{in} = \frac{1}{1 - i\omega RC} \tilde{V}_{in}. \quad (2.42)$$

This is a linear relation between the input and output voltage amplitudes, so to simplify the discussion a bit, let's define the **transfer function**

$$\tilde{T}(\omega) := \frac{\tilde{V}_{\text{out}}}{\tilde{V}_{\text{in}}}, \quad (2.43)$$

(transfer function)

which for the low-pass filter is

$$\tilde{T}(\omega) = \frac{1}{1 - i\omega RC} \quad (2.44)$$

(transfer function)

from Eq. (2.42). To simplify even more, we can also consider the **amplitude transfer function**, which discards the phase information:

$$T(\omega) := \left| \frac{\tilde{V}_{\text{out}}}{\tilde{V}_{\text{in}}} \right|. \quad (2.45)$$

(amplitude transfer function)

Then for the low-pass filter, from Eq. (2.42) we have

$$T(\omega) = \frac{1}{\sqrt{1 + (\omega RC)^2}}. \quad (2.46)$$

(amplitude transfer function, low-pass filter)

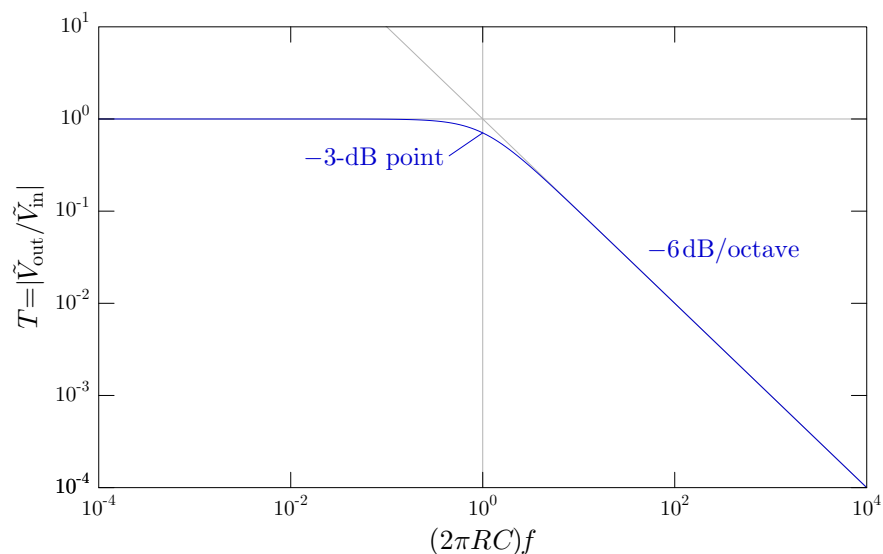
Looking at the asymptotics of the low-pass filter, note that

- As  $\omega \rightarrow 0$  ( $\omega \ll 1/RC$ ), note that  $\tilde{T}(\omega) \rightarrow 1$ , which means that the integrator does not change the signal at low frequencies.
- For large  $\omega$  ( $\omega \gg 1/RC$ ),  $\tilde{T}(\omega) \sim i/\omega RC$ . Note that  $i = e^{i\pi/2}$ , and so  $ie^{-i\omega t} = e^{-i(\omega t - \pi/2)}$ . This means the output phase **lags** the input phase by  $90^\circ$ . Also, the output amplitude is reduced by a factor  $\omega^{-1}$ . This **power-law behavior** appears as a straight line on a log-log plot, with slope  $-1$ .
- The high-frequency scaling of  $\omega^{-1}$  is usually called **-6 dB/octave**. Remember that decibels are defined such that the ratio of two powers in decibels is of the form  $10 \log_{10}(P/P_0)$ , and the ratio of two *amplitudes* is  $20 \log_{10}(V/V_0)$ . One **octave** means a doubling of frequency (from the musical term), and  $\omega^{-1}$  scaling means that doubling the frequency cuts the amplitude in half. In dB, this is  $20 \log_{10}(1/2) = -6$  dB.
- The **transition point** between the low- and high-frequency behavior is called the **3-dB point**, or more properly, the **-3-dB point**. The convention is to define the transition point as the point where  $T(\omega)$  drops from 1 to  $1/\sqrt{2}$  (so that the *power* transferred drops to  $1/2$ ). Then  $20 \log(1/\sqrt{2}) = -3$  dB. We can find the corresponding frequency by setting  $T(\omega_{3\text{ dB}}) = 1/\sqrt{2}$ , which has the solution

$$\omega_{3\text{ dB}} = \frac{1}{RC}, \quad f_{3\text{ dB}} = \frac{1}{2\pi RC}. \quad (2.47)$$

(3-dB frequency)

All this is summarized in the plot below.



Because the integrator “passes” low frequencies without attenuation, and “rolls off” high frequencies, it is called a **low-pass filter**.

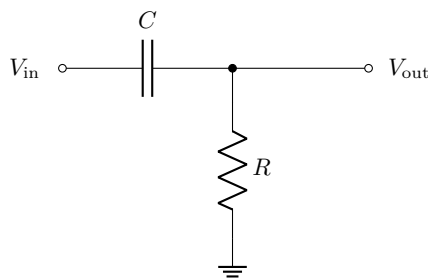
### 2.3.6 Example Problem: Alternate Scaling

What is the scaling of  $-6$  dB/octave, expressed in dB/decade?

**Solution.** This is still a scaling of  $\omega^{-1}$ . A decade is a factor of 10, which means a factor of 10 reduction in amplitude, or  $20 \log_{10}(1/10) = -20$  dB, so  $-20$  dB/decade.

### 2.3.7 Example Problem: High-Pass Filter

Consider the differentiator from Section 2.2.2.



In doing this problem you should see why this is also called a **high-pass filter**.

- Compute  $\tilde{T}(\omega)$ .
- Compute  $T(\omega)$ .
- Work out the low- and high-frequency asymptotics of  $\tilde{T}(\omega)$ .
- Find  $f_{3\text{ dB}}$ .

**Solution.**

- Using the voltage-divider formula again,

$$\tilde{T}(\omega) = \frac{\tilde{V}_{\text{out}}}{\tilde{V}_{\text{in}}} = \frac{R}{R + X_C} = \frac{R}{R + i/\omega C} = \frac{\omega RC}{\omega RC + i}. \quad (2.48)$$

- Taking the modulus,

$$T(\omega) = \frac{\omega RC}{\sqrt{1 + (\omega RC)^2}}. \quad (2.49)$$

(c) Small  $\omega$ :

$$\tilde{T}(\omega) \sim -i\omega RC. \quad (2.50)$$

Since  $-i = e^{-i\pi/2}$ , this *advances* the phase, and is a  $90^\circ$  **phase lead**, like the derivative  $d/dt = -i\omega$  (hence, differentiator).

For large  $\omega$ :

$$\tilde{T}(\omega) \approx 1. \quad (2.51)$$

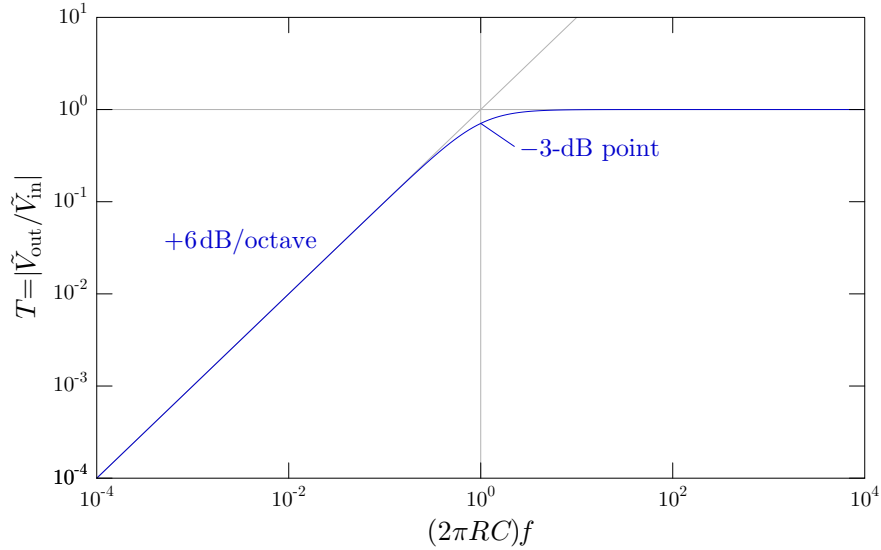
Hence, the high-pass filter.

(d) The 3-dB point occurs where  $T(\omega) = 1/\sqrt{2}$ , which has the same solution

$$f_{3\text{ dB}} = \frac{1}{2\pi RC} \quad (2.52)$$

as the low-pass filter.

Overall, the plot for the high-pass filter is basically a mirror image of the plot for the low-pass filter.



## 2.4 Phase

One thing that we haven't paid much attention to yet is the **phase shift** due to a linear circuit. (We did this just a bit, in looking at the low- and high-pass filters, in the asymptotic limits of low and high frequency, where the phase shift ended up being nothing or  $\pm 90^\circ$ .) In general, the complex transfer function [from Eq. (2.43)],

$$\tilde{T}(\omega) := \frac{\tilde{V}_{\text{out}}}{\tilde{V}_{\text{in}}}, \quad (2.53)$$

gives information about both the amplitude [in  $T(\omega)$ ] and phase (via the complex phase). In other words, we can always write

$$\tilde{T}(\omega) := T(\omega) e^{-i\phi(\omega)}, \quad (2.54)$$

where  $\phi(\omega)$  is the frequency-dependent phase shift (remember the minus sign here is because of the phase convention in  $e^{-i\omega t}$ ). If  $\phi > 0$  at some frequency, we call this a **phase lead**, whereas  $\phi < 0$  is a **phase lag**.

Now remember that for an arbitrary complex number  $z$ , we can write it in polar and cartesian forms as

$$z = r e^{-i\phi} = r \cos \phi - ir \sin \phi =: x + iy, \quad (2.55)$$

where  $r$ ,  $\phi$ ,  $x$ , and  $y$  are all real. Equating real and imaginary parts, we get  $x = r \cos \phi$  and  $y = -r \sin \phi$ , and dividing these equations gives

$$\frac{y}{x} = -\frac{\sin \phi}{\cos \phi} = -\tan \phi. \quad (2.56)$$

Solving for  $\phi$  gives

$$\phi = -\tan^{-1} \frac{y}{x}, \quad (2.57)$$

remembering that  $\tan x$  is an odd function. Now  $x$  and  $y$  are respectively the real and imaginary parts of  $\tilde{T}(\omega)$ , so

$$\phi(\omega) = -\tan^{-1} \left( \frac{\text{Im}[\tilde{T}(\omega)]}{\text{Re}[\tilde{T}(\omega)]} \right). \quad (2.58)$$

(phase shift of linear circuit)

Thus, we have the phase shift at any frequency in terms of the complex transfer function.

### 2.4.1 Example: Low-Pass Filter

In the low-pass filter, we had from Eq. (2.44)

$$\tilde{T}(\omega) = \frac{1}{1 - i\omega RC}. \quad (2.59)$$

Multiplying upstairs and downstairs by  $1 + i\omega RC$ , we can rewrite this as

$$\tilde{T}(\omega) = \frac{1 + i\omega RC}{1 + (\omega RC)^2}. \quad (2.60)$$

Now the real and imaginary parts are more obvious, and if we put these into Eq. (2.58), we get

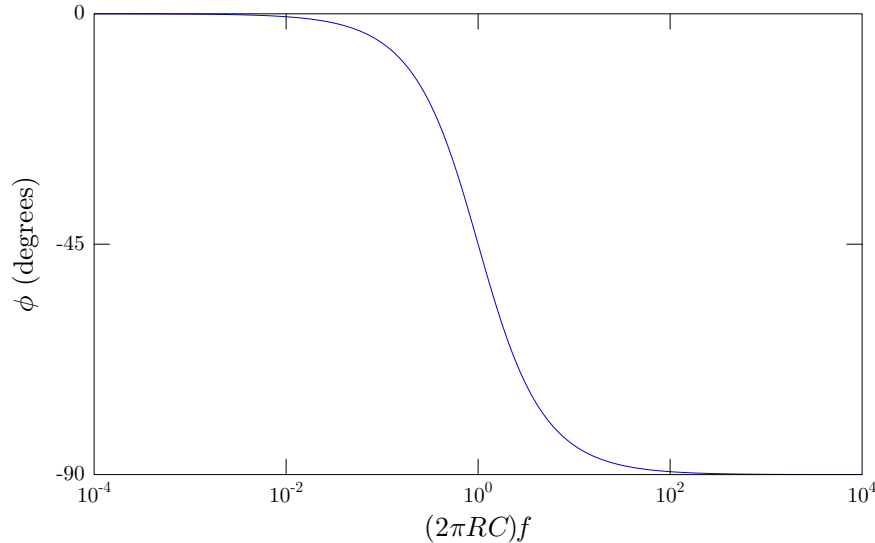
$$\phi(\omega) = -\tan^{-1}(\omega RC). \quad (2.61)$$

(phase shift, low-pass filter)

How does this behave? In the extreme limits:

- For small frequencies [ $\omega \ll (RC)^{-1}$ ],  $\phi \approx -\omega RC$ , which is a small (close to  $0^\circ$ ) **phase lag**.
- For large frequencies [ $\omega \gg (RC)^{-1}$ ],  $\phi \approx -\pi/2$ , which is a  $90^\circ$  **phase lag**.

In between, the phase lag moves smoothly between  $0$  and  $90^\circ$ , as shown in the plot below.



The high-pass filter is very similar, but the phase is  $90^\circ$  for small frequencies (in the “stop band”), and changes to  $0^\circ$  for large frequencies (in the “pass band”).

## 2.5 Power

Now we arrive at the real meaning of why capacitors and inductors seem a lot like resistors, but with “imaginary resistance.” Remember that the power dissipated in a circuit is  $P = IV$ . This is still true in an ac circuit, but

$$P(t) = I(t) V(t) \quad (2.62)$$

is the *instantaneous* power. This can either be positive or negative; positive means dissipating energy or perhaps storing energy in a capacitor or inductor, while negative means we are getting some stored energy back.

What we’re interested in here is the *time-averaged* behavior, so we know the *net* effect of everything that happens over a cycle. So let’s assume a monochromatic voltage and a current, with a possible phase shift  $\phi$  in the current. In real notation, we have

$$\begin{aligned} V(t) &= V_0(\cos \omega t) \\ I(t) &= I_0 \cos(\omega t + \phi). \end{aligned} \quad (2.63)$$

Then the power, time-averaged over one period  $T = 2\pi/\omega$  is

$$\langle P \rangle = \frac{1}{T} \int_0^T V(t) I(t) dt. \quad (2.64)$$

Expanding the cosine in the current using the sum-angle formula,

$$I(t) = I_0[\cos(\omega t) \cos(\phi) - \sin(\omega t) \sin(\phi)] \quad (2.65)$$

and using this with the above expression for  $V(t)$ , we get

$$\langle P \rangle = \frac{I_0 V_0}{T} \int_0^T [\cos^2(\omega t) \cos(\phi) - \cos(\omega t) \sin(\omega t) \sin(\phi)] dt. \quad (2.66)$$

In the second term the  $\cos(\omega t) \sin(\omega t) = (1/2) \sin(2\omega t)$  averages to zero. In the first term  $\cos^2(\omega t) = (1/2) + (1/2) \cos(2\omega t)$ , which averages to just  $1/2$ . Thus,

$$\langle P \rangle = \frac{1}{2} I_0 V_0 \cos(\phi). \quad (2.67)$$

Now to simplify this a bit more, we want to compare this to time-averaged values of  $V(t)$  and  $I(t)$  separately. It doesn’t make sense to time average them directly, because they average to zero. But we can compute the **rms** or **root-mean-square** values. This just means: square it, time-average it, take the square root, done. For the voltage, if we square it,

$$V^2(t) = V_0^2 \cos^2 \omega t, \quad (2.68)$$

then average it,

$$\frac{1}{T} \int_0^T V^2(t) dt = \frac{V_0^2}{2}, \quad (2.69)$$

then take the square root, we get the rms voltage:

$$V_{\text{rms}} = \frac{V_0}{\sqrt{2}}. \quad (2.70)$$

(rms voltage)

Similarly, for current,

$$I_{\text{rms}} = \frac{I_0}{\sqrt{2}}. \quad (2.71)$$

(rms current)

Note that these two rms expressions are valid *only* for sine waves. Then we can rewrite Eq. (2.67) as

$$\langle P \rangle = I_{\text{rms}} V_{\text{rms}} \cos(\phi). \quad (2.72)$$

(rms current)

It is common to define the **power factor** as

$$\text{power factor} := \frac{\langle P \rangle}{I_{\text{rms}} V_{\text{rms}}} = \cos(\phi). \quad (2.73)$$

(power factor)

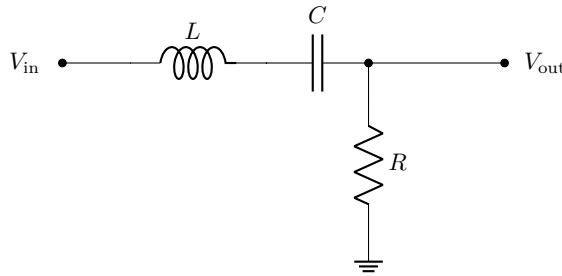
What does this mean? Breaking this down into cases:

- $\cos \phi = 1$  ( $\phi = 0$ ): this occurs for a purely resistive load; the maximum power is dissipated here.
- $\cos \phi = 0$  ( $|\phi| = \pi/2$ ): this is a purely reactive (capacitive or inductive) load, *no power is dissipated* (it is only stored and retrieved over each cycle).
- $0 < \cos \phi < 1$  ( $0 < |\phi| < \pi/2$ ): there is some reactive component to the load impedance, so some power is dissipated, but not as much as for an equivalently large but real impedance.
- $\cos \phi < 0$  ( $|\phi| > \pi/2$ ): it's also possible to have a *negative* power factor, which means the “load” is in fact a generator or EMF source.

Light bulbs and toasters are good examples of resistive loads. An example of a capacitive load is a piezo speaker or buzzer. Examples of inductive (reactive) loads are electric motors or lighting transformers or ballasts for fluorescent lights (the “magnetic” kind, not the “electronic” kind). Inductive loads are important in high-power applications, and the problem is that the power factor can be very small, so that large voltages and currents are needed to drive a motor. This isn't efficient, because the large voltages and currents cause wasted power to be dissipated *elsewhere*. A trick to help here is to “correct” the phase of the load impedance by connecting a parallel capacitor, which increases the power factor. It is common in air-conditioning compressor motors to have *two* capacitors, a “start” capacitor and a “run” capacitor. The start capacitor increases the parallel capacitance and hence the power factor when the motor first powers on, to give the motor extra startup torque.

## 2.6 Resonant Circuits

As another example of mixed impedances, we can consider a resonant filter with a resistor, inductor, and capacitor, as shown below.



Resonant LC circuits like this are also often called **tank circuits**. To analyze this, let's lump the series inductor and capacitor together into a single element, of impedance

$$\begin{aligned} Z_{LC} &= X_L + X_C = -i\omega L + \frac{i}{\omega C} \\ &= -\frac{iL}{\omega} \left( \omega^2 + \frac{1}{LC} \right) \\ &= -\frac{iL}{\omega} (\omega^2 - \omega_0^2), \end{aligned} \quad (2.74)$$



where we have defined the **resonant (LC) frequency**

$$\omega_0 := \frac{1}{\sqrt{LC}}. \quad (2.75)$$

(LC frequency)

Then we can treat what is left as a voltage divider. The transfer function is

$$\begin{aligned} \tilde{T}(\omega) &= \frac{R}{R + Z_{LC}} \\ &= \frac{R}{R - (iL/\omega)(\omega^2 - \omega_0^2)} \\ &= \frac{i\omega R/L}{(\omega^2 - \omega_0^2) + i\omega R/L}. \end{aligned} \quad (2.76)$$

Defining the damping constant

$$\gamma := \frac{R}{L}, \quad (2.77)$$

the transfer function becomes

$$\tilde{T}(\omega) = \frac{i\omega\gamma}{(\omega^2 - \omega_0^2) + i\omega\gamma}. \quad (2.78)$$

(transfer function, RLC circuit)

This is the same as the response function for a damped harmonic oscillator, with resonant frequency  $\omega_0$  and damping rate  $\gamma$ .

$$\tilde{T}(\omega) = \frac{\omega\gamma}{\sqrt{(\omega^2 - \omega_0^2)^2 + \omega^2\gamma^2}}. \quad (2.79)$$

(amplitude transfer function, RLC circuit)

Note that on resonance ( $\omega = \omega_0$ ),  $T(\omega = \omega_0) = 1$  [and in fact  $\tilde{T}(\omega = \omega_0) = 1$ , so the signal is transmitted without amplitude reduction or any phase shift. Away from resonance,  $T(\omega) < 1$ , leading to a transmission “peak” around  $\omega_0$ .

### 2.6.1 Q Factor

A common way to quantify the width of the resonance peak is the **Q factor**. The idea is to find the  $-3$  dB points of the resonance peaks (compared to the peak) as a measure of the width. Thus, setting  $T(\omega) = 1/\sqrt{2}$ ,

$$\frac{1}{\sqrt{2}} = \frac{\omega\gamma}{\sqrt{(\omega^2 - \omega_0^2)^2 + \omega^2\gamma^2}}, \quad (2.80)$$

we can square this and rearrange to find

$$(\omega^2 - \omega_0^2)^2 = \omega^2\gamma^2. \quad (2.81)$$

This is the square of a quadratic equations, with four solutions

$$\omega = \pm \sqrt{(\gamma/2)^2 + \omega_0^2} \pm \frac{\gamma}{2}. \quad (2.82)$$

We only want the positive solutions, because  $T(\omega) \geq 0$ , and  $T(\omega)$  has the same sign as  $\omega$ . Thus,

$$\omega_{3\text{dB}} = \sqrt{(\gamma/2)^2 + \omega_0^2} \pm \frac{\gamma}{2}. \quad (2.83)$$

If we define the width of the peak to be the difference between the  $-3$ -dB points,

$$\delta\omega_{3\text{dB}} = \gamma. \quad (2.84)$$

(full width at half maximum)

This is also called the **full width at half maximum (FWHM)** (“half” here refers to the *power* transmission  $|T(\omega)|^2$ ).

$$\omega^2 - \omega_0^2 = \omega\gamma. \quad (\text{amplitude transfer function, RLC circuit}) \quad (2.85)$$

Then we can write down the  $Q$  factor as the ratio of the resonance frequency to the FWHM:

$$Q = \frac{\omega_0}{\delta\omega_{3\text{dB}}} = \frac{\omega_0}{\gamma}. \quad (Q \text{ factor, RLC circuit}) \quad (2.86)$$

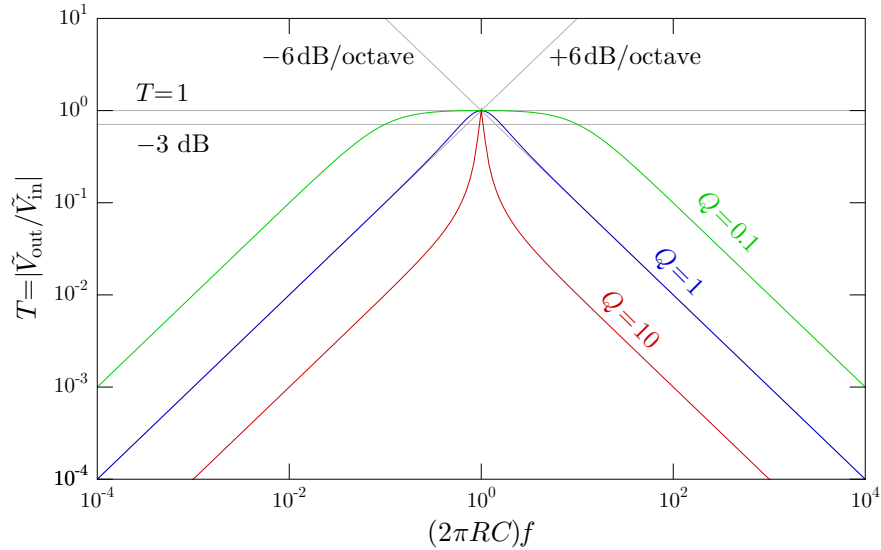
The “ $Q$ ” here indicates the “quality” of the resonator: a large  $Q$  means a small  $\gamma$  compared to  $\omega_0$ , and thus a narrow resonance. Actually, this is *not* the definition of the  $Q$  factor, but it’s the one that physicists use when dealing with resonances—we’ll come back to the formal definition below.

Also, it is useful to write the  $Q$  factor in terms of the electronic parameters. Using Eqs. (2.75) and (2.77),

$$Q = \frac{\sqrt{L}}{R\sqrt{C}}. \quad (Q \text{ factor, RLC circuit}) \quad (2.87)$$

Thus, for example, an increased capacitance not only lowers the resonance frequency, but also the  $Q$  factor (because the  $Q$  factor depends on the resonance frequency).

The transfer function  $T(\omega)$  is plotted below for three different values of  $Q$ .



Note that asymptotically, the filter rises/falls at  $\pm 6$  dB/octave, like the low- and high-pass filters.

### 2.6.1.1 Fundamental Definition

Now to justify the simple formula (2.86) for the  $Q$  factor that we used. The more fundamental definition of the  $Q$  factor is defined in terms of stored energy in the oscillating circuit and the energy dissipated as follows:

$$Q := 2\pi \cdot \frac{(\text{maximum}) \text{ stored energy}}{\text{energy loss per oscillation cycle}}. \quad (\text{definition of } Q \text{ factor}) \quad (2.88)$$

The capacitor stored energy is, on average,

$$\langle E_C \rangle = \left\langle \frac{1}{2} CV^2 \right\rangle = \frac{1}{2} CV_{\text{rms}}^2, \quad (2.89)$$

while the inductor's stored energy is

$$\langle E_L \rangle = \left\langle \frac{1}{2} L I^2 \right\rangle = \frac{1}{2} L I_{\text{rms}}^2. \quad (2.90)$$

These expressions are actually the same; using  $I = C(dV/dt)$  and time-averaging gives  $I_{\text{rms}}^2 = \omega^2 C^2 V_{\text{rms}}^2$ , and so, for example,

$$\langle E_C \rangle = \frac{1}{2} C V_{\text{rms}}^2 = \frac{1}{2\omega_0^2 C} I_{\text{rms}}^2 = \frac{1}{2} L I_{\text{rms}}^2, \quad (2.91)$$

where we used  $\omega^2 = \omega_0^2 = 1/LC$ . The power dissipated is

$$\langle P \rangle = I_{\text{rms}} V_{\text{rms}}, \quad (2.92)$$

and so the energy lost per cycle is

$$T \langle P \rangle = T I_{\text{rms}} V_{\text{rms}}, \quad (2.93)$$

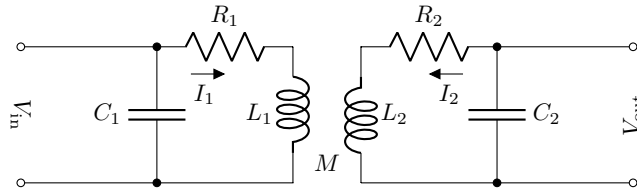
where  $T$  is the period. Putting this together,

$$Q = 2\pi \left( \frac{L I_{\text{rms}}^2}{T I_{\text{rms}} V_{\text{rms}}} \right) = \omega_0 \left( \frac{L I_{\text{rms}}}{V_{\text{rms}}} \right) = \omega_0 \frac{L}{R} = \frac{\omega_0}{\gamma}, \quad (2.94)$$

in agreement with Eq. (2.86).

## 2.7 Coupled, Resonant Circuits

Some more interesting behavior occurs if we have a *pair* of resonant circuits, coupled via a **transformer** (a pair of inductors whose fields couple to each other). Consider the circuit below, consisting of two such circuits. In the resonant case that we will analyze, this circuit is useful for boosting the amplitude of ac signals, similar to ordinary transformer circuits. This is also the key to how the Tesla coil (Section 2.8.1) works.



The two inductors are coupled by a **mutual inductance**  $M$ , such that, for example, the voltage drop across the inductor is

$$V_{L_1} = L_1 \frac{dI_1}{dt} + M \frac{dI_2}{dt}. \quad (2.95)$$

Let  $q_1$  and  $q_2$  be the charges on capacitors  $C_1$  and  $C_2$ , respectively (we will reserve the symbols  $Q_1$  and  $Q_2$  for the  $Q$  factors of each circuit later).

First, we will derive a coupled pair of differential equations for  $q_1$  and  $q_2$ . Starting with the voltage across the input capacitor, we have

$$V_{\text{in}} = \frac{q_1}{C_1}, \quad (2.96)$$

and the equivalent voltage via the inductor is

$$V_{\text{in}} = I_1 R_1 + L_1 \dot{I}_1 + M \dot{I}_2. \quad (2.97)$$

Equating these and using  $I_{1,2} = -\dot{q}_{1,2}$  (given the orientations of the currents, which tend to drain charge from the capacitors) gives

$$\frac{q_1}{C_1} = -\dot{q}_1 R_1 - L_1 \ddot{q}_1 - M \ddot{q}_2. \quad (2.98)$$

Then simplifying, we have

$$\begin{aligned}\omega_1^2 q_1 + \gamma_1 \dot{q}_1 + \ddot{q}_1 + \frac{M}{L_1} \ddot{q}_2 &= 0 \\ \omega_2^2 q_2 + \gamma_2 \dot{q}_2 + \ddot{q}_2 + \frac{M}{L_2} \ddot{q}_1 &= 0,\end{aligned}\tag{2.99}$$

where we have written the matching equation for the output circuit. Here, the resonant frequencies of the circuits are

$$\omega_{1,2} := \frac{1}{\sqrt{L_{1,2}C_{1,2}}},\tag{2.100}$$

and the damping rates are

$$\gamma_{1,2} := \frac{R_{1,2}}{L_{1,2}}.\tag{2.101}$$

Considering a sinusoidal input signal at frequency  $\omega$ , The differential equations (2.99) also apply to complex charges  $\tilde{q}_{1,2}(t)$ , with time dependence of the form

$$\tilde{q}_{1,2}(t) = q_{1,2} e^{-i\omega t}.\tag{2.102}$$

This convention allows us to eliminate the derivatives in the coupled ODEs. That is, in complex notation, we just replace every derivative by  $(-i\omega)$ , so that

$$\begin{aligned}\omega_1^2 \tilde{q}_1 - i\omega\gamma_1 \tilde{q}_1 - \omega^2 \tilde{q}_1 - \omega^2 \frac{M}{L_1} \tilde{q}_2 &= 0 \\ \omega_2^2 \tilde{q}_2 - i\omega\gamma_2 \tilde{q}_2 - \omega^2 \tilde{q}_2 - \omega^2 \frac{M}{L_2} \tilde{q}_1 &= 0.\end{aligned}\tag{2.103}$$

These equations become

$$\begin{aligned}(\omega_1^2 - \omega^2 - i\omega\gamma_1) q_1 - \omega^2 \frac{M}{L_1} \tilde{q}_2 &= 0 \\ (\omega_2^2 - \omega^2 - i\omega\gamma_2) q_2 - \omega^2 \frac{M}{L_2} \tilde{q}_1 &= 0.\end{aligned}\tag{2.104}$$

after a bit of simplification. Note that in matrix form, we can rewrite these equations as

$$\begin{bmatrix} (\omega_1^2 - \omega^2 - i\omega\gamma_1) & -\omega^2(M/L_1) \\ -\omega^2(M/L_2) & (\omega_2^2 - \omega^2 - i\omega\gamma_2) \end{bmatrix} \begin{bmatrix} \tilde{q}_1 \\ \tilde{q}_2 \end{bmatrix} = 0. \quad \text{(evolution equations)}\tag{2.105}$$

For any matrix equation of the form

$$\begin{bmatrix} A & B \\ C & D \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 0\tag{2.106}$$

to hold the determinant of the  $2 \times 2$  matrix must vanish; i.e.,  $AD - BC = 0$ . Thus, we must have

$$(\omega_1^2 - \omega^2 - i\omega\gamma_1)(\omega_2^2 - \omega^2 - i\omega\gamma_2) = \omega^4 k^2,\tag{2.107}$$

(frequency condition)

where

$$k := \frac{M}{\sqrt{L_1 L_2}}\tag{2.108}$$

is the **coupling coefficient** for the inductor pair (satisfying  $0 \leq k \leq 1$ ). This is a quartic equation in  $\omega$ ; the physical solutions determine the oscillation frequencies of the coupled-oscillator system. For now, note that in the limit of zero coupling  $k$  and zero damping  $\gamma_{1,2}$ , the solutions to this equation are  $\omega = \pm\omega_0$  (note that the frequencies will appear as positive/negative pairs; the interpretation is that only the positive solution is physical). In the more general case, there will be *two* physical frequencies, and in general the system

will oscillate at both frequencies simultaneously—because this is a linear system, the general solution is a superposition of the individual frequencies.

The general condition (2.107) is somewhat cumbersome to analyze in general, so we will simplify to the ideal-resonator limit, where  $R_{1,2} = 0$  (and thus  $\gamma_{1,2} = 0$ ), for matched resonators ( $\omega_1 = \omega_2 =: \omega_0$ ). In this case, Eq. (2.107) becomes

$$(\omega_0^2 - \omega^2)^2 = \omega^4 k^2. \quad (2.109)$$

Again, this condition has four solutions, but the only two positive (and thus physical) frequencies are

$$\omega_{\pm} = \frac{\omega_0}{\sqrt{1 \pm k}}. \quad (2.110)$$

Now Eqs. (2.104) become

$$\begin{aligned} (\omega_0^2 - \omega_{\pm}^2) \tilde{q}_1 &= \omega_{\pm}^2 \frac{M}{L_1} q_2 \\ (\omega_0^2 - \omega_{\pm}^2) \tilde{q}_2 &= \omega_{\pm}^2 \frac{M}{L_2} q_1 \end{aligned} \quad (2.111)$$

at the permitted frequencies. Dividing these equations gives

$$\left| \frac{\tilde{q}_1}{\tilde{q}_2} \right| = \sqrt{\frac{L_2}{L_1}}. \quad (2.112)$$

Using  $\tilde{q}_{1,2} = C_{1,2} \tilde{V}_{1,2}$ , we then have a circuit voltage “gain” of

$$G := \left| \frac{\tilde{V}_{\text{out}}}{\tilde{V}_{\text{in}}} \right| = \frac{q_2 C_1}{q_1 C_2} = \frac{\sqrt{L_1} C_1}{\sqrt{L_2} C_2} = \sqrt{\frac{C_1}{C_2}} = \sqrt{\frac{L_2}{L_1}}. \quad (2.113)$$

(resonant “gain”)

Here, we used  $L_1 C_1 = L_2 C_2$ , since the resonances are matched ( $\omega_1 = \omega_2$ ). This shows that the coupled resonances have a large step-up effect in voltage if the inductance of the second coil is large. Recall that for an ideal solenoid  $L \propto N^2$ , so that the gain appears to be just  $N_2/N_1$ , as we might expect for a transformer circuit. However, note that there are extra geometric factors involved that may differ between the two inductors that obscure this interpretation. Note also that this solution is independent of the coupling coefficient  $k$ —a side effect of the coincident circuit resonances.

Besides the step-up, there is some dynamical behavior. Recall that the oscillation frequencies are  $\omega_{\pm}$ , and so the general solution is

$$\tilde{q}_1(t) = A e^{-i\omega_+ t} + B e^{-i\omega_- t}, \quad (2.114)$$

for example, for the part of the solution applying to the first circuit. If we assume that  $q_1(0) = q_0$  is real, then the coefficients are real and equal, and we can assume a real solution of the form

$$q_1(t) = \frac{q_0}{2} (\cos \omega_+ t + \cos \omega_- t). \quad (2.115)$$

Now to obtain a more precise version of Eq. (2.112), note that Eq. (2.110) may be written as  $\omega^2/\omega_0^2 = 1 \pm k$ , and so the matrix equation (2.105) may be rewritten

$$\begin{bmatrix} \pm k & -M/L_1 \\ -M/L_2 & \pm k \end{bmatrix} \begin{bmatrix} \tilde{q}_1 \\ \tilde{q}_2 \end{bmatrix} = 0, \quad (2.116)$$

corresponding to frequencies  $\omega_{\pm}$ . This means that at frequency  $\omega_{\pm}$ ,  $\pm k \tilde{q}_1 = (M/L_1) \tilde{q}_2$ , so

$$\tilde{q}_2 = \pm \frac{k L_1}{M} \tilde{q}_1 = \pm \sqrt{\frac{L_1}{L_2}} \tilde{q}_1 = \pm G \tilde{q}_1. \quad (2.117)$$

Thus the solution for  $q_2$  corresponding to Eq. (2.115) is

$$q_2(t) = \frac{q_0 G}{2} (\cos \omega_+ t - \cos \omega_- t), \quad (2.118)$$

because the solution is larger by an overall factor  $G$ , and the phase reverses on the  $\omega_-$  component of the solution. Now examining the standard trigonometric relations

$$\begin{aligned}\cos \alpha \cos \beta &= \frac{1}{2} \cos(\alpha - \beta) + \frac{1}{2} \cos(\alpha + \beta) \\ \sin \alpha \sin \beta &= \frac{1}{2} \sin(\alpha - \beta) + \frac{1}{2} \sin(\alpha + \beta),\end{aligned}\tag{2.119}$$

we can rewrite the above solutions as

$$\begin{aligned}q_1(t) &= q_0 \cos(\omega_B t/2) \cos(\bar{\omega} t) \\ q_2(t) &= q_0 G \sin(\omega_B t/2) \sin(\bar{\omega} t),\end{aligned}\tag{2.120}$$

(beating solution)

where the “beat frequency”  $\omega_B$  is

$$\omega_B := \omega_- - \omega_+ = \omega_0 \left( \frac{1}{\sqrt{1-k}} - \frac{1}{\sqrt{1+k}} \right),\tag{2.121}$$

and the mean frequency  $\bar{\omega}$  is

$$\bar{\omega} := \frac{\omega_- + \omega_+}{2} = \frac{\omega_0}{2} \left( \frac{1}{\sqrt{1-k}} + \frac{1}{\sqrt{1+k}} \right).\tag{2.122}$$

The interpretation of Eqs. (2.120) is more straightforward in the weak-coupling limit  $k \ll 1$ , where

$$\omega_B = k\omega_0 + O(k^2), \quad \bar{\omega} = \omega_0 + O(k^2).\tag{2.123}$$

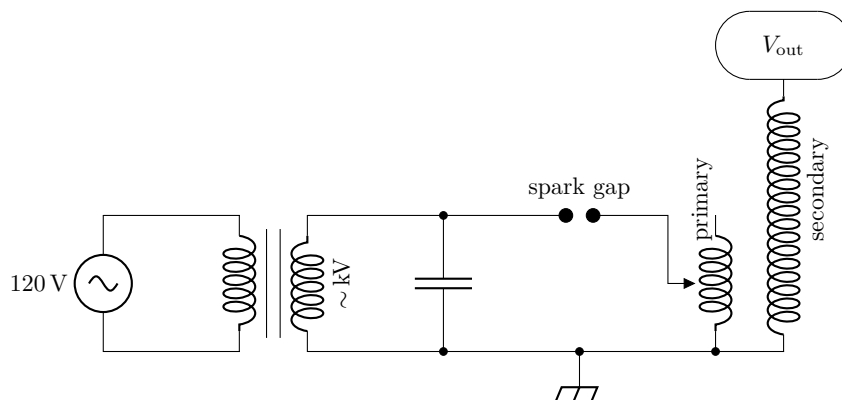
Here, the coupled system oscillates at the original frequency  $\bar{\omega} \approx \omega_0$ , with the two circuits oscillating  $90^\circ$  out of phase. On top of this, the amplitudes are modulated slowly at half the beat frequency. The amplitudes are likewise out of phase, representing the complete transfer of energy between the two circuits. (This is the reason for the factor of 2 in defining the beat frequency, which is the frequency at which the amplitude is maximized in one of the circuits, ignoring any signs in the amplitude.) For larger  $k$ , this still happens, but the mean frequency increases monotonically and arbitrarily as  $k \rightarrow 1$ , as does the beat frequency. Note that for stronger couplings, the interpretation of complete transfer may not hold, if the beat and mean frequencies become comparable.

## 2.8 Circuit Practice

### 2.8.1 Tesla Coil

A nice example of a resonant circuit is the **Tesla coil**. There are many variations, but the simplest circuit is below. The idea is to put in a fairly “low” voltage (line voltage), and generate an output voltage as large as possible, preferably at least hundreds of kV or even MV.

Here we will go through a basic tesla coil circuit to give you some experience in “reading” a more complex circuit diagram.



Typical designs use a neon-sign transformer in the first stage to boost the voltage to  $\sim 6$  kV or more. This drives an RLC circuit (the “R” is the wire resistance), which oscillates at a tunable resonant frequency (by tuning the variable inductor). A spark gap also interrupts the RLC circuit; the nonlinear sparking increases  $dI/dt$  to help get a larger output voltage in the secondary transformer. The secondary coil is just a long air-core coil with many turns, with a large electrode on top. The capacitance of the secondary circuit is the capacitance due to the windings of the coil, as well as the output terminal (the other electrode being ground). So the secondary circuit is also an RLC oscillator. The primary is tuned to resonate with the secondary. The oscillators are generally tuned to radio frequencies (hundreds of kHz), which makes the output (relatively) safe for humans.

The coupled-oscillator analysis of Section 2.7 applies here, such that the voltage step-up factor can be written as either  $\sqrt{C_{\text{primary}}/C_{\text{secondary}}}$  or  $\sqrt{L_{\text{secondary}}/L_{\text{primary}}}$ . Thus, for example, a large primary capacitance leads to a large output voltage. A simple, heuristic derivation of this same result comes from noting that the maximum energy stored in the primary capacitor has the form  $E_1 = (1/2)C_{\text{primary}}V_1^2$ . If we assume all the energy is transferred to the secondary, then we can equate this to the analogous relation for the secondary capacitance energy, obtaining that  $V_2/V_1 = \sqrt{C_{\text{primary}}/C_{\text{secondary}}}$ . However, note that this circuit does not act as a simple transformer—for example, on resonance the secondary circuit acts as a quarter-wave antenna (so that the node of a standing wave occurs at ground, and an antinode occurs at the output terminal), so that the voltage and current are not constant through the secondary winding.

## 2.9 Exercises

### Problem 2.1

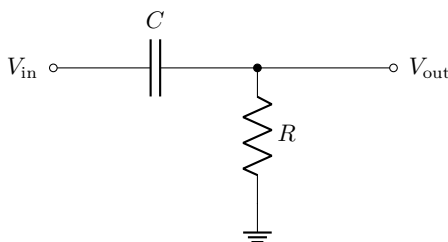
Use the capacitor law  $Q = CV$  to show that the effective capacitance of two *parallel* capacitors is  $C_{\text{eff}} = C_1 + C_2$ , where  $C_1$  and  $C_2$  are the capacitances of the individual capacitors.

### Problem 2.2

Show that the effective capacitance  $C_{\text{eff}}$  of two capacitors  $C_1$  and  $C_2$  in *series* is given by  $C_{\text{eff}}^{-1} = C_1^{-1} + C_2^{-1}$ .

### Problem 2.3

Consider the differentiator circuit shown below, where  $V_{\text{in}}(t)$  is an arbitrary input voltage.



(a) Show that the differential equation for this circuit is given by

$$\frac{dV_{\text{out}}}{dt} = -\frac{V_{\text{out}}}{RC} + \frac{dV_{\text{in}}}{dt}. \quad (2.124)$$

(We did this in class.)

(b) Use the integrating-factor trick that we used for the integrator circuit (i.e., define  $\tilde{V} := V_{\text{out}}e^{t/RC}$ , simplify the equation, and integrate from 0 to  $t$ ) to find the general solution

$$V_{\text{out}}(t) = V_{\text{in}}(t) + [V_{\text{out}}(0) - V_{\text{in}}(0)]e^{-t/RC} - \frac{e^{-t/RC}}{RC} \int_0^t V_{\text{in}}(t') e^{t'/RC} dt'. \quad (2.125)$$

(c) Write down the solution in the case where  $V_{\text{in}}$  is turned on suddenly to a constant value from 0, *just after*  $t = 0$ . Give a brief, qualitative description of the solution (use a sketch if you need to).

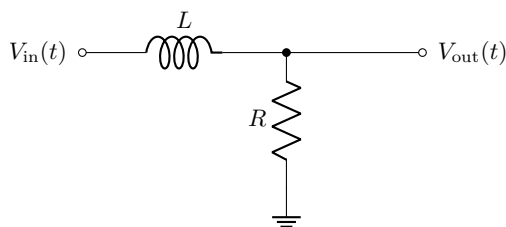
*Note:* the reason for assuming the turn-on comes just after  $t = 0$  is a bit technical, but it's necessary to get the boundary terms in the general solution to come out right. Basically, the “leading edge” of the input signal is an “interesting” part of the signal (the most interesting, actually), and so we must make sure to capture it in our interval of integration. Equivalently, we could have taken our integration range from  $-\infty$  to  $t$ , taking our boundary condition to be  $V_{\text{in}}(-\infty) = 0$ .

*Another note:* From your solution here, you should see the reason why this circuit is called a “dc block.”

### Problem 2.4

Consider the LR circuit below.

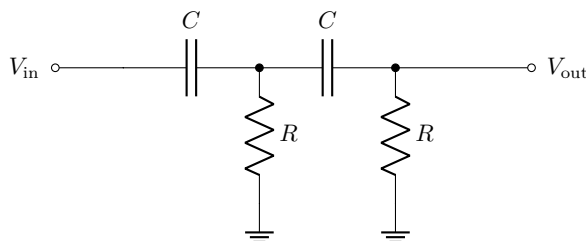




- Is this a high-pass or a low-pass filter? Give a *qualitative* argument for your answer.
- Write down a differential equation relating  $V_{in}(t)$  and  $V_{out}(t)$ .
- Solve the equation for an arbitrary input  $V_{in}(t)$  (*not* necessarily a single frequency!), by using the integrating-factor method. You should end up with a solution in the form of an integral over  $V_{in}(t)$ . (*Hint*: it may help to define  $\tilde{V} := V e^{(R/L)t}$  for one of the voltages.
- Write down the solution for the case where the input voltage is turned on suddenly from 0 to  $V_0$  just after  $t = 0$ , and then held at  $V_0$  for all  $t > 0$ . You can assume the input was zero for all times in the past. What is the time constant of your solution?

### Problem 2.5

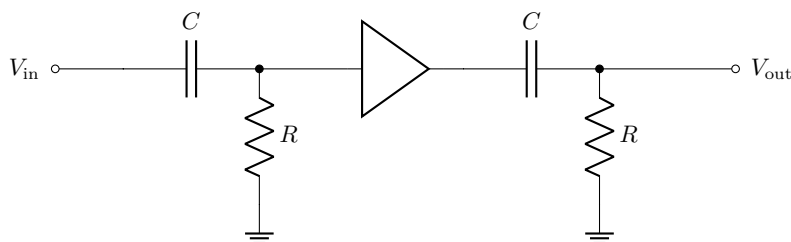
Consider the circuit shown below of two cascaded high-pass filters.



Find the (complex) transfer function  $\tilde{V}_{out}/\tilde{V}_{in}$ , assuming an input frequency of  $\omega$ . How does this scale for large and small frequency? (Also give the scaling at small frequencies in dB/octave.)

### Problem 2.6

Consider the circuit shown below of two cascaded high-pass filters, separated by a **buffer amplifier**.



The buffer amplifier has two important properties: the input (left-hand-side connection) draws no current (and thus produces no load on the first  $RC$  filter), and the output voltage is equal to the voltage at the input.

- Find the (complex) transfer function  $\tilde{V}_{out}/\tilde{V}_{in}$ , assuming an input frequency of  $\omega$ . Examine the scaling behavior for large and small frequencies, and give the scaling at small frequencies in dB/octave.
- Make a (log-log) plot of the amplitude transfer function  $T(\omega) = |\tilde{V}_{out}(\omega)/\tilde{V}_{in}(\omega)|$ . Also include on the same plot the corresponding transfer functions from the Problem 2.5, and the transfer function for a simple high-pass filter.

(c) Make another plot of the phase of the output compared to the input for the same three circuits as in (b). That is, if we write out the amplitude and phase of the transfer function as

$$\tilde{T}(\omega) = T(\omega) e^{-i\phi}, \quad (2.126)$$

then make a plot of  $\phi$  vs.  $\omega$ . Here, use a logarithmic frequency axis, and a linear phase axis. Be clear about the nature of the phase shift (lead vs. lag).

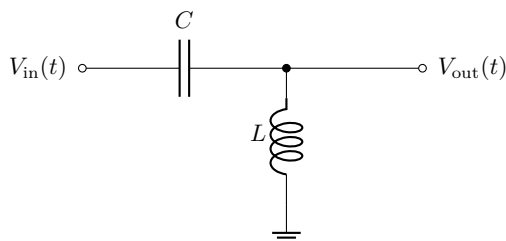
*Note:* use any program you like for plotting. For a plotting tutorial in Mathematica, see

<http://atmoptics.uoregon.edu/~dsteck/teaching/mathematica/plot-tutorial.nb>  
<http://atmoptics.uoregon.edu/~dsteck/teaching/mathematica/plot-tutorial.pdf>

In Mathematica, use `LogLogPlot` for a log-log plot and `LogLinearPlot` for a log plot in the  $x$ -direction instead of `Plot`.

### Problem 2.7

Consider the circuit below, consisting of a capacitor of capacitance  $C$ , and an inductor of inductance  $L$ .



- Is this a high-pass or a low-pass filter? Give a *qualitative* argument for your answer.
- Derive the transfer function  $[\tilde{T}(\omega)]$  for the circuit.
- Derive the amplitude transfer function  $[T(\omega)]$  for the circuit.
- Give the scaling of the transfer function in the “stop band” of the circuit [i.e., the asymptotic region where  $T(\omega)$  is small]. Express your answer in dB/octave.

### Problem 2.8

A “schmapacitor” of “schmapacitance”  $S$  is a bipolar component whose operation is defined by the relation

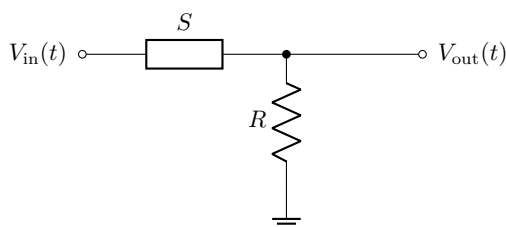
$$I = S \frac{d^3 V}{dt^3}, \quad (2.127)$$

where  $I$  is the schmapacitor current and  $V$  is the voltage drop across the schmapacitor.

Derive expressions for the effective schmapacitance of two schmapacitors in series and for two schmapacitors in parallel.

### Problem 2.9

Consider the circuit below, consisting of a resistor of resistance  $R$ , and a “schmapacitor” of “schmapacitance”  $S$ .



The schmapacitor is defined as in Problem 2.8 by the relation

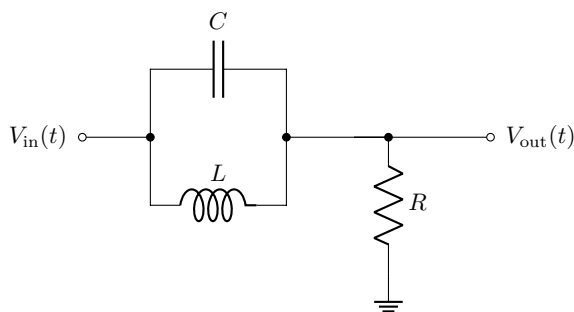
$$I = S \frac{d^3 V}{dt^3}, \quad (2.128)$$

where  $I$  is the schmapacitor current and  $V$  is the voltage drop across the schmapacitor.

- Is this a high-pass or a low-pass filter? Give a *qualitative* argument for your answer.
- Derive the transfer function  $[\tilde{T}(\omega)]$  for the circuit.
- Derive the amplitude transfer function  $[T(\omega)]$  for the circuit.
- Give the scaling of the transfer function in the “stop band” of the circuit [i.e., the asymptotic region where  $T(\omega)$  is small]. Express your answer in dB/octave.
- What is the asymptotic phase in the stop band? Is it a phase lead or a phase lag?
- Derive an expression for the  $-3$ -dB frequency  $f_{3\text{dB}}$  of the circuit.

### Problem 2.10

Consider the RLC circuit below.



- Derive the transfer function  $\tilde{T}(\omega)$  for the circuit.
- Derive the amplitude transfer function  $[T(\omega)]$  for the circuit.
- What is the behavior of  $\tilde{T}(\omega)$  for small and large frequencies? What happens at  $\omega = 1/\sqrt{LC}$ ? Based on your answers, give a *qualitative* description of how this circuit functions as a frequency-dependent filter.

### Problem 2.11

A “schmesonator” is a bipolar component defined by the relation

$$V(t) = S_1 \frac{dI(t)}{dt} + S_2 \int_{t_0}^t I(t') dt', \quad (2.129)$$

where  $S_1$  and  $S_2$  are constant parameters, and  $t_0$  is the time at which the schmesonator circuit is “schmactivated.”

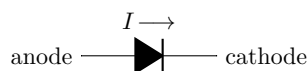
Derive an expression for the schmesonative reactance (i.e., the impedance of the schmesonator assuming voltage and current oscillate at frequency  $\omega$ ).

# Chapter 3

## Diodes

### 3.1 Ideal Diode

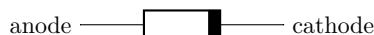
**Diodes** are useful circuit elements: in the simplest sense, they act as one-way valves or “trapdoors” for electrons and thus for electrical current. Schematically, the following symbol represents them:



The two terminals are called the **anode** and **cathode**. The rules that govern the ideal diode are as follows:

1. If  $V_{\text{anode}} > V_{\text{cathode}}$ , then the diode acts like a **short circuit**: lots of current can flow. In this case, the diode is said to be **forward-biased**.
2. If  $V_{\text{anode}} \leq V_{\text{cathode}}$ , then the diode acts like an **open circuit**: no current flows at all. In this case, the diode is said to be **reverse-biased**.

Of course, the situation with real diodes is more complicated, and we’ll get to that. For now, a good mnemonic to remember the direction of current flow is that the diode symbol makes an “arrow” in the direction of the current flow (towards the cathode). Also, in the lab, a typical “signal diode” (for mA-level currents) looks schematically like this:



There is usually a band (not always dark) that marks the cathode end; you can think of the band as being a “minus sign” that marks the cathode.

The diode is obviously a *nonlinear* device, since voltage is not simply proportional to current. Since Ohm’s law  $V = IR$  does not hold in a simple way with constant  $R$ , so a diode is an example of a **non-Ohmic** device.

### 3.2 Vacuum Diodes

The name of the diode comes from the original vacuum-tube realization of a diode, which has two electrodes (hence, the “di”). The cathode is a heated electrode that “boils off” electrons; whether the electrons make it to the anode—and hence whether current flows—is determined by the relative voltage on the anode, because the anode with either repel or attract the electrons from the cathode. (Think about it to see that the current flows only with the correct voltage polarity.)

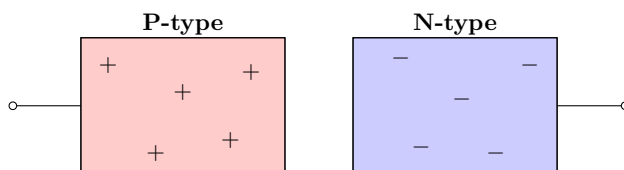
### 3.3 Semiconductor Diodes

By far the most important realization of a diode is using semiconductor materials, so that's what we'll concentrate on here. Roughly speaking, a *semiconductor* is a material somewhere in between a conductor and an insulator. In an insulator, electrons are bound in place, so they can't move around to form a current. In a conductor, electrons move freely, and current flows (there is little resistance). In a semiconductor, the electrons are *mostly* bound, but a few are thermally activated into conducting states, so conduction can happen.

But more important are **doped** semiconductors, where impurities are added to enhance conduction. There are two types:

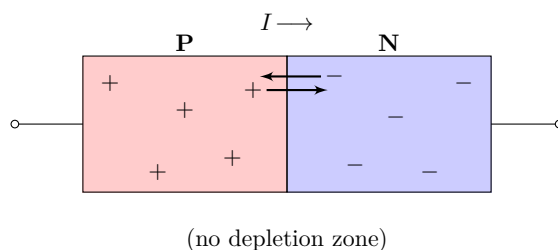
1. In **n-type semiconductors**, the semiconductor is doped with impurities that introduce *excess* electrons (called **n-type carriers**, “n” for for “negatively charged”). The n-type carriers do the conducting.
2. In **p-type semiconductors**, the semiconductor is doped with impurities that introduce a *deficit* of electrons. The *absence* of an electron is something like a positive charge, and is called an **electron hole**, or just **hole**. (These are then called **p-type carriers**, “p” for for “positively charged”). The p-type carriers (holes) do the conducting here.

A **semiconductor diode** is a **junction** between p-type and n-type semiconductors (called a **p-n junction**). If we first consider *separate* p- and n-type semiconductors, we get something like this:

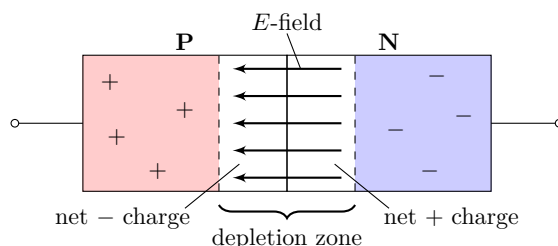


Note that despite the presence of the carriers in each semiconductor, the semiconductors are *electrically neutral*.

Now let's smooch these together and see what happens. There is now a p-n junction, the charges are free to diffuse across due to random, thermal motion.



Remember that the semiconductors were neutral *before* we put them together, so now that carriers are moving across, this sets up net charges on either side of the junction, which creates an electric field. All this is summarized in the diagram below.

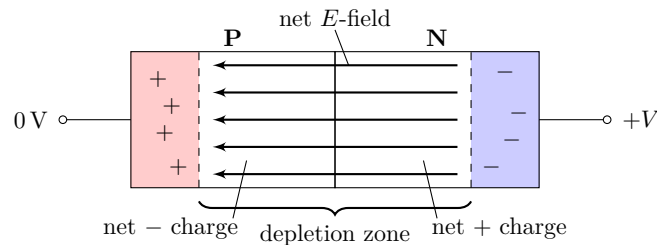


Note that for any carriers that diffused across the junction, the electric field tries to “restore” them back to their original home. So the diffusion continues until the electric field builds up to the point where just balances the tendency for more carriers to diffuse across the junction. Meanwhile, the n-type carriers that make it across to the p-type material will **annihilate** the p-type carriers (electron + hole = nothing), and the same thing happens for p-type carriers that make it across the junction to the n-type material. So there is a region around the junction called the **depletion zone** or **depletion layer**, where there are *no* (net) carriers. Note that in this equilibrium state, no net current can flow, because there are no carriers to transport charge (current) across the depletion zone.

Now let’s argue that this p-n junction realizes the diode as shown below.

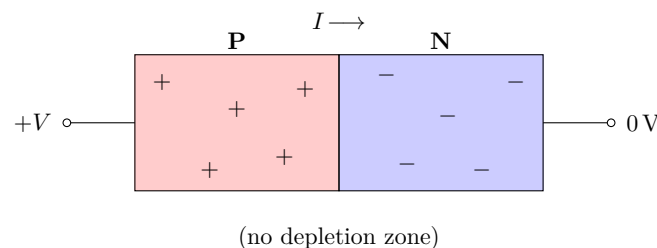


First, let’s do the **reverse-biased case**, where no current should flow. That is, suppose we set the anode to 0 V, and bring the cathode up by  $+V$  relative to the anode. Then the situation is shown in the diagram below.



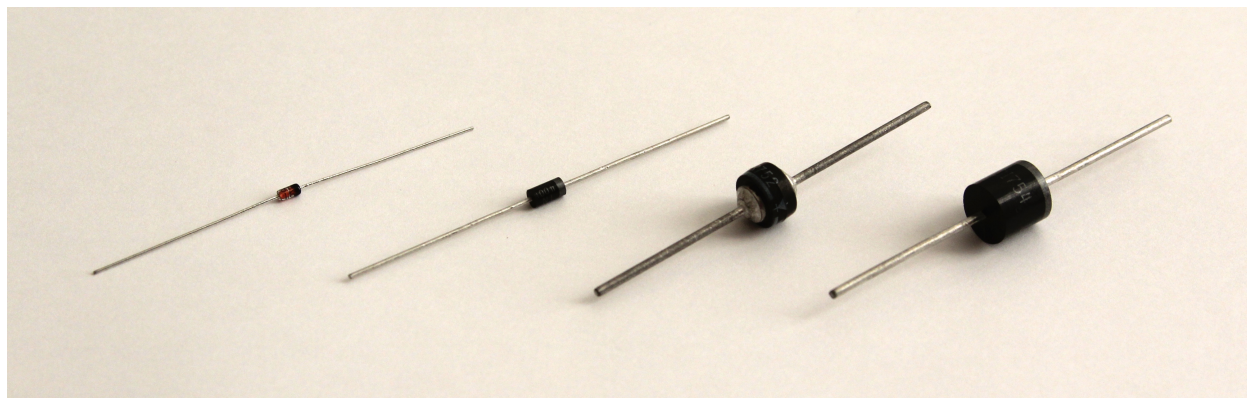
In this case, the electric field due to the applied voltage *adds* with the electric field from the carrier diffusion. This also pulls the p-type carriers towards the anode (the terminal on the p-type material), and n-type carriers towards the cathode. The net effect is just to *increase* the width of the depletion zone.

Now let’s do the **forward-biased case**, where the diode should conduct. So we bring the anode to a voltage  $+V$  relative to the cathode. Now the external field due to this potential difference points from left to right, and *cancels* the electric field due to carrier diffusion, and it tries to make the carriers in each material move towards the junction. The net effect is to *shrink* the depletion zone. Once a *sufficiently strong voltage* is in place, the depletion zone completely disappears, and a net current can flow, as shown below.



Note that the current here involves each type of carrier moving towards the junction, where it annihilates one of the opposite type.

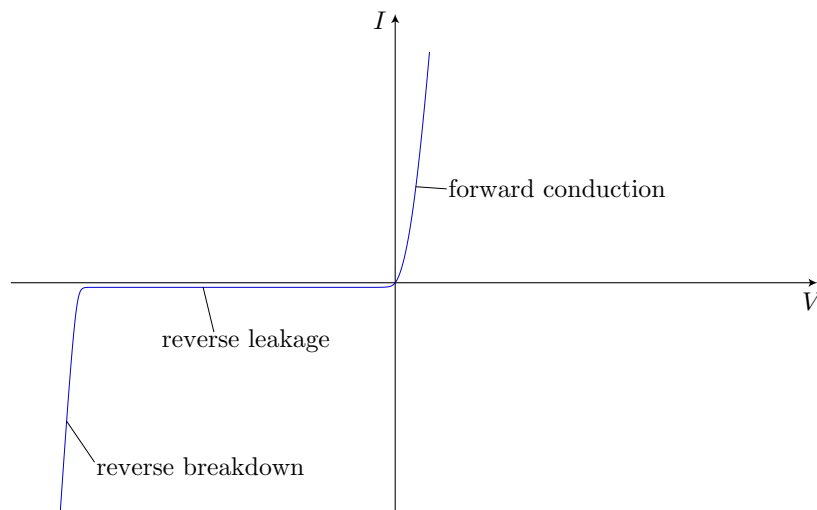
A few typical diodes are shown in the photo below.



From left to right, these are: 1N914B (signal diode), 1N4001 (1-A, general-purpose rectifier), MR752 (6-A, 200-V rectifier), GI754 (6-A, 400-V rectifier). All cathodes point to the upper-right-hand corner.

### 3.4 Current–Voltage Characteristics

All these effects lead to a more complicated relation between voltage and current. Below is a schematic plot of the  $V$ – $I$  relation for a diode.



This is the *forward* current  $I$  plotted against the *forward* voltage  $V$ ; *negative* values of course indicate that the voltage/current are going in the reverse-biased direction. There are a few features to notice here.

1. **Forward conduction.** As the diode is forward-biased, the forward current rapidly (in fact, exponentially) increases. But for any forward current, there is a **forward voltage drop**—remember this is needed to squish the depletion region down to nothing so the diode can conduct. A handy number to remember for quick calculations is 0.7 V for the voltage drop, for forward-biased silicon diodes. The drop is somewhat more for high-power diodes, somewhat less for Schottky and germanium diodes. For example, the common 1N914 (silicon signal diode) has a forward voltage drop of 0.7 V at 10 mA, dropping to 0.6 V at 1 mA and going up to 0.9 V at 100 mA. The 1N5711 Schottky diode (small-signal diode, rated for 15-mA maximum current) has a similar drop of about 0.7 V at 10 mA, but a comparatively lower drop of 0.4 V at 1 mA. These numbers are all temperature-dependent, but the values here are at 25°C, and vary from device to device.
2. **Reverse leakage.** When the diode is reverse-biased, the current is not *completely* blocked, but some current flows. This is roughly constant over a wide range of reverse-bias voltages, and is usually labelled

$I_s$ , with the “S” for **saturation current**. This is also called, more obviously, the **reverse-leakage current**. This is typically  $\sim 10$  nA for silicon and Schottky diodes.

3. **Reverse breakdown.** If the reverse-bias voltage is sufficiently large, the insulating properties of the diode break down (like any insulator), and the diode conducts. The voltage at which the diode starts to conduct is the **reverse-breakdown voltage**, and is at least 100 V for the 1N914 and at least 70 V for the 1N5711.

### 3.4.1 Diode Law

Neglecting the reverse breakdown, the diode  $V$ – $I$  relation is reasonably well-described by the **diode law**,<sup>1</sup>

$$I = I_s \left( e^{V/nV_T} - 1 \right), \quad (3.1)$$

(diode law)

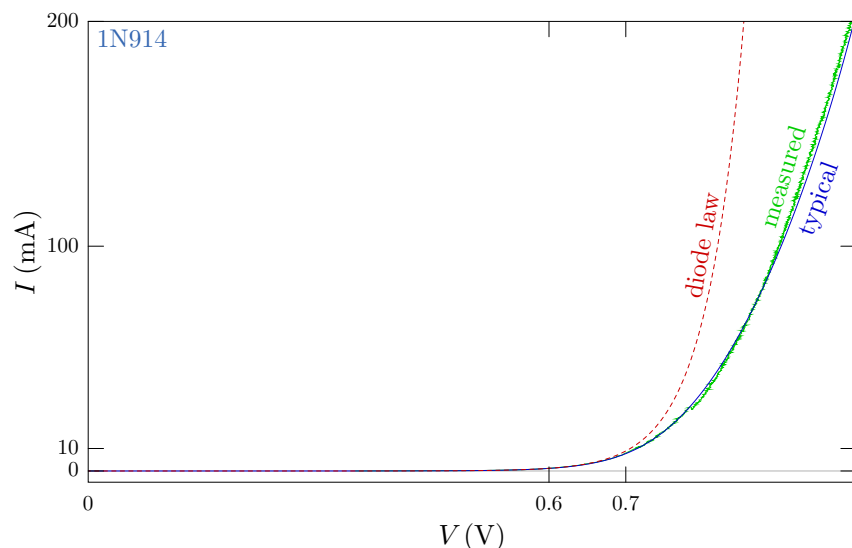
where  $I_s$  is the saturation current that we already discussed,  $V_T$  is the **thermal voltage**

$$V_T := \frac{k_B T}{e}, \quad (3.2)$$

(diode law)

$k_B = 1.381 \times 10^{-23}$  J/K is the **Boltzmann constant**,  $e = 1.602 \times 10^{-19}$  C is the **fundamental charge** (magnitude of the electron charge), and  $T$  is the absolute temperature. The thermal voltage is 25.3 mV at 20°C, or 25.7 mV at 25°C. Also,  $n$  is the **ideality factor**, which typically falls in the range of 1 to 2, and is a sort of “fudge factor” for real junctions. Often this is just set to 1 for simplicity, in which case the diode law becomes the “ideal diode law.”

Note that the diode law is only really valid below the “knee” of the  $V$ – $I$  curve, where the current starts to really take off. As an example, consider the plot below of two models for the 1N914 at 25°C.



The “diode law” (dashed) curve shows the diode law with  $I_s = 6.2229 \times 10^{-9}$  and  $n = 1.9224$ . The “typical” (solid) curve shows a more complete diode model, which models “high-injection” effects (i.e., the modification of diode behavior due to a high concentration of injected carriers) that become significant on the scale of the “knee current” ( $I_{\text{knee}} = 43$  mA here), and a small internal resistance (of about  $0.3 \Omega$  here).<sup>2</sup> The knee

<sup>1</sup>Also called the **Shockley diode law**, after William Shockley, “The Theory of p-n Junctions in Semiconductors and p-n Junction Transistors,” *The Bell System Technical Journal* **28**, 435 (1949) (doi: 10.1002/j.1538-7305.1949.tb03645.x) <https://archive.org/details/bstj28-3-435>, Eq. (3.13).

<sup>2</sup>This is the SPICE (circuit-simulator) model for the 1N914 from Central Semiconductor; the  $I_s$  and  $n$  values in the diode law are the same as the parameters in the SPICE model.



current models “high-injection effects,” and it is typically incorporated via the phenomenological equation

$$I = \frac{I_{\text{diode}}}{\sqrt{1 + I_{\text{diode}}/I_{\text{knee}}}}, \quad (3.3)$$

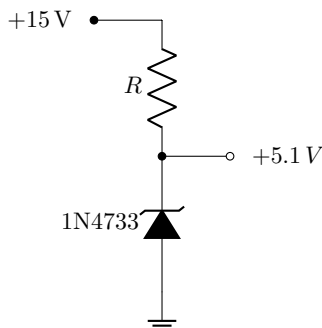
where  $I_{\text{diode}}$  is the exponential (diode-law) approximation. The simple exponential overestimates the conducted current at larger forward voltages, but it is a reasonable approximation up to 0.6 or 0.7 V (up to around 10 mA). The green curve represents data measured from an actual 1N914; note that the knee-current model matches it quite closely.

### 3.5 Zener Diodes

A **Zener diode** (pronounced ZEE-ner) is a regular diode with a carefully engineered reverse-breakdown voltage, typically in the range of 3–100 V. The Zener diode has a special symbol, as shown below.

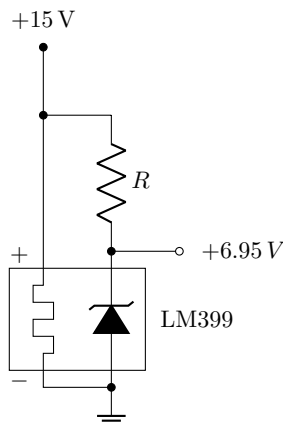


The Zener diode is mainly useful as a voltage regulator. A typical circuit is shown below.



The resistance  $R$  here depends on the intended load. The idea is that the Zener diode “wants” to operate at the reverse-breakdown voltage (5.1 V for the 1N4733), and it draws just the right amount of current to make the voltage drop across  $R$  equal to the difference between the supply (+15 V) and output (+5.1 V) voltages. How does the Zener diode “know” how much current to draw? You should try thinking this through: Think of the reverse-biased diode as a variable resistor, and think of the circuit as a voltage divider. The diode has large resistance for small currents, and small resistance for large currents. The only self-consistent point is for the diode to drop the breakdown voltage, which is the transition between these two regimes.

The main problem with this circuit is that, as we have seen, the properties of diodes depend on temperature. A good solution when you need a precise voltage reference is a **temperature-stabilized Zener diode**, like the LM399, which has a breakdown voltage of 6.95 V, and an integrated oven to keep the Zener’s temperature constant.

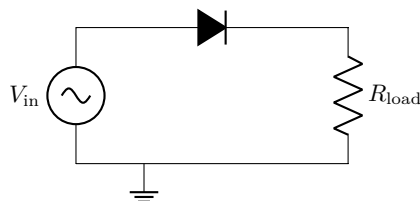


The extra two connections are the power-supply leads for the oven heater. Note that this circuit can draw large (up to 200 mA) for the first few seconds after the circuit is turned on, while the oven temperature stabilizes. You can then obtain other reference voltages using a voltage divider or an op-amp circuit to multiply the voltage by a known factor—something we will get to later.

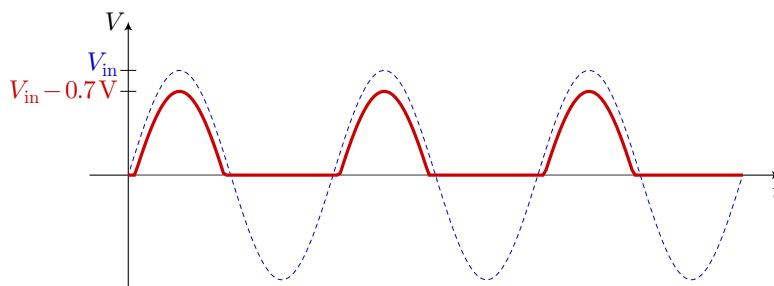
## 3.6 Rectifier Circuits

### 3.6.1 Half-Wave Rectifier

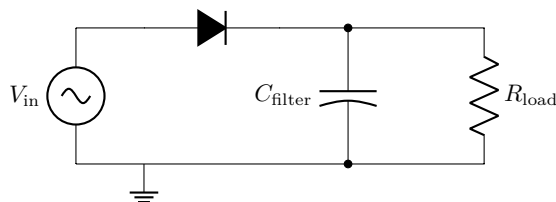
One of the main uses of a diode is as a **rectifier**, or something that changes an alternating-current (ac) signal into a direct-current (dc) signal. This is especially useful for changing the line voltage (120 V in the U.S.) into a dc voltage (e.g., for a power supply), after stepping the line voltage down (or up) by some factor using a transformer. The simplest example of a rectifier circuit is the **half-wave rectifier**, which uses only a single diode.



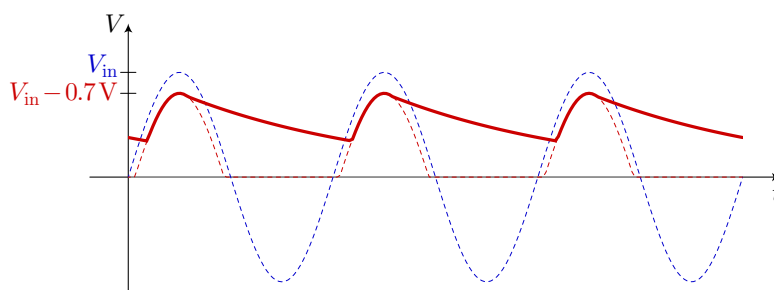
Here, the supply (with *amplitude*  $V_{in}$ ) creates a voltage across the load (represented schematically here by a load resistor  $R_{load}$ ); since current only flows in one direction, the supply can only impose a positive voltage across the resistor, and the voltage during forward conduction is always a diode drop *below* the voltage of the power supply.



This is not very much like a dc voltage, and not a very good dc power supply. The trick to getting a better dc signal is to add a smoothing capacitor across the load, as in the schematic below.



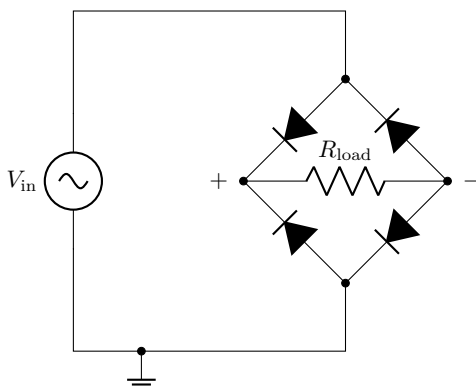
Then the circuit charges the capacitor through the diode to a maximum of  $V_{in}$ , less the diode drop. When the rectified voltage falls away, the capacitor “props up” the output voltage, which decays exponentially with a  $1/e$  time of  $R_{load}C_{filter}$ . The net result is shown below.



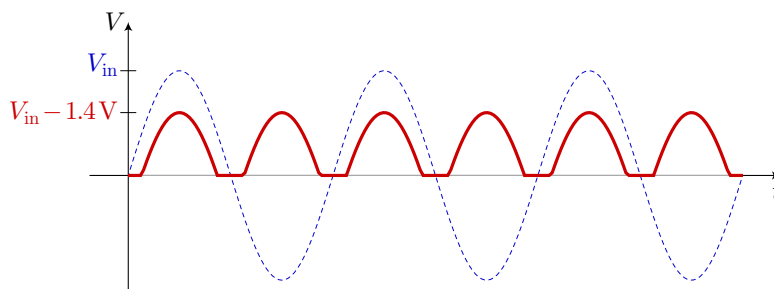
The resulting wave has **ripple**, and choosing the RC time to be as long as possible makes the ripple very small (for good power supplies, the ripple should be of order mV or tens of mV, depending on the current and application). Smaller load resistances (i.e., loads drawing more current) produce more ripple, while larger smoothing capacitors reduce the ripple.

### 3.6.2 Full-Wave Rectifier

A major problem with the half-wave rectifier is the long time between rectified peaks, making it difficult to get small ripple. A solution to this is the **full-wave rectifier**, which uses four diodes as shown below.

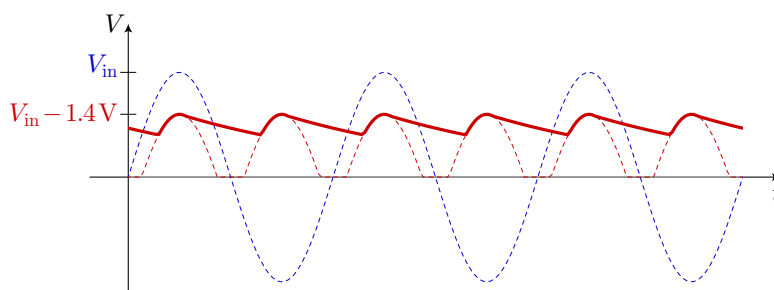


This passes both the positive and negative peaks, as shown below. Trace the current through the diodes on both phases of the input to convince yourself that this works out.



We don't "waste" the negative peaks here, but the price is that we lose *two* diode drops from the input voltage. Also, note that we are plotting the **voltage difference** across the load resistor; the absolute voltages on either end are more complicated, because the output voltage is referenced to the input voltage in a way that depends on which diodes are conducting at the moment. If the input voltage is **floating** (i.e., not referenced to ground, which is usually the case of a transformer output), we can instead ground the  $-$  side of the load, in which case the *absolute* output voltage goes from 0 to  $V_{in} - 1.4\text{ V}$ .

Adding a capacitor to this setup, we get a decent, filtered power supply.



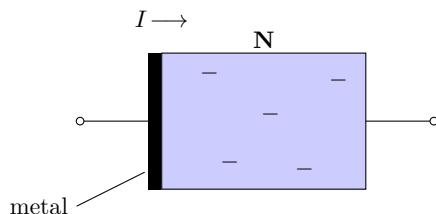
Again, we lose almost another volt compared to the half-wave rectifier, but for the same load and filter capacitor, we get much less ripple. Yet better power supplies (with smaller ripple) can be realized by adding a linear regulator IC.

### 3.7 Schottky Diodes

In some diode circuits, the forward-voltage drop of silicon diodes (of roughly 0.7 V at small currents) can be a disadvantage. (In power-supply rectifiers, for example, the diode voltage drop “wastes” some of the input voltage and thus also wastes power.) We already mentioned that germanium diodes have a lower voltage drop, but nowadays these diodes are fairly uncommon, having been replaced by the much more popular **Schottky diode**. The symbol for the Schottky diode is shown below; there are extra “tips” as compared to the usual diode symbol.



A Schottky diode is formed at the junction of a *metal* and a semiconductor—typically an n-type semiconductor. The basic connection is shown below.



The diode mechanism here is more complicated to explain than in the case of the p-n junction. Without going into much detail, the interaction depends strongly on the material properties of the metal and the semiconductor, specifically on the work function of the metal and how it aligns with the forbidden and allowed energies of the semiconductor. (Remember that there are metal–semiconductor junctions needed just to wire up a semiconductor diode, and these junctions do *not* act as diodes—these junctions are called **ohmic contacts**.) The idea is that the n-type carriers can freely diffuse into the metal. However, the metal modifies the allowed carrier energies in the vicinity of the junction and this creates an effective barrier (called a **Schottky barrier**) for the carriers trying to return from the metal to the semiconductor.

The salient difference between the Schottky and p-n diodes is that there is essentially only one carrier involved in the Schottky diode (the majority n-type carrier)—no depletion layer forms in the same sense as a p-n junction because both carrier types are not present to “annihilate.” However, a thin depletion layer forms at the junction from the diffusion of carriers into the metal, leaving behind a net positive charge on the semiconductor side of the junction (and a net accumulation of electrons on the metal side). A reverse voltage increases this depletion layer, just as in the p-n junction, while a forward voltage overcomes the Schottky barrier, allowing forward current to flow. There are several important resulting characteristics of Schottky diodes compared to p-n diodes:

1. The **forward voltage** is typically smaller than for Si p-n diodes, around 0.3 V for small-signal devices. This forward voltage is tunable by choosing different metals for the junction, but cannot be too small as a very small forward voltage (small Schottky barrier) leads to excessive reverse leakage.

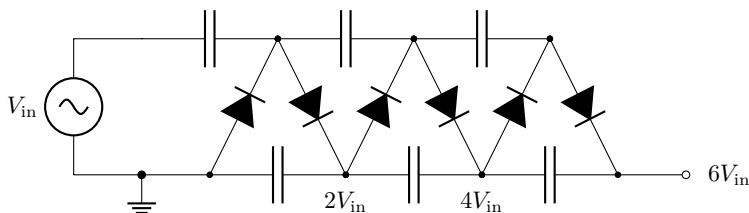
2. **Reverse leakage** is typically worse in Schottky diodes.
3. **Reverse breakdown** also typically occurs at lower voltages in Schottky diodes.
4. A main advantage of Schottky diodes is **fast switching** or **fast reverse recovery**. The idea is that when a p-n diode goes from reverse to forward bias, there is a delay before forward conduction starts, because the minority carriers on each side of the junction (e.g., the n-type carriers in the p-type region) must be returned to their original side so that conduction can begin. Because Schottky diodes only use majority carriers, there is no such issue with “minority carrier storage,” and so their reverse recovery is much faster. For example, the p-n reverse recovery times range from 4 ns (1N914 small-signal diode) to 30  $\mu$ s (1N4011, 1-A rectifier). The 1N5711 small-signal Schottky diode specifies a 1-ns maximum reverse recovery time, and the 1N5819 small-signal Schottky does not even specify a maximum recovery time (just advertising “extremely fast switching”).

Other than the differences noted here, Schottky diodes are functionally equivalent to p-n diodes, and are available in anything ranging from small-signal diodes to high-current power rectifiers.

## 3.8 Circuit Practice

### 3.8.1 Cockcroft–Walton Multiplier

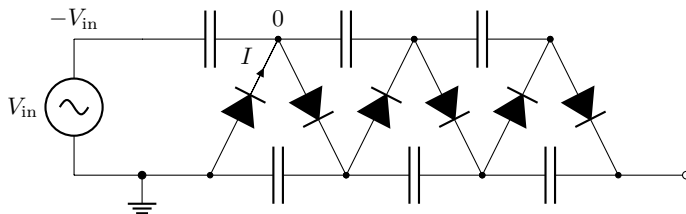
Explain how the circuit below, the **Cockcroft–Walton multiplier**, works. The multiplier is driven by an ac source of amplitude  $V_{in}$ , and each “stage” of the multiplier consists of two capacitors and two diodes. The output after  $N$  stages is (in steady state) a dc voltage of  $2NV$ . Three stages are shown in the example below.



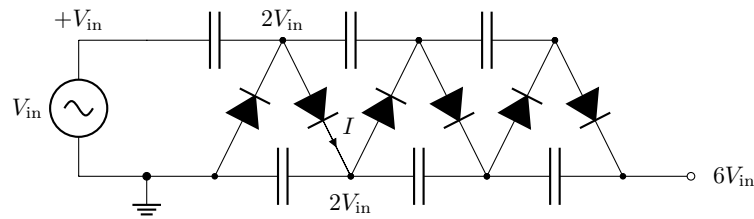
Assume ideal capacitors and diodes (and think about why it is reasonable to ignore the forward voltage drop of the diodes in this kind of circuit).

**Solution.** The idea in multiplying voltages is to get *high* voltages, so the input ac voltage would be of the order of 1 kV, in which case a diode drop of 0.7 V makes little difference.

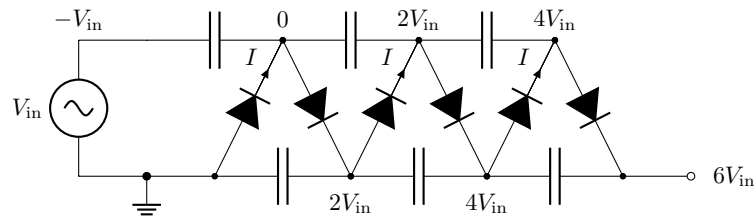
To analyze the circuit, we will first trace the voltages when the input voltage has swung low. Current flows across the first diode to charge the first capacitor.



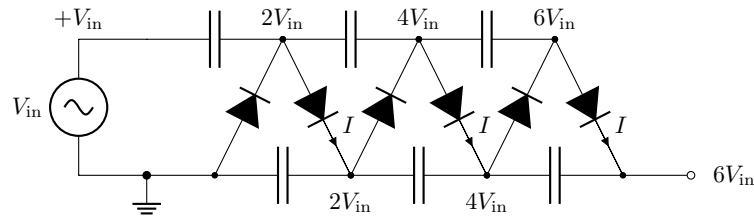
Then, on the positive-voltage phase, the absolute voltage of the first capacitor increases, but the voltage *across* it stays the same. Current flows to charge the second capacitor.



The process continues. In steady state, to balance any leakage of current from the output terminal, current would flow as follows in the negative-input phase,



which switches to the following during the positive-input phase:



### 3.9 Exercises

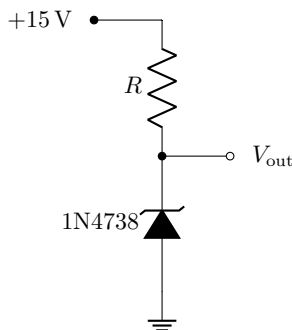
#### Problem 3.1

Every diode has some intrinsic capacitance.

- What is the origin of this capacitance in a p-n junction diode? Be *brief*.
- If the diode is reverse-biased, does the capacitance increase, decrease, or stay the same with increasing reverse voltage? More importantly: Why? (Be *brief*.)

#### Problem 3.2

Consider the zener-diode voltage regulator shown below.

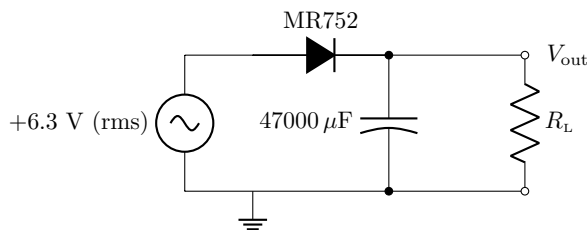


Suppose you design this circuit to drive (supply) a load resistance  $R_L$ .

- What is  $V_{out}$ ?
- If  $R = 1\text{ k}\Omega$ , what is the smallest  $R_L$  that the circuit can handle without “sagging” the output voltage?
- What should be the power rating of the resistor  $R$ ? Be explicit about your assumptions.

#### Problem 3.3

Consider the half-wave rectifier circuit shown below. The ac input voltage to the circuit is a 60 Hz signal from a power transformer.



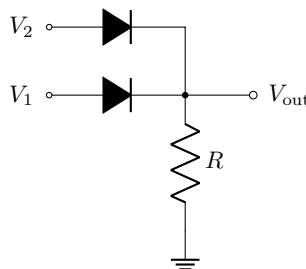
- What is the peak output voltage across the load resistor, assuming  $R_L = 10\text{ }\Omega$ ? Account for the voltage drop across the diode (look it up!).
- Estimate the voltage ripple, assuming the same load resistance.

#### Problem 3.4

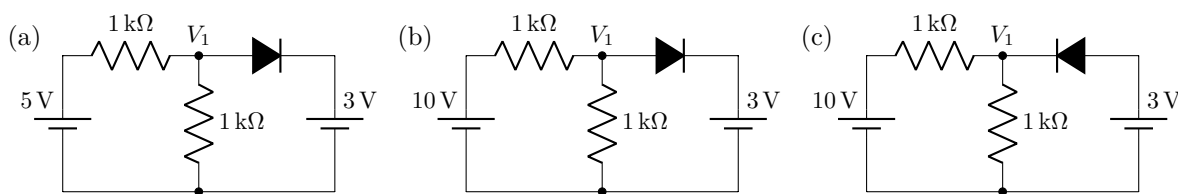
Sketch the analogous full-wave-rectifier-bridge circuit to the circuit in Problem 3.3, and repeat the calculations.

**Problem 3.5**

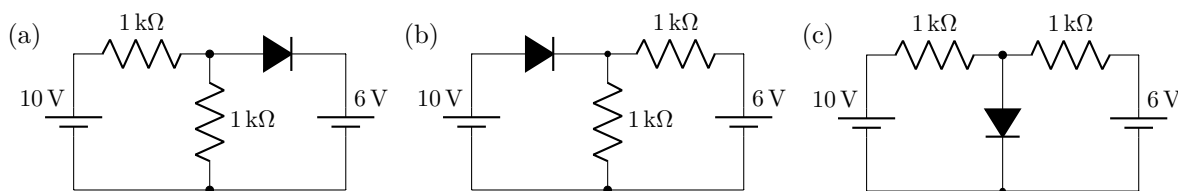
Consider the circuit below. Give the output voltage  $V_{\text{out}}$  in terms of the input voltages  $V_1$  and  $V_2$ . You can assume ideal diodes (no forward voltage, breakdown, or leakage).

**Problem 3.6**

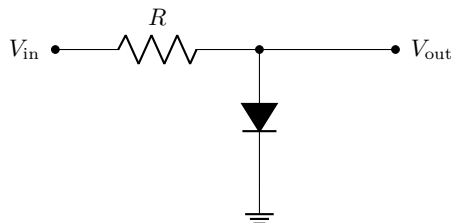
In each circuit, compute  $V_1$  (relative to the negative battery terminals). Account for the forward-voltage drop across the diode as needed.

**Problem 3.7**

In each circuit, compute the current through the diode. Account for the forward-voltage drop across the diode as needed, but ignore any reverse-leakage current. (You can skip any division if you don't have a calculator, just simplify your answer as much as you can.)

**Problem 3.8**

(a) Consider the “voltage-divider” circuit below, consisting of a resistor and a **forward-biased** (signal) diode.



Assuming  $V_{\text{in}} > 0$ , use the diode law to show that

$$V_{\text{in}} + I_s R = V_{\text{out}} + I_s R e^{V_{\text{out}}/nV_T}, \quad (3.4)$$



which determines  $V_{\text{out}}$  in terms of  $V_{\text{in}}$ ,  $R$ , and  $T$ .

(b) Does the idea of a Thévenin-equivalent circuit make sense for this circuit? If yes, give the Thévenin-equivalent voltage and resistance. If no, explain *briefly* why not.

(c) Does  $V_{\text{out}}$  increase or decrease with  $T$ ? (Assume  $R$ ,  $I_{\text{s}}$ , and  $n$  are independent of temperature.)

(d) Obtain an explicit solution to the relation you derived in (a) as follows. Rearrange the equation to combine the parts that depend on  $V_{\text{out}}$  in the form  $x e^x$ , where

$$x = \frac{V_{\text{in}} - V_{\text{out}} + I_{\text{s}}R}{nV_T}. \quad (3.5)$$

That is, you should obtain an equation of the form  $x e^x = y$ , where  $y$  is independent of  $V_{\text{out}}$ . Then use the defining relation

$$W(z) e^{W(z)} = z \quad (3.6)$$

of the **Lambert W function** to show that

$$V_{\text{out}} = V_{\text{in}} + I_{\text{s}}R - (nV_T) W\left(\frac{I_{\text{s}}R}{nV_T} e^{(V_{\text{in}} + I_{\text{s}}R)/nV_T}\right) \quad (3.7)$$

is the explicit solution of the relation (3.4).

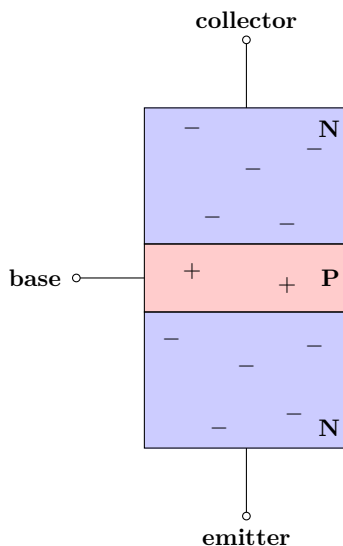
## Chapter 4

# Bipolar Junction Transistors

### 4.1 Overview

A **bipolar junction transistor (BJT)** is our first example of a device that is both *nonlinear* and *active*—active, in the sense that the device should be “powered,” or to say it another way, it uses one signal to modify another signal. At first, it’s a bit counterintuitive to have a device with *three* terminals, but roughly speaking, you can think of it functionally as having a pair of input terminals and a pair of output terminals, but one terminal is “shared” between the input and the output.

BJTs come in two flavors: NPN and PNP, which refers to the stack of doped semiconductors that form the transistor. The schematic construction of the NPN transistor is shown below: the name just gives the order of the layers from top to bottom (or bottom to top).



That is, the NPN transistor is a p-type semiconductor sandwiched in between two n-type semiconductors. However, the important thing is that the inner p-type layer is **thin** and **lightly doped**. The light doping means that there are relatively few p-type carriers in the p-type layer. The PNP transistor is pretty much the same thing, except for interchanging p-type and n-type layers. All the analysis here will go through to that case under this interchange, provided we also reverse all the currents and voltage differences. Thus we’ll stick to the NPN case here (which is the more common case; usually, PNPs are only found when they are paired in some way with an NPN, in part because PNPs are easier to make well).

The schematic symbols for NPN and PNP transistors are shown below.



6. For any device, there are limits to how large  $I_B$ ,  $I_C$ , and  $V_{CE}$  can be, and these limits are different for each species of transistor.
7. Under all the conditions above, the two currents are *proportional*:

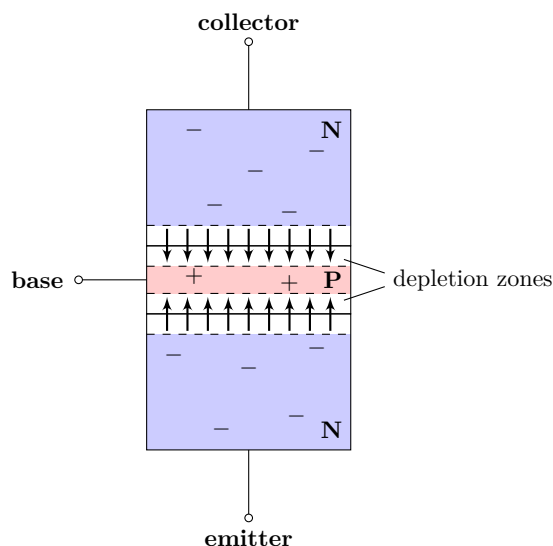
$$I_C = \beta I_B. \quad (4.1)$$

(transistor current-control relation)

The parameter<sup>1</sup>  $\beta$  is roughly constant, and as a simple value for estimating what happens in transistor circuits, you can assume  $\beta \sim 100$ . This parameter varies among different transistor species and even among individuals of one species. It also depends on other circuit parameters (for example,  $\beta$  could vary by almost a factor of 2 in some transistors if  $I_C$  changes by an order of magnitude). On transistor data sheets,  $\beta$  is often denoted by  $h_{FE}$  (or  $h_{fe}$ ).

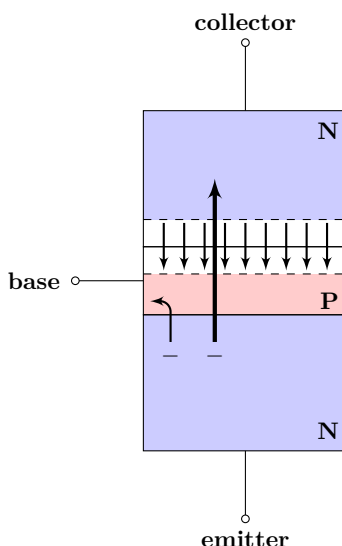
### 4.3 Mechanism

To understand why the transistor works the way it does, let's go back to the diagram of the stack of n- and p-type materials in the NPN transistor. The two p-n junctions set up two depletion zones, with corresponding electric fields as shown below.



The two depletion zones block any C-E current from flowing (in either direction). Then with voltages set up as  $V_{BE} > 0.6\text{ V}$  and  $V_{CE} > 0$  (actually  $V_{CE}$  needs to be at least roughly  $0.2\text{ V}$ ), then the B-E depletion zone disappears, and the C-E depletion zone grows a bit.

<sup>1</sup>The name “ $\beta$ ” for the ratio  $I_B/I_C$  makes a bit more sense by noting that there is another parameter  $\alpha = I_C/I_E$ . These are related by  $\alpha = \beta/(\beta + 1)$ , so  $\alpha$  is close to unity when  $\beta \gg 1$ . Of the two parameters,  $\beta$  is typically the more useful one in circuit analysis.



Then, basically two things happen with the carriers, as in the diagram above.

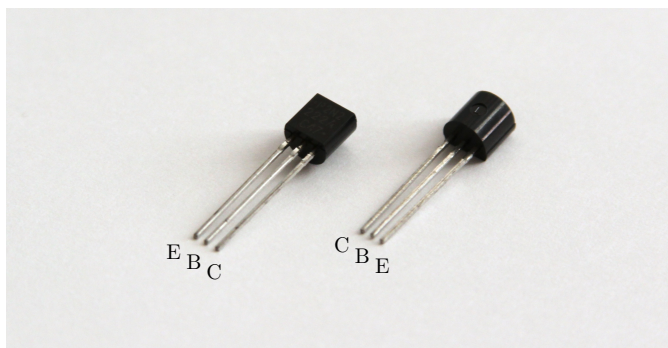
1. Since current is flowing from base to emitter, there are n-type carriers flowing from the *emitter* to the *base* (the negative charge means they are flowing *against* the current). Remember the p-type region is lightly doped, so it is mostly n-type carriers that are transporting the current.
2. Many, or possibly most of, the n-type carriers never make it to the base terminal. What happens is that, once they are pulled into the p-type region by  $V_{BE}$ , they can diffuse into the C–B depletion zone, in which case the electric field in the depletion zone sweeps them into the collector’s n-type region, leading to  $I_C > 0$  (provided  $V_{CE} > 0$ , so the swept-up n-type carriers are removed through the collector terminal).

The transistor is a really beautiful device.

Note that from our description and diagrams, it would seem that the emitter and collector are basically equivalent, and it is true that they are very similar. However, in a real transistor, the geometry is different, and the emitter material is heavily doped compared to the collector material (since it “produces” the carriers needed to transport the current). So while it is *possible* to operate a transistor with emitter and collector interchanged, it would not work nearly as well (the  $\beta$  in this configuration is much smaller, for example).

## 4.4 Packaging

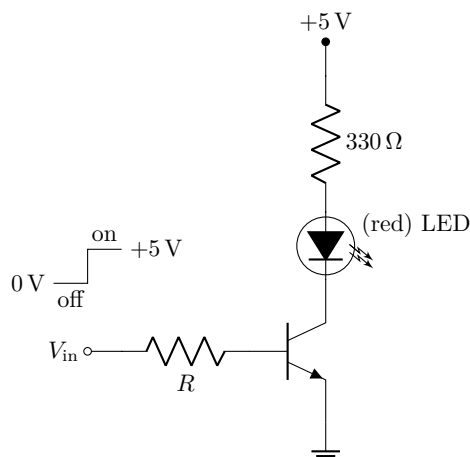
Transistors come in many shapes and sizes. The common TO-92 plastic case (appropriate for low-power signal transistors) is shown in the photograph below, with typical connections.



Note that these connections are common with 2N-series transistors in TO-92 packages, but they are *not* universal.<sup>2</sup>

## 4.5 Transistor Switch

Below is the first example circuit we'll do with transistors: using a transistor as a switch, here to turn a **light-emitting diode (LED)** on or off.



This switch works as advertised (for the stated input voltages) for a resistance  $R$  around  $1\text{ k}\Omega$ . To learn more about transistor operation, we will consider this resistance and a larger value for the base resistance to see two different regimes of transistor operation. This circuit behaves the same for any of a number of small-signal NPN transistors (2N2222A, 2N3904, 2N4401, etc.).

### 4.5.1 Saturation Mode

First, let's consider the case where  $R = 1\text{ k}\Omega$ , which would be a typical way to design this switch. The “off” case is pretty easy: with  $0\text{ V}$  input,  $I_B = 0$ , and thus  $I_C = 0$  from Eq. (4.1), so the LED is off.

Now for the “on” case, with  $5\text{ V}$  input. To start, the B–E junction of the transistor acts like a diode. Usually, the base current will be low, because it will switch a current  $\beta$  times larger, from Eq. (4.1). So we will assume diode drop across the B–E junction, but on the low end, say  $V_{BE} = 0.6\text{ V}$ . Then the  $1\text{-k}\Omega$  resistor drops the rest of the input voltage, or  $4.4\text{ V}$ . This gives a base–emitter current

$$I_B = \frac{4.4\text{ V}}{1\text{ k}\Omega} = 4.4\text{ mA}, \quad (4.2)$$

which will control the collector current  $I_C$ .

Before considering the transistor action, let's consider the LED. LEDs have a higher forward voltage than signal diodes. Typical values for “bright” operation of a “T-1<sup>3</sup>/<sub>4</sub>” sized LED (the most common size, with a plastic bulb  $5\text{ mm}$  in diameter) depend on color:

- For red, orange, yellow, green-yellow: forward current  $I_F = 20\text{ mA}$ ; forward voltage  $V_F = 1.8\text{ V}$ .
- For green, blue, white, UV:  $I_F = 20\text{ mA}$ ;  $V_F = 3.3\text{ V}$ .

“High-brightness” LEDs can have larger forward currents, but these values are good for “normal” LEDs. We have a red LED, so  $V_F = 1.8\text{ V}$ . If we think of the transistor as a switch that connects C–E (this is not quite true, more on this shortly), then the resistor drops  $5 - 1.8 = 3.2\text{ V}$ ; with a  $330\text{-}\Omega$  resistor, this gives  $I_F$  as a bit under  $10\text{ mA}$ , which is reasonable, but hardly pushing the LED's brightness.

<sup>2</sup>In fact, the specific transistors in the photo should be wired *backwards* from this labeling. These are P2N2222A transistors by ON Semiconductor, and for some reason they decided to wire these CBE as viewed from the front, despite all *other* 2N2222A's in this package being wired EBC.

But what does the transistor *want* to do? Eq. (4.1) says that  $I_C$  should be  $\beta$  times  $I_B$ ; assuming  $\beta \sim 100$ , then  $I_C \sim 440 \text{ mA}$ . But, as we found out, with the LED and resistor voltage drops,  $I_C$  can't possibly be this high (the transistor can't have a *negative* voltage drop, for example). In fact, the transistor drop  $V_{CE}$  can't be less than about  $0.2 \text{ V}$ , so the resistor actually drops about  $3.0 \text{ V}$ . With the  $330\text{-}\Omega$  resistor, this gives  $I_C = I_F \approx 9.1 \text{ mA}$ .

This mode of transistor operation is called **saturation mode**: as a switch, the transistor is “wide open,” and the current  $I_C$  is limited by external elements (here, LED and  $330\text{-}\Omega$  resistor), not the transistor. Of course, the current-limiting resistor could be reduced somewhat (say, to  $150\text{-}\Omega$ ) for a brighter LED here. [The term “saturation” comes from the base region being saturated by (minority) carriers from both the collector and emitter regions; note that in saturation the base–collector junction is also slightly forward-biased.)

The important thing to notice in this example is that a high-impedance source ( $1 \text{ k}\Omega$ , from the base resistor, plus any impedance of the input voltage source) controls a higher-current load via the transistor.

### 4.5.2 Forward-Active Mode

Now suppose  $R = 100 \text{ k}\Omega$  in the same circuit. Then the base current is two orders of magnitude smaller, or  $44 \mu\text{A}$ . The transistor tries to set  $I_C \sim 100I_B = 4.4 \text{ mA}$ , and now this is certainly possible for the transistor. Just to double check, the voltage drop across the  $330\text{-}\Omega$  resistor is  $330\text{-}\Omega \times 4.4 \text{ mA} = 1.5 \text{ V}$ . Then removing the resistor and LED drops,  $V_{CE} = 5 \text{ V} - 1.8 \text{ V} - 1.5 \text{ V} = 1.7 \text{ V}$ . Now it is the *transistor* that is regulating the current via its voltage drop. This is called the **forward-active mode** of the transistor, and it is in this regime that the transistor can act as an amplifier with some gain (the saturation mode is a kind of “infinite-gain” amplification).

### 4.5.3 Summary

Let's just summarize the difference between saturation and forward-active modes in the switching circuit, because this can be a bit confusing the first time through, and it's important to understand the difference *intuitively*.

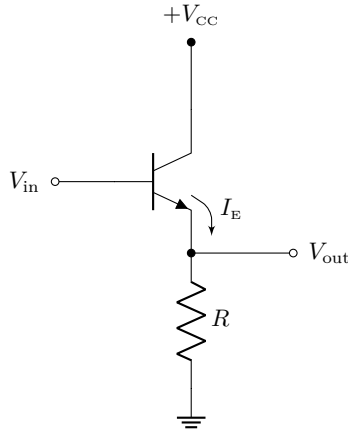
In **forward-active mode**, the base current  $I_B$  controlled the collector current  $I_C$  (and thus the LED current) via the relation  $I_C = \beta I_B$ . The transistor “regulates”  $I_C$  by adjusting the voltage  $V_{CE}$  to the proper amount. We know the current is set by the voltage drop across the  $330\text{-}\Omega$  resistor: the larger the resistor drop, the larger the current. This means that  $V_{CE}$  is *smaller* for larger currents, because more of the supply voltage must be taken up by the resistor.

In **saturation mode**, the transistor relation *doesn't* hold. That's because the base current makes the transistor “want” more current than the  $330\text{-}\Omega$  resistor will allow. The transistor normally tries to drive more current by reducing  $V_{CE}$ , but it can't do so below zero (or more precisely, below about  $0.2 \text{ V}$  if  $I_C > 0$ ), so  $I_C < \beta I_B$ .

So you can really think of the transistor relation as being more like  $I_C \leq \beta I_B$ . Both the transistor and an external element (resistor) will want to limit the current; whichever wants *less* current is the one that wins.

## 4.6 Emitter Follower

As our second transistor circuit, we'll continue with a transistor in forward-active mode. Below is the **emitter follower**, where the output voltage is intended to “follow,” or match, the input voltage. It is an amplifier in the sense of being an amplifier for *current*.



Note that we are introducing some new notation here. In addition to the base and collector currents in the diagram on p. 92, we are introducing the emitter current  $I_E$ . Also, the power-supply voltage is labelled  $+V_{CC}$ ; the plus denotes a positive voltage with respect to ground and the “C” subscript denotes this is intended to power the collector terminal (two C’s distinguishes this from the voltage  $V_C$  at the collector, which is the same in this circuit, but not always).

We will consider the signals to be *ac* signals, with some dc bias that we don’t care much about. This is kind of a fact of life with transistors, which can only work with current flowing in one direction—if we want to handle a bipolar signal (like an audio signal), these biases ensure that the net signal is compatible with the way the transistor works. The bias voltage/current is typically rejected at some point near the output of the circuit, for example using a high-pass filter.

Thus, let’s set

$$V(t) = V_0 + v(t), \quad I(t) = I_0 + i(t), \quad (\text{dc and ac components}) \quad (4.3)$$

where  $V_0$  and  $I_0$  are the dc biases, and the (small) ac signals are  $v(t)$  and  $i(t)$ . This simplifies the analysis a bit, because for example we can write the output voltages as

$$V_{\text{out}} = V_E = V_B - 0.6 \text{ V}, \quad (4.4)$$

but dropping the dc offsets, this becomes

$$v_{\text{out}} = v_E = v_B. \quad (4.5)$$

The **ac voltage gain** for the circuit is defined by

$$G := \frac{v_{\text{out}}}{v_{\text{in}}} = \frac{v_E}{v_B} = 1, \quad (\text{ac voltage gain}) \quad (4.6)$$

so this is indeed a unity-gain circuit (voltage follower).

To compute the current gain, we can relate the emitter current to the base and collector currents by

$$I_E = I_B + I_C = I_B + \beta I_B = I_B(\beta + 1), \quad (4.7)$$

where we used Eq. (4.1). In terms of ac components, this is

$$i_E = i_B(\beta + 1), \quad (4.8)$$

so the **ac current gain** is  $(\beta + 1)$ .

### 4.6.1 Input and Output Impedance

Since the emitter follower has current gain, it can allow a source with high output impedance to drive a lower-impedance load. To quantify this, let’s define  $Z_{\text{load}}$  to be the impedance at the transistor emitter,



which is  $R$  in parallel with any other load impedance attached to the  $V_{\text{out}}$  terminal. Then we can calculate

$$i_B = \frac{i_E}{\beta + 1} = \frac{v_E}{Z_{\text{load}}(\beta + 1)} = \frac{v_B}{Z_{\text{load}}(\beta + 1)}. \quad (4.9)$$

Here, we used the current gain (4.8), then  $v_E = i_E R$ , then  $v_E = v_B$ . Then we can define the **input impedance** of the amplifier by

$$Z_{\text{in}} := \frac{v_{\text{in}}}{i_{\text{in}}}. \quad (\text{input impedance, definition}) \quad (4.10)$$

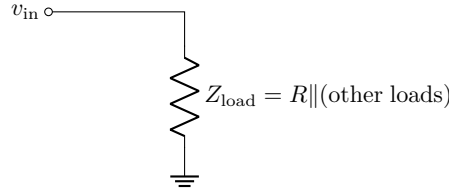
Writing the input voltage and current in terms of the base voltage and current,

$$Z_{\text{in}} = \frac{v_B}{i_B}, \quad (4.11)$$

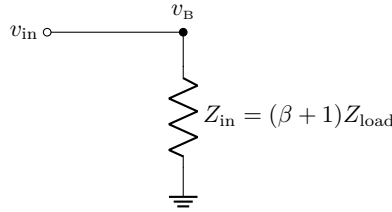
and then using Eq. (4.9),

$$Z_{\text{in}} = Z_{\text{load}}(\beta + 1). \quad (\text{input impedance, emitter follower}) \quad (4.12)$$

We can interpret this as follows: *Without* the transistor, the input  $v_{\text{in}}$  would have to “drive” the load impedance (the resistor  $R$  and other connected stuff) directly, as shown below



The effect of the transistor is to multiply the resistor value by  $(\beta + 1)$ , so effectively the impedance is about 100 times larger, and thus much easier to drive—the transistor is supplying most of the current via the collector, so the voltage source driving the resistor doesn’t have to supply much current at all. Thus, the equivalent circuit *with* the transistor is shown below.



Of course, we can think of this (Thévenin) equivalent circuit because the transistor is acting as a linear amplifier.

We can also define an **output impedance** as

$$Z_{\text{out}} := \frac{v_{\text{out}}}{i_{\text{out}}}. \quad (\text{output impedance, definition}) \quad (4.13)$$

Manipulating this a bit,

$$Z_{\text{out}} = \frac{v_{\text{in}}}{i_E} = \frac{v_{\text{in}}}{i_B(\beta + 1)}, \quad (4.14)$$

and then defining the source impedance

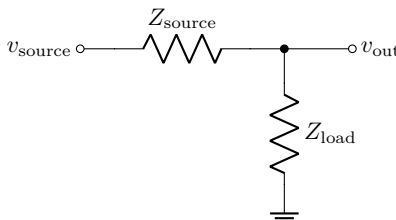
$$Z_{\text{source}} := \frac{v_{\text{in}}}{i_B} \quad (\text{impedance of input source}) \quad (4.15)$$

(i.e., the impedance of whatever is driving the transistor base at voltage  $V_{in}$ ), the output impedance becomes

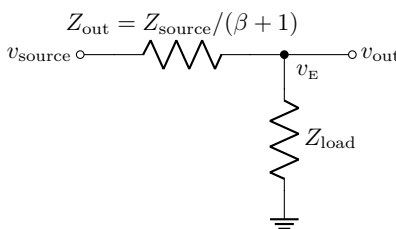
$$Z_{out} = \frac{Z_{source}}{\beta + 1}. \quad (4.16)$$

(output impedance, emitter follower)

We can interpret this as follows: *Without* the transistor in place, the load would be connected directly to the source, loading it down somewhat:



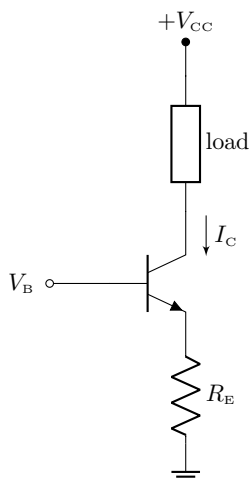
That is,  $v_{source}$  and  $Z_{source}$  represent the Thévenin-equivalent voltage and impedance of the signal source. The transistor effectively divides the source impedance by  $(\beta + 1)$ , making the source effectively “stiffer” (closer to an ideal voltage source) by about 100-fold, as “seen” by the load. Equivalently, the transistor is supplying most of the load current, so the source only has to supply a small fraction of the current that it otherwise would.



Again, from this point of view, the transistor is “lumped into” the source circuit, which is appropriate because it is acting as a linear device.

## 4.7 Transistor Current Source

The next circuit we will consider is the **transistor current source**. This is more of a prelude to the transistor amplifier that is coming up next, but this is also useful in its own right. A **current** source is a circuit that maintains a constant (programmed) current, independent of voltage, within limits of course. In the circuit below, the goal is to maintain a constant current through the load, independent of the load impedance.



In this case, the transistor is maintaining the collector current  $I_C$ . To see that this works, note that the base and emitter voltages are related as usual by

$$V_E = V_B - 0.6 \text{ V}. \quad (4.17)$$

Then the emitter resistor  $R_E$  sets the emitter current via

$$I_E = \frac{V_E}{R_E} = \frac{V_B - 0.6 \text{ V}}{R_E}. \quad (4.18)$$

Remember that the base and collector currents add to form the emitter current. Since  $I_C = \beta I_B$ , then

$$I_E = I_C + I_B = I_C(1 + \beta^{-1}) \approx I_C, \quad (4.19)$$

provided  $\beta$  is large. Thus,

$$I_C \approx \frac{V_B - 0.6 \text{ V}}{R_E}. \quad (4.20)$$

(transistor current source)

That is, the load current  $I_C$  is programmed by the input base voltage  $V_B$ , as well as the emitter resistor  $R_E$ . Importantly, the question of the load resistance never came into the analysis, so the current is independent of the load resistance.

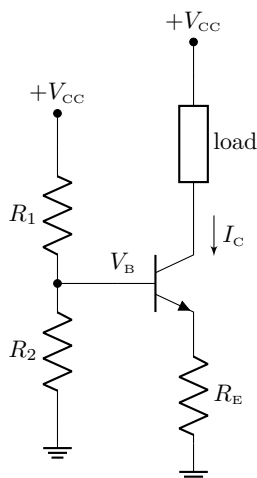
### 4.7.1 Compliance

Well, that is, within limits. What are the limits? That is, what is the **compliance** of the current source? Specifically, the compliance refers to the range of output voltages for which the transistor is properly regulating the current—the output voltage here meaning the collector voltage  $V_C$ , which the transistor “presents” to the load.

For proper transistor operation, we need  $V_B > V_E$  to switch the collector current, and we also need the collector  $V_C \gtrsim V_E + 0.2 \text{ V}$ . On the upper end (i.e., lower-current end), we can have  $V_C$  going all the way up to  $+V_{CC}$ , assuming this doesn't exceed the breakdown voltage of the transistor. So, for a given input voltage  $V_B$ , the range of output voltages  $V_C$  is from  $V_E + 0.2 \text{ V} = (V_B - 0.6 \text{ V}) + 0.2 \text{ V} = V_B - 0.4 \text{ V}$ , on up to  $+V_{CC}$ . This corresponds to a range of current from 0 on up to  $(V_{CC} - V_B + 0.4 \text{ V})/R_{\text{load}}$ , in terms of the load resistance.

### 4.7.2 Bias Network

Of course, to make the current source *work*, we need to set the base voltage  $V_B$ . How do we do this if we only have ground and the power-supply voltage? The answer, of course, is a voltage divider, as in the circuit below.



The voltage divider here could even be a potentiometer (variable resistor), so we can fine-tune the regulated current (our treatment is only approximate, for example, because we are using the approximate  $V_{BE}$  drop of 0.6 V). The voltage divider gives an *unloaded* voltage of

$$V_B \approx \frac{R_2}{R_1 + R_2} V_{CC}, \quad (4.21)$$

but we are “loading” the divider with the transistor, so we have to be careful. Remember that the Thévenin equivalent circuit for the voltage divider (Section 1.4.1) is as follows:

$$V_{Th} = \frac{R_2}{R_1 + R_2} V_{CC} \quad \bullet \quad \begin{array}{c} R_{Th} = R_1 \parallel R_2 \\ \text{---} \text{---} \text{---} \end{array} \rightarrow V_B$$

In our discussion of the emitter follower, we saw that the input impedance of the transistor is about  $\beta R_E$ . That is, as “seen” from the base-side of the transistor, the equivalent input circuit is

$$V_{Th} = \frac{R_2}{R_1 + R_2} V_{CC} \quad \bullet \quad \begin{array}{c} R_{Th} = R_1 \parallel R_2 \\ \text{---} \text{---} \text{---} \end{array} \begin{array}{c} V_B \\ \text{---} \\ \beta R_E \\ \text{---} \\ \text{---} \end{array}$$

Here, we have the input voltage divider, still connected to the transistor base, but with the transistor replaced by its effective input impedance. Thus,  $V_B$  is determined by another voltage divider, but  $V_B$  is essentially the unloaded value above (i.e.,  $V_{Th}$ ) under the condition

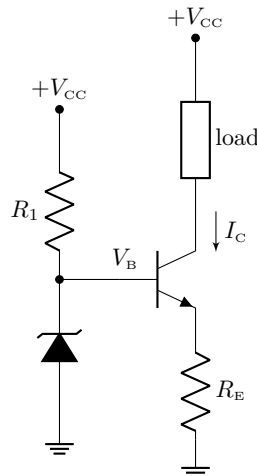
$$R_{Th} = R_1 \parallel R_2 \ll \beta R_E. \quad (4.22)$$

(design condition for bias network)

This condition ensures that the voltage divider is “stiff”—that is, it acts like an ideal voltage source. If this condition is not fulfilled, the divider’s voltage “sags” under the load of  $R_E$  via the transistor (i.e., the input impedance of the transistor).

Note that in designing this circuit, you might think that we can just *calculate* the effect of  $R_E$  on  $V_B$ , and so we can just tweak the ratio of  $R_1$  to  $R_2$  to compensate. However, the resulting voltage depends on  $\beta$ , which can vary from device to device, or with temperature, or with collector current, etc. That is, in this circuit design (and other circuit designs), it’s important to be in the regime where we can neglect the effect of  $R_E$  on  $V_B$ —that is, in the regime where  $\beta$  is large, but the *particular value* of  $\beta$  is not critical.

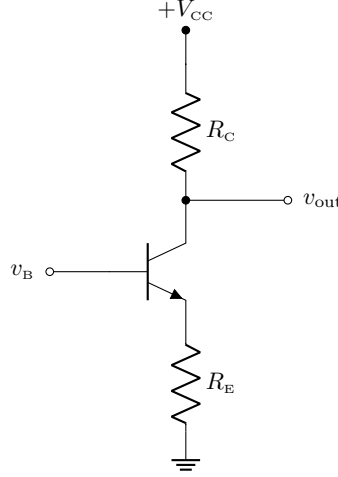
One other solution to the problem of a sagging voltage divider is to use the Zener-diode voltage regulator (Section 3.5) to set  $V_B$ .



In this case, we can get a fairly coarse choice of  $V_B$  by choosing the Zener diode, and fine-tune the current by setting  $R_E$ .

## 4.8 Common-Emitter Amplifier

And now for our first voltage amplifier that has **gain**—the output can be larger than the input. The basic circuit is simple: just take a transistor current source, and use a collector resistance  $R_C$  as the load. Heuristically, the input voltage programs a collector current  $I_C$ , and  $R_C$  acts to convert the current back into a voltage, which serves as the output.



This amplifier only works as advertised on ac signals, so let's consider small ac variations  $v(t)$  and  $i(t)$  with respect to dc biases  $V_0$  and  $I_0$ , as before:

$$V(t) = V_0 + v(t), \quad I(t) = I_0 + i(t). \quad (4.23)$$

The current-source result (4.20), dropping dc biases, becomes

$$i_C = \frac{v_B}{R_E}. \quad (4.24)$$

The voltage drop across  $R_C$  is

$$V_{CC} - V_C = I_C R_C, \quad (4.25)$$

which, in terms of ac quantities, is

$$v_C = -i_C R_C. \quad (4.26)$$

Combining this with Eq. (4.24),

$$v_C = -\frac{R_C}{R_E} v_B. \quad (4.27)$$

Identifying these voltages with the input and output voltages,

$$v_{out} = -\frac{R_C}{R_E} v_{in}. \quad (4.28)$$

(common-emitter amplifier)

Defining the **gain** of the amplifier by

$$G := \frac{v_{out}}{v_{in}}, \quad (4.29)$$

(ac amplifier gain)

the common-emitter gain is

$$G = -\frac{R_C}{R_E}. \quad (4.30)$$

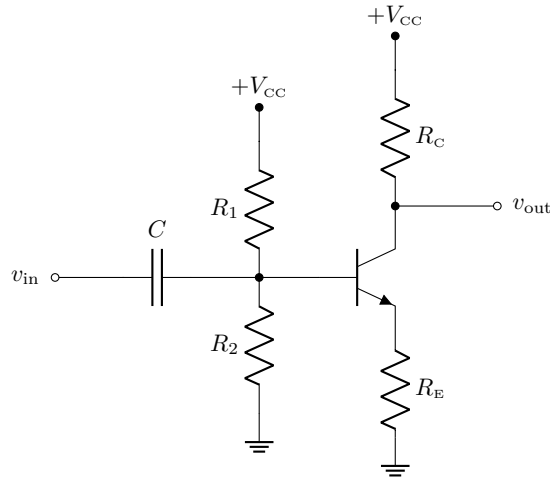
(common-emitter gain)

Then we see that the gain is *negative*, meaning the input ac signal is inverted at the output, and the ratio of collector and emitter resistors controls the gain.

For the sample-circuit numbers, we have  $G = -10$ .

## 4.9 Bias Network (AC Coupling)

As in the current source, we need an input network to set the dc bias. A nonzero bias is critical for the amplification of the ac signal: the B–E junction can only conduct in the forward direction, and so the input signal can't cross through zero and still be amplified without a lot of distortion.



Here, the voltage divider sets the bias voltage, and the input capacitor only passes the ac part of the input signal.

Let's go through the different parts of the circuit, and work out all the relevant parameters. As a concrete example, we will use the parameters:

$$\begin{aligned}
 V_{CC} &= 15 \text{ V}, \\
 R_1 &= 56 \text{ k}\Omega, \\
 R_2 &= 5.6 \text{ k}\Omega, \\
 C &= 0.1 \mu\text{F}, \\
 R_E &= 330 \Omega, \\
 R_C &= 3.3 \text{ k}\Omega.
 \end{aligned} \tag{4.31}$$

1. **AC input impedance.** The input circuit is a high-pass filter, with capacitance  $C$ . The Thévenin resistance is the parallel resistance of  $R_1$  and  $R_2$ , but we must also include the input impedance  $\beta R_E$  of the amplifier as an additional parallel resistance. Thus, at high frequencies when the capacitor acts as a short-circuit, the input impedance of the circuit is the Thévenin resistance, or  $R_1 \parallel R_2 \parallel \beta R_E$ .

For this circuit,  $R_1 \parallel R_2 = 5.1 \text{ k}\Omega$ . Also,  $\beta R_E \approx 100 \text{ k}\Omega$ , which is much larger than  $R_1 \parallel R_2$ , so  $R_1 \parallel R_2 \parallel \beta R_E \approx R_1 \parallel R_2 = 5.1 \text{ k}\Omega$ .

2. **High-pass input.** The corner frequency of the input network is

$$f_{3\text{dB}} = \frac{1}{2\pi(R_1 \parallel R_2 \parallel \beta R_E)C}. \tag{4.32}$$

(input corner frequency)

For the circuit here,  $f_{3\text{dB}} = 310 \text{ Hz}$ .

3. **Loading condition.** It is good practice for the input impedance of the transistor to have negligible effect. That is, we should have

$$\beta R_E \gg R_1 \parallel R_2, \quad (4.33)$$

in which case the input impedance and corner frequency of the input network do not depend on  $\beta$  (and thus on temperature, etc.).

For the sample numbers here, we have already seen that we satisfy this condition.

4. **Gain.** From Eq. (4.30), the gain of the amplifier is

$$G = -\frac{R_C}{R_E}. \quad (4.34)$$

5. **Output impedance.** Because the transistor acts like a current source (with  $I$  insensitive to  $V$ ), the transistor effectively presents a very large impedance at the output—a large change in voltage causes *little* change in the current. Thus the impedance at the output is the Thévenin resistance at that point, or the parallel impedance of the collector with  $R_C$ . Thus, we can take the output impedance of the amplifier to be approximately  $R_C$ .

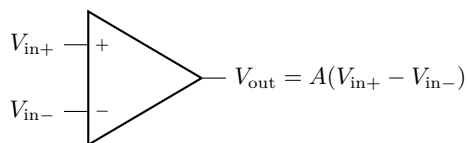
In this example, the output impedance of the amplifier is 3.3 k $\Omega$ .

6. **Bias points.** The bias at the input is just the voltage-divider voltage, including any correction from  $\beta R_E$ , which we usually want to ignore. In the circuit here, the 15 V supply is divided down to 1.36 V.

Why does this make sense? It's important to note that the bias at the input isn't critical in the sense of needing to be near the middle of the supply range. If we are doing a lot of amplification and want to avoid **clipping**, though, the output should be close to the center of the supply. Let's see if this is the case. The emitter voltage is  $V_E = V_B - 0.6 \text{ V} = 0.76 \text{ V}$ . This programs a collector current  $I_C \approx (V_B - 0.6 \text{ V})/R_E = 2.3 \text{ mA}$ , from Eq. (4.20). Then the voltage drop across  $R_C$  is 7.6 V which is about half of  $V_{CC}$ . Notice how the matching ratios of  $R_1$  to  $R_2$  and  $R_C$  to  $R_E$  lead to a sensible bias voltage here, but because the (fixed) voltage  $V_{BE}$  enters here, this rule does not always apply.

## 4.10 Transistor Differential Amplifier

Now we come to a more sophisticated transistor amplifier, the **transistor differential amplifier**. Schematically, a differential amplifier has the following form:

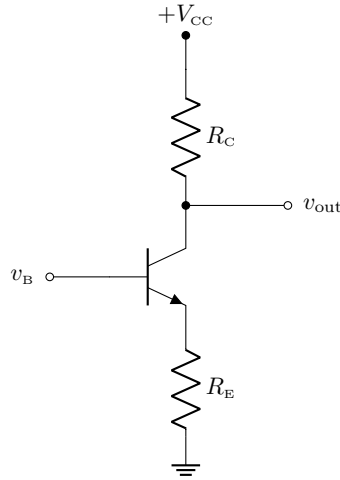


That is, the output is the difference of the inputs (the signs at the inputs tell you which is subtracted from which), and multiplied by the voltage **gain factor**  $A$ .

Why is a differential amplifier useful? For example:

1. In the transmission of signals, differential amplifiers give you **noise immunity**— if you transmit a signal on two lines (signal and ground), the same noise appears on each line, and gets cancelled out by a differential amplifier on the receiving end.
2. Differential amplifiers offer a straightforward way to implement **negative feedback**, which is a nice way to achieve close-to-ideal behavior in nearly every respect. We will go over this in much more detail when we get to op-amps.

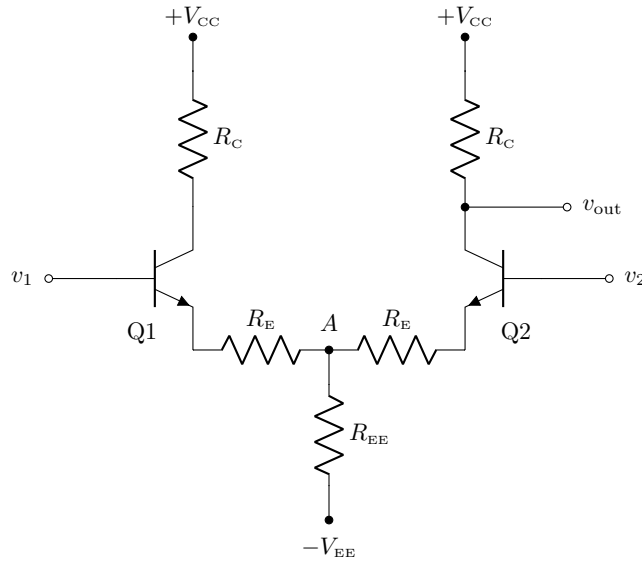
The simplest realization of a differential amplifier uses two transistors. To understand this, first recall the common-emitter amplifier.



In this amplifier, the result for the ac signal was

$$v_{\text{out}} = -\frac{R_C}{R_E} v_{\text{in}}. \quad (4.35)$$

The differential amplifier is basically a “stack” of two common-emitter amps, sharing a common resistance  $R_{EE}$  at the negative end.



From the layout of the schematic, you can see why this circuit is also called a **long-tailed pair** (the “long” here refers to the magnitude of  $R_{EE}$ , which as we will see is typically large compared to  $R_E$ ). Note also that the emitter-end of the circuit is powered by a negative supply, so we can have positive or negative outputs.

To analyze this circuit, we will again concentrate on the ac components of the inputs,  $v_1$  and  $v_2$ . Now let’s think of  $v_1$  and  $v_2$  as being deviations from the mean voltage  $\bar{v}$ :

$$v_1 = \bar{v} + \frac{\Delta v}{2}, \quad v_2 = \bar{v} - \frac{\Delta v}{2}. \quad (4.36)$$

Here,  $\bar{v}$  is the **common-mode signal**,

$$\bar{v} := \frac{v_1 + v_2}{2}, \quad (4.37)$$

and  $\delta v$  is the **differential signal**,

$$\Delta v := v_1 - v_2. \quad (4.38)$$



Ideally, a differential amplifier responds *only* to the differential signal, and is insensitive to the common-mode signal.

#### 4.10.1 Differential-Only Input

We will show that the circuit responds linearly to the inputs, so we can treat the differential and common-mode signals *separately*, just add them together to handle an *arbitrary* input signal. So first, let's concentrate on only the differential signal. That is, setting  $\bar{v} = 0$ , we have

$$v_1 = \frac{\Delta v}{2}, \quad v_2 = -\frac{\Delta v}{2}. \quad (4.39)$$

Then, just as in the emitter follower,

$$v_{E1} = v_{B1} = v_1, \quad (4.40)$$

and similarly

$$v_{E2} = v_2. \quad (4.41)$$

At point  $A$  in the circuit, ignoring dc offsets, we have a 50% voltage divider between  $v_{E1}$  and  $v_{E2}$ , so the voltage is the average of these:

$$v_A = \frac{v_{E1} + v_{E2}}{2} = \frac{v_1 + v_2}{2} = 0. \quad (4.42)$$

We can interpret this as follows: The point  $A$  acts like the “ground” point for the two common-emitter amplifiers in this circuit. Since this point is stable with respect to differential inputs, the common-emitter results carry through here, and in particular for the right-hand common-emitter amp,

$$v_{\text{out}} = -\frac{R_C}{R_E} v_2, \quad (4.43)$$

or in terms of the differential signal,

$$v_{\text{out}} = \frac{R_C}{2R_E} \Delta v. \quad (4.44)$$

(differential response)

Note that the minus sign disappeared here, which is why it is sensible to take the output from the right-hand transistor. Thus, we have

$$G_{\text{diff}} := -\frac{R_C}{2R_E}. \quad (4.45)$$

(differential gain factor)

as the **differential gain factor** for the amplifier.

#### 4.10.2 Common-Mode-Only Input

Now we can focus on *just* the common-mode signal. That is, we take  $\Delta v = 0$ , so that

$$v_1 = v_2 = \bar{v}. \quad (4.46)$$

Then applying Kirchoff's law to the currents at point  $A$ ,

$$i_{EE} = i_{E1} + i_{E2} = 2i_E, \quad (4.47)$$

since the two emitter currents are the same. Here,  $i_{E1}$  and  $i_{E2}$  are the (ac components of the) currents out each emitter (both equal to  $i_E$ ), and  $i_{EE}$  is the current through  $R_{EE}$ . Then applying Ohm's law at point  $A$ ,

$$v_A = i_{EE} R_{EE} = 2i_E R_{EE}, \quad (4.48)$$

and now Ohm's law across either emitter resistor gives

$$i_E = \frac{v_E - v_A}{R_E} = \frac{\bar{v} - 2i_E R_{EE}}{R_E}, \quad (4.49)$$

after eliminating  $v_A$ . Solving for  $i_E$ ,

$$i_E = \frac{\bar{v}}{R_E + 2R_{EE}}. \quad (4.50)$$

Then the output is, just as we had in the common-emitter amplifier,

$$v_{\text{out}} = -i_C R_C \approx -i_E R_C, \quad (4.51)$$

if we assume  $\beta$  large, and then using Eq. (4.50),

$$v_{\text{out}} = - \left( \frac{R_C}{R_E + 2R_{EE}} \right) \bar{v}, \quad (4.52)$$

(common-mode response)

such that we can define

$$G_{\text{CM}} := - \frac{R_C}{R_E + 2R_{EE}} \quad (4.53)$$

(common-mode gain factor)

as the **common-mode gain factor**.

### 4.10.3 General Input and Common-Mode Rejection

Since the circuit responds linearly to the inputs, we can take a general pair of inputs  $v_1$  and  $v_2$ , decompose them into differential and common-mode components via Eqs. (4.37) and (4.38), and then the output is

$$v_{\text{out}} = G_{\text{diff}} \Delta v + G_{\text{CM}} \bar{v}. \quad (4.54)$$

(common-mode gain factor)

Again, for a “good” differential amplifier,  $G_{\text{CM}}$  should be zero, or at least small compared to  $G_{\text{diff}}$ .

One way to quantify the “goodness” of the differential amplifier is the **common-mode rejection ratio (CMRR)**, which we define as

$$\text{CMRR} := \left| \frac{G_{\text{diff}}}{G_{\text{CM}}} \right| = \frac{R_E + 2R_{EE}}{2R_E}. \quad (4.55)$$

(common-mode rejection ratio)

Thus, typical differential-amplifier designs will be such that  $R_{EE} \gg R_E$ , in which case the CMRR reduces to the simple ratio

$$\text{CMRR} \approx \frac{R_{EE}}{R_E}. \quad (4.56)$$

Again, an ideal differential amplifier has a large CMRR. Typically the CMRR is large enough that it is usually measured in dB.

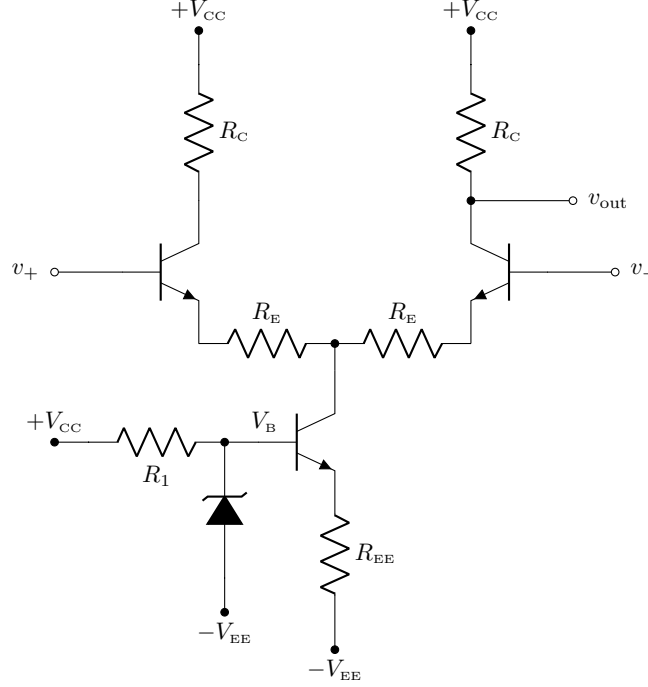
To give an intuitive recap, a *differential* input signal produces an output signal by changing the currents  $i_{E1}$  and  $i_{E2}$  through the emitter resistors  $R_E$ . But these currents cancel at point  $A$ , the current through  $R_{EE}$  stays fixed, and the two transistors act like common-emitter amplifiers. To maximize the differential gain, we want  $R_E$  to be much smaller than  $R_C$ . A *common-mode* signal, however, *does* change  $i_{EE}$ , because the common input change affects  $v_A$ . This effect is mitigated by making  $R_{EE}$  large, so that the change in current  $i_E$  due to a common-mode voltage input is small. The change in the collector current  $i_C$  and thus output voltage is correspondingly small for a common-mode input, making the amplifier relatively insensitive to common-mode signals.

### 4.10.4 Improving the Differential Amplifier

One problem with the differential amplifier is that it is usually desirable to have a large differential gain  $G_{\text{diff}}$ , which means  $R_C$  should be large. However,  $R_C$  is also the output impedance, and we usually *don't* want a large output impedance (we want the output to act more like an ideal voltage source). A simple solution to

this is to buffer the output using an emitter follower, which will reduce the effective output impedance by a factor of  $\beta$ .

One trick to improving the CMRR comes from the observation from Eq. (4.55) that is that  $R_{EE}$  should be large for a high CMRR. The idea is to replace  $R_{EE}$  by a *current source*, forcing  $i_{EE}$  to be constant in the common-mode analysis above. This could be, for example, the transistor current source from Section 4.7, as shown in the diagram below.



Intuitively, a current source is like a very large  $R_{EE}$ , because it yields a small (essentially zero) change in  $i_{EE}$  due to a change in  $v_A$ . In fact one way to view the current source is that it regulates  $v_A$  precisely to maintain a constant  $i_{EE}$ . In doing so, it prevents any change in collector currents due to a common-mode input, and thus the circuit is insensitive to  $\bar{v}$ . In the circuit above, note that the Zener diode sets the input voltage to the current source, which sets the current in combination with  $R_{EE}$  (the dc current through  $R_{EE}$  should be set to give the desired dc voltage bias, typically 0 V, at the output).

Finally, note that we can obtain a good CMRR by having  $R_E \ll R_{EE}$ . In fact, why can't we just set  $R_E = 0$ ? We *can* replace the resistors  $R_E$  by shorts, but still these resistances will not be zero. In this case, the resistances will be the **intrinsic resistance**  $r_e$  of the transistors, to which we will return below.

## 4.11 Ebers–Moll Model

So far, we have been understanding transistor circuits using the crude **current-control model**, centered on the equation [Eq. (4.1)]

$$I_C = \beta I_B, \quad (4.57)$$

with  $\beta$  roughly constant. The crudeness is, of course, that  $\beta$  is *not* constant. Usually we try to design in such a way that this doesn't matter, but it is still useful to develop a better model of the transistor.

Recall the diode law (3.1), relating the diode current and voltage:

$$I = I_s \left( e^{V/nV_T} - 1 \right). \quad (4.58)$$

Here,  $I_s$  is the saturation current (reverse leakage current),  $n$  is the ideality,  $V_T$  is the **thermal voltage**

$$V_T := \frac{k_B T}{e}, \quad (4.59)$$

$k_B = 1.381 \times 10^{-23}$  J/K is the Boltzmann constant, and  $e = 1.602 \times 10^{-19}$  C is the fundamental charge. The base–emitter junction of the transistor works much like a diode (recall that the collector current arises as a “side effect” of the base–emitter current, since carriers headed initially towards the base are swept into the collector), and so a similar law serves as a simple model of transistor operation, the **Ebers–Moll equation**.<sup>3</sup>

$$I_C = I_S \left( e^{V_{BE}/nV_T} - 1 \right). \quad (4.60)$$

(Ebers–Moll equation)

This relates the collector current  $I_C$  to the base–emitter voltage  $V_{BE} = V_B - V_E$ , via a *transistor* saturation current  $I_S$ . Unlike the current-control view of the  $\beta$  model of the transistor, the Ebers–Moll equation gives a **voltage-control view** of transistor operation. The control of  $I_C$  via  $I_B$  is only *indirect*, since  $I_B$  is controlled by  $V_{BE}$  via diode-like conduction. Thus, the transistor is a **transconductance device** (meaning a device that converts voltage to current). The Ebers–Moll model turns out to be accurate over a wide range of currents, typically from nA–mA. Note that in what follows, we will drop the ideality factor  $n$  to simplify notation somewhat without significantly affecting the calculations.

### 4.11.1 Magnitudes

To compare with what we know before, for a typical circuit analysis we assumed  $V_{BE} \approx 0.6$  V, which is much bigger than the thermal voltage  $V_T$  (25.69 mV at 25°C). Thus the exponential in the Ebers–Moll equation is large compared to the 1, and so Eq. (4.60) becomes (setting  $n \approx 1$  to simplify the discussion)

$$I_C \approx I_S e^{V_{BE}/V_T} \quad (4.61)$$

and

$$I_C \gg I_S. \quad (4.62)$$

(To give a sense of scale for this latter statement,  $I_S \approx 1.75$  pA for the common 2N4401, so this condition is certainly fulfilled for typical transistor circuits when the B–E junction is forward-biased.) This also means that, as in the diode, that  $V_{BE}$  is relatively insensitive to variations in the collector current, if we think about the relation in a “backwards” way. That is, if  $I_C$  changes by a factor of 10, and  $V_{BE}$  changes by  $\Delta V_{BE}$ , then

$$10 \approx e^{\Delta V_{BE}/V_T}, \quad (4.63)$$

so that the voltage change is

$$\Delta V_{BE} \approx V_T \log 10, \quad (4.64)$$

which is about 59 mV at 25°.

### 4.11.2 Relation to $\beta$

Then how to we recover the current-control relation (4.57)? If we think of the base–emitter junction as a diode, then we simply apply the diode law (4.58) to obtain

$$I_B = I_{SE} \left( e^{V_{BE}/nV_T} - 1 \right), \quad (4.65)$$

where  $I_{SE}$  is the saturation current of the base–emitter junction. Then solving this for  $V_{BE}$ ,

$$V_{BE} = nV_T \log \left( \frac{I_B}{I_{SE}} + 1 \right), \quad (4.66)$$

and putting this into the Ebers–Moll equation (assuming the same ideality in both relations), we find

$$I_C = \left( \frac{I_S}{I_{SE}} \right) I_B. \quad (4.67)$$

<sup>3</sup>J. J. Ebers and J. L. Moll, “Large-Signal Behavior of Junction Transistors,” *Proceedings of the IRE* **42**, 1761 (1954) (doi: 10.1109/JRPROC.1954.274797).

This has the form of the transistor  $\beta$  relation (4.57), with a  $\beta$  of  $I_S/I_{SE}$ . Note, however, that  $\beta$  is temperature-dependent and somewhat dependent on current, which is not reflected in this simple derivation (for example, we assumed that the same ideality  $n$  applies in both the B–E diode law and the Ebers–Moll equation, but even a slight difference in ideality leads to imperfect cancellation of the exponential factors and thus some dependence on  $V_{BE}$  and  $T$ ).

### 4.11.3 Intrinsic Emitter Resistance

The Ebers–Moll equation is also useful in establishing an **intrinsic emitter resistance** of the transistor. For small ac signals, the intrinsic resistance in the emitter is defined by

$$v_{BE} = i_E r_e. \quad (4.68)$$

More generally, we can regard the resistance to be defined by the  $I$ – $V$  slope

$$\frac{1}{r_e} := \frac{dI_E}{dV_{BE}}, \quad (\text{intrinsic emitter resistance: definition}) \quad (4.69)$$

thinking of the current as a response to the voltage. That is, an ac emitter current  $i_E$  modulates the base-emitter voltage  $v_{BE}$  by an amount controlled by  $r_e$  (the “little r” emphasizes that this resistance applies to ac voltages or small voltage changes, not to large dc voltages). Given that these ac signals are small, this relation is given by the derivative of the  $v_{BE}$ – $i_E$  relation, or the Ebers–Moll equation. We can calculate this via

$$\frac{1}{r_e} = \frac{i_E}{v_{BE}} \approx \frac{i_C}{v_{BE}}, \quad (4.70)$$

since we assume  $\beta$  is large, as usual. Then at some bias collector current  $I_C$ ,

$$\frac{1}{r_e} \approx \frac{dI_C}{dV_{BE}}. \quad (4.71)$$

Differentiating Eq. (4.61), this becomes

$$\frac{1}{r_e} \approx \frac{I_C}{nV_T}, \quad (4.72)$$

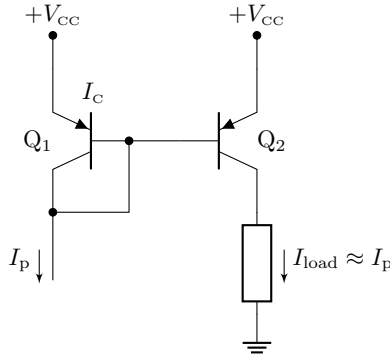
or

$$r_e \approx \frac{nV_T}{I_C} \approx \frac{25.3 \text{ mV}}{I_C} \text{ at } 20^\circ\text{C} \quad (\text{intrinsic emitter resistance}) \quad (4.73)$$

if  $n \approx 1$ . This gives a useful expression for the intrinsic resistance of the transistor, which depends on both current and temperature.

### 4.11.4 Current Mirror

A good example of a simple transistor circuit that can best be understood via the Ebers–Moll equation is the **current mirror**. Essentially,  $I_p$  in the circuit below is the “program current,” and the mirror “copies” the current through the load, independent of the load impedance (within limits, of course). For example, the program current can be set with a resistor connecting the programming terminal (collector of  $Q_1$ ) to ground. This circuit is shown using PNP transistors, but can also work with NPN transistors if the voltages are reversed.



How does this work?  $I_C \approx I_P$  for  $Q_1$  if we assume a large value of  $\beta$ . Then this sets  $V_{BE}$  via the Ebers–Moll equation, and thus  $V_B$  for  $Q_1$ . The transistor bases are connected, so this also sets  $V_B$  for  $Q_2$ , and thus  $V_{BE}$  for  $Q_2$ , since  $Q_2$ 's emitter is connected to the same supply as  $Q_1$ 's. Then, since  $Q_2$  satisfies the same Ebers–Moll equation as  $Q_1$ ,  $Q_2$ 's collector current must be the same as  $Q_1$ 's. Note that we didn't need the exact *form* of the Ebers–Moll equation, just that it relates  $I_C$  to  $V_{BE}$ , and that it is the *same* for both transistors. This is *only* the case if the two transistors are identical and at the same temperature; if these conditions are not true, the analysis is more complicated. In practice, to make sure the properties and temperature match, the current mirror could be implemented using a matched transistor pair in a monolithic package.

Note that we ignored the base current in the above treatment. The base–collector connection on  $Q_1$  shunts the base current for *both* transistors through  $I_P$ . Thus, a better expression for the load current is

$$\left(1 + \frac{2}{\beta}\right) I_{\text{load}} = I_P, \quad (4.74)$$

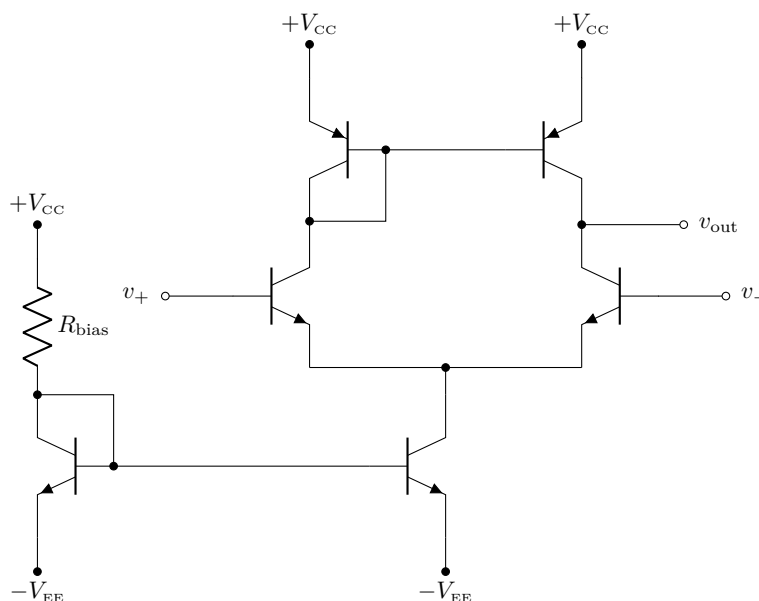
since the program current is the load current plus two base currents. This correction leads to a 2% or less discrepancy between the program and load currents. (There are also some nice tricks for better balancing this basic current mirror.)

This is a circuit that is used, for example, as a common building block in integrated circuits. For example, the (now obsolete) OPA622 op-amp uses an external resistor connected to internal current mirrors to set the **quiescent current** (current when the circuit is idling)—this allows the user to set the trade off between power efficiency and high speed.<sup>4</sup>

#### 4.11.4.1 Application to the Differential Amplifier

The current mirror is also useful in one more refinement of the transistor differential amplifier, continuing the discussion from Section 4.10.4. In order to build an amplifier with very high differential gain, from Eq. (4.45) we should have  $R_C \gg R_E$ . One strategy that we discussed before is to set the external  $R_E = 0$ , which replaces  $R_E$  by the intrinsic emitter resistance  $r_e$ . However,  $R_C$  should be as large as possible. The scheme that we discussed for making  $R_{EE}$  effectively large was to replace the resistor by a current source, making the effective resistance very high. A similar strategy works for the collector resistors  $R_C$ , which can be replaced by a PNP current mirror as in the circuit shown below (where in seeking a high gain we have also omitted the emitter resistors  $R_E$ ). One other difference compared to the previous differential amplifier in Section 4.10.4 is that we have replaced  $R_{EE}$  by one half of an NPN current mirror (instead of a transistor current source). The resistor  $R_{\text{bias}}$  sets the total dc current for the PNP current mirror, and thus it sets the output bias level.

<sup>4</sup>See [http://pdf1.alldatasheet.com/datasheet-pdf/view/137153/BURR-BROWN/OPA622/+037275VvyRSCEuybTFVZKhE.+/  
datasheet.pdf](http://pdf1.alldatasheet.com/datasheet-pdf/view/137153/BURR-BROWN/OPA622/+037275VvyRSCEuybTFVZKhE.+/). This is a good exercise: find the current mirrors in the schematic diagram, Fig. 2 (p. 10).



The current mirror is rather flexible in terms of the *common-mode* currents through the signal-input transistors, but it fights *differences* in these currents (i.e., the currents  $i_e$  in the previous analysis of Section 4.10.2). Thus, the current mirror gives a very large effective value of  $R_C$ , but only for *differential* currents—the *common-mode* current is determined by the current-source transistor at the bottom of the diagram, as before. This circuit is a simple version of what you might find in an **operational amplifier (op-amp)**, which we will study in much more detail in Chapter 7. Generally speaking, it’s much better to use this kind of high-gain differential amplifier in the form of a pre-packaged op-amp, where the details of the circuit have already been carefully engineered and characterized. For example, this circuit has high output impedance (because the output only “sees” the collectors of two transistors), and in normal use the output would be buffered so it can drive a practical load.

### 4.11.5 Other Refinements to the Transistor Model

The BJT models that we have considered so far (the “ $\beta$  model” and the Ebers–Moll model) are fairly simple,<sup>5</sup> although they are quite useful for circuit-design purposes. To close out our discussion of BJTs, we will note some other complications that are useful to keep in mind.

#### 4.11.5.1 Temperature Dependence of the Base–Emitter Voltage

First, remember that due to the diode-like nature of the BJT, the “input voltage”  $V_{BE}$  depends on temperature. We can get idea of the strength of this dependence by getting the temperature slope from the diode law. For example, differentiating the diode law in the form (4.66) gives

$$\frac{dV_{BE}}{dT} = \frac{V_T}{T} \log \left( \frac{I_B}{I_S} + 1 \right) = \frac{V_{BE}}{T}. \quad (4.75)$$

Assuming  $V_{BE} \approx 0.6$  V, this is about  $2 \text{ mV}/^\circ\text{C}$  at  $25^\circ\text{C}$ . However, this is wrong! It turns out that  $I_S$  *increases* exponentially with temperature, which tends to counteract the temperature dependence in  $V_T$ . The net effect is that

$$\frac{dV_{BE}}{dT} \approx -2.1 \text{ mV}/^\circ\text{C}, \quad (4.76)$$

<sup>5</sup>For comparison, see for example the much more detailed model of J. M. Early, “Design Theory of Junction Transistors,” *Bell Labs Technical Journal* **32**, 1271 (1953) (archive: bstj32-6-1271) (doi: 10.1002/j.1538-7305.1953.tb01462.x).

or different from the naïve calculation by about a minus sign, and the proportionality to  $T^{-1}$  still approximately holds.<sup>6</sup>

#### 4.11.5.2 Early Effect

The **Early effect** says that  $V_{BE}$  also depends on  $V_{CE}$  at fixed  $I_C$ , and conversely  $I_C$  depends on  $V_{CE}$  at fixed  $V_{BE}$ .<sup>7</sup> The origin of this effect is apparent from the model diagram of the BJT on p. 91. Since the base–collector junction is reverse-biased, an increasing  $V_{CE}$  (assuming  $V_{CE} > 0$ ) leads to a larger depletion zone, thus effectively reducing the size of the base region (recall that this region is *thin*, increasing its sensitivity to variations in the depletion zone). Roughly speaking, this makes it easier for carriers entering the base (from the emitter) to be swept into the collector by the depletion-zone field.

A simple model that accounts for this effect is a modification of the Ebers–Moll model<sup>8</sup>

$$I_C = I_S \left( e^{V_{BE}/nV_T} - 1 \right) \left( 1 + \frac{V_{CE}}{V_A} \right), \quad (\text{modified Ebers–Moll model}) \quad (4.77)$$

where  $V_A$  is the **Early voltage**. This is a lowest-order model in  $V_{CE}$  (i.e.,  $V_{CE}$  should not be too large compared to  $V_A$ , and the model doesn't quite get the transistor behavior correct at the threshold of saturation. Note that the base current should not be affected by  $V_{CE}$  in the same way; thus the relation  $I_C = \beta I_B$  should be modified to read  $I_C = \beta(1 + V_{CE}/V_A)I_B$  to cancel the effect of the extra factor in the Early-modified Ebers–Moll equation (4.77) on  $I_B$ .

For the relatively common 2N4401, for example,  $V_A \approx 360$  V, and so it is certainly true that  $V_{CE} \ll V_A$  in typical circuits. Thus, the dependence on  $V_{CE}$  is relatively weak. At fixed  $V_{BE}$ , a change  $\Delta V_{CE}$  leads to a fractional change of  $I_C$  of

$$\frac{\Delta I_C}{I_C} = \frac{\Delta V_{CE}}{V_A + V_{CE}}, \quad (4.78)$$

according to the model (4.77). For  $V_{CE} = 15$  V, for example, this is a fractional change in  $I_C$  of about 0.3% for a 1-V change in  $V_{CE}$ . At fixed  $I_C$ , it works out that an increase in  $V_{CE}$  *decreases*  $V_{BE}$  slightly. From Eq. (4.77), we find

$$\Delta V_{BE} \approx -\frac{nV_T}{V_A + V_{CE}} \Delta V_{CE} \sim -10^{-4} \Delta V_{CE} \quad (4.79)$$

by expanding to lowest order in  $\Delta V_{BE}$ . Again, the magnitude of the effect here is small.

Within the model (4.77) for the Early effect, we can also estimate the effective collector resistance by rearranging Eq. (4.78) to read

$$r_c = \frac{\Delta V_{CE}}{\Delta I_C} = \frac{V_A + V_{CE}}{I_C}, \quad (\text{intrinsic collector resistance}) \quad (4.80)$$

which gives the change in collector voltage given a change in collector current. For  $I_C = 1$  mA, for example, this works out to over 360 k $\Omega$ . This expression is useful as an estimate for the effective “output” resistances of the current-mirror transistors and the current-source transistor in the “fancy” differential transistor amplifier of Section 4.11.4.1.<sup>9</sup>

#### 4.11.5.3 Miller Effect

The **Miller effect** says that the collector–base junction, which normally acts like a reverse-biased diode, acts as if it has a small parallel capacitance, on the order of a few pF. (Remember that a reverse-biased junction

<sup>6</sup>Paul Horowitz and Winfield Hill, *The Art of Electronics*, 2nd ed. (Cambridge, 1989), p. 81 (ISBN: 0521370957).

<sup>7</sup>J. M. Early, “Effects of Space-Charge Layer Widening in Junction Transistors,” *Proceedings of the IRE* **40**, 1401 (1952) (doi: 10.1109/JRPROC.1952.273969).

<sup>8</sup>F. A. Lindholm and D. J. Hamilton, “Incorporation of the early effect in the Ebers–Moll model,” *Proceedings of the IRE* **59**, 1377 (1971) (doi: 10.1109/PROC.1971.8435).

<sup>9</sup>This effective resistance is often discussed in terms of the reciprocal of the effective resistance, termed the **output admittance**,  $h_{oe}$  (defined at zero base current). For example, the 2N4401 data sheet from ON Semiconductor specifies a guaranteed range of 1 to 30  $\mu\text{S}$  for  $h_{oe}$  at  $I_C = 1$  mA and  $V_{CE} = 10$  V (dc); this is consistent with our estimate in the text of  $r_c^{-1} \approx I_C/(V_A + V_{CE}) \approx 27 \mu\text{S}$ .



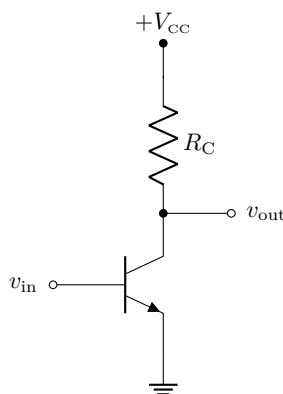
has a depletion region, which acts as a thin, insulating layer.) The main problem is that if a transistor circuit has voltage gain  $G$ , then this Miller capacitance  $C_{CB}$  gets “transferred” to the input as an effective capacitance of  $(1 + |G|)C_{CB}$  (Problem 4.10). With any input impedance, this forms a low-pass filter, so for fast circuits, either the input impedance needs to be kept small, or the gain  $G$  must be small.

#### 4.11.5.4 Variation of $\beta$

Finally, we have already noted that  $\beta$  is not *really* a constant, but this is worth reiterating. It varies between transistors and with temperature. The *only* thing to rely on should be that  $\beta$  is large (100 or more), not that it has any particular value. Note that in terms of gains,  $\beta$  drops out of all the circuits we analyzed; it only appears in the impedance expressions, where its exact value is not critical.

## 4.12 Common-Emitter Amplifier Revisited: Emitter Degeneration

We commented before in Section 4.11.4.1 that shorting out the emitter resistors in the differential transistor amplifier leads to a large gain, because it is the (small) intrinsic emitter resistance  $r_e$  that takes the place of the emitter resistors  $R_E$ . The same holds true of the basic common-emitter amplifier from Section 4.8. In the case that we omit the emitter resistor, we obtain the **grounded-emitter amplifier** shown below.



In this case the ac voltage gain is simply the common-emitter gain (4.30) with  $R_E$  replaced by  $r_e$ ,

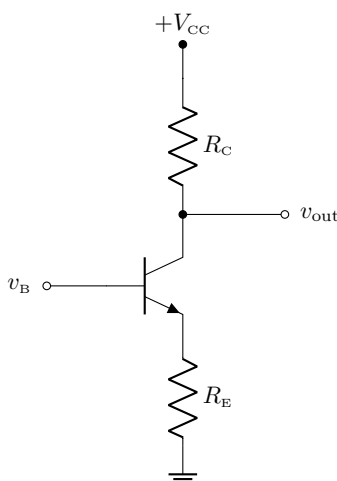
$$G = -\frac{R_C}{r_e} = -\frac{R_C I_C}{nV_T}, \quad (4.81)$$

where we used the expression (4.73) for the emitter resistance. The advantage of this configuration is that the gain is as high as it can get, given  $R_C$  (see, e.g., Problem 4.11). The *disadvantages* include a gain that depends on temperature via  $V_T$ , as well as the collector current. What is worse, this circuit is hard to *bias* properly. From the temperature dependence (4.76) of  $V_{BE}$  in the discussion of the Ebers–Moll model, at fixed  $I_C$  the base voltage changes by  $-2.1 \text{ mV}/^\circ\text{C}$ . Conversely, at fixed  $V_{BE}$ ,  $I_C$  changes by a factor of  $e^{(2.1 \text{ mV}/V_T)}$  for a temperature increase of  $1^\circ\text{C}$ , corresponding to an increase of 9% around  $20^\circ\text{C}$ . This means that a temperature rise of 8 to  $9^\circ\text{C}$  causes the collector current to *double*. This means that a grounded-emitter amplifier with a dc base voltage set to bias the output at  $V_{CC}/2$ , will have a bias that drops to zero with a rise in temperature of less than  $10^\circ\text{C}$ —too unreliable to be useful in practice.

The addition of the emitter resistor in the common-emitter amplifier, shown again below, stabilizes all this. The gain is now

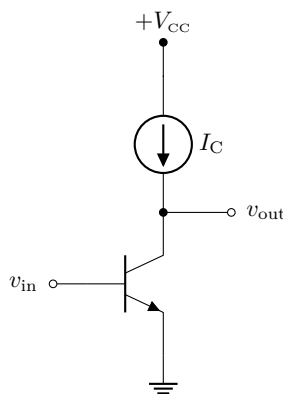
$$G = -\frac{R_C}{r_e + R_E} \approx -\frac{R_C}{R_E}, \quad (4.82)$$

provided  $R_E \gg r_e$ , so the gain is stable with respect to temperature and current. Also,  $R_E$  (with the base bias voltage) sets  $I_C$  and thus the bias voltage (with  $R_E$ ); provided the voltage drop across  $R_E$  is large compared to  $V_{BE}$  (or at least the potential variations in  $V_{BE}$ ), the output bias is stable with respect to temperature, collector current,



This addition of the emitter resistor in the grounded-emitter amplifier is an important idea, and is called **emitter degeneration**. In fact, another poor feature of the grounded-emitter amplifier is substantial distortion for large output signals (see Problem 4.11), because the gain depends on  $I_C$ . The output with a sufficiently large  $R_E$ , however, is virtually distortion-free because the gain ( $-R_C/R_E$ ) is constant. This reduction in distortion (and other annoying behaviors)—at the cost of a reduced gain—is the hallmark of negative feedback (we will discuss this in more depth in Section 7.7). In fact,  $R_E$  provides a kind of negative feedback: for an increase in the input voltage, the emitter voltage likewise increases, suppressing any changes in  $V_{BE}$  (changes that would otherwise occur in the grounded-emitter configuration).

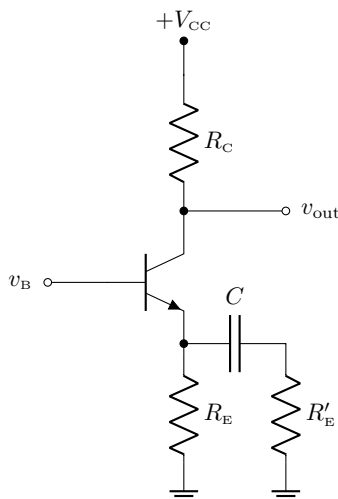
It's worth noting that as in the transistor differential amplifier, it is a common trick to replace the “load”  $R_C$  by a current source, as shown below. Again, this is not an *ideal* current source, which would be completely inflexible about its current  $I_C$ —this corresponds to arbitrarily large resistance, and thus an arbitrarily large amplifier gain (which is not nearly as handy as it might sound). Rather, a *real* current source has an effective large resistance: as a more specific model, you can think of the current source in the diagram below as being bypassed by a large resistor  $R_C$  that sets the ac gain.



Despite the large effective resistance, there is no problem in setting the dc collector current, and thus  $r_e$  (and thus  $G$ ). That is, of course, at constant temperature. This gives good linearity as long as the output is not near an edge of the supply range, under the assumption that the output is unloaded (anything that saps current from the collector will obviously affect the gain). This amplifier still suffers from the problem of thermal drift of the bias voltage, however; an emitter-degeneration resistor can still help here.

Another trick in the same vein is to go back to the basic common-emitter amplifier and bypass  $R_E$  with a capacitor (see Problem 4.6). Then ac signals see a large gain set by  $r_e$ , whereas the dc bias is stabilized by  $R_E$ , taking advantage of emitter degeneration. Even with the stabilized emitter current, the emitter current varies with input (ac) voltage, so this configuration still suffers from distortion (which is okay if it is used as

a high-gain amplifier within a negative-feedback loop). When improved linearity via emitter degeneration is still desirable at ac frequencies, the circuit below is a common amplifier.

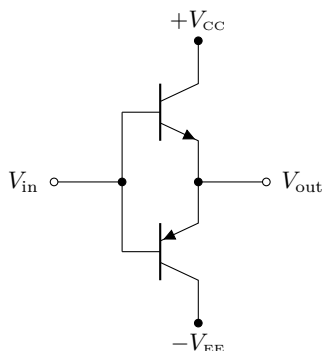


At ac frequencies the emitter resistance is  $R_E \parallel R'_E$  (which is just  $R'_E$  if  $R'_E \ll R_E$ ). This circuit allows the amount of emitter degeneration to be set independently at ac and dc by choosing the two resistors appropriately.

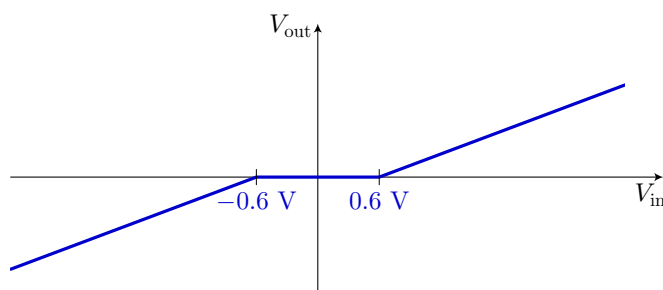
### 4.13 Biasing the Push–Pull Pair

The single-transistor amplifiers that we have studied so far have required a dc bias to keep the desired ac signal away from the lower power supply. In high-power circuits, this is typically a bad thing: the quiescent (dc) current, which is flowing even when there is no ac signal, wastes a lot of power. Thus a high-power emitter follower is typically not a good idea (not without a huge heatsink, anyway).

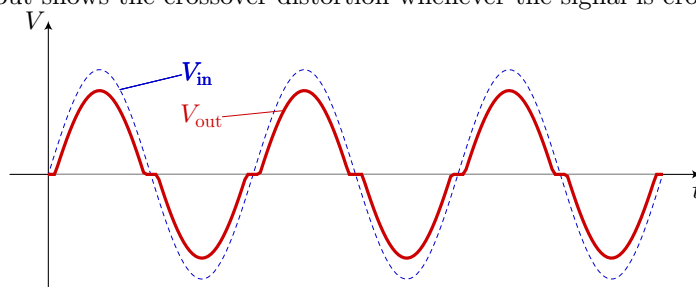
One attempt to improve this is to stack a complementary pair of emitter followers, in a **push–pull** current amplifier, as shown below.



The idea is that, when  $V_{out}$  is connected to ground via a load, the NPN transistor conducts when  $V_{in} > 0$ , and the PNP takes over when  $V_{in} < 0$ . The basic *problem* with this circuit is that to make one of the transistors conduct, the emitter (output) voltage must be a diode drop closer to ground than the base (input) voltage. That is, a graph of the output voltage responding to input voltage looks schematically as in the graph below, if we assume the simple model that the forward-biased base–emitter voltage drop is 0.6 V.



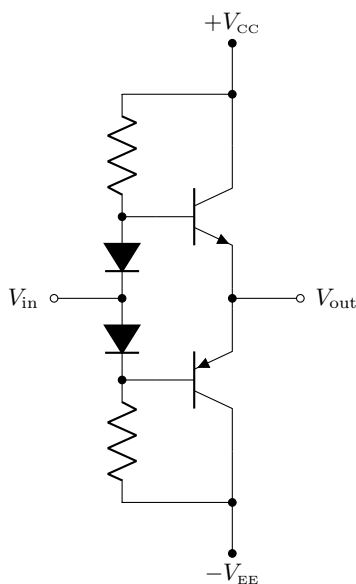
The problem in using this as an amplifier is that it leads to **crossover distortion**, because the base-emitter drop changes as the signal crosses through zero. If you put a relatively large (compared to  $0.6\text{ V}$ ) sine wave into the input, the output shows the crossover distortion whenever the signal is crossing through zero.



Yeesh, that's not good.

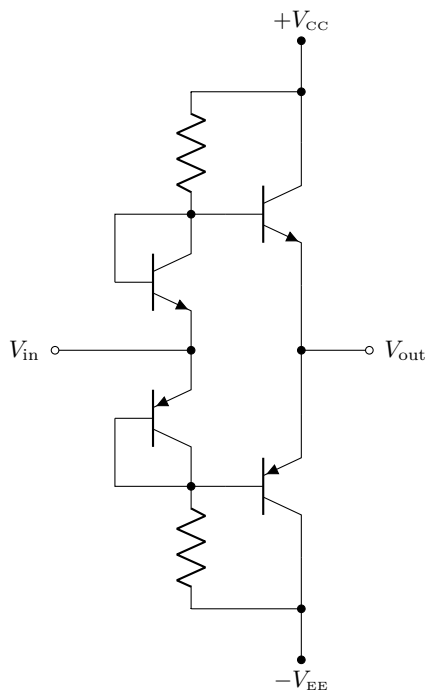
One solution to this problem is to “bias” the transistor bases into conduction: basically, boost the NPN's base by  $0.6\text{ V}$ , and “unboost” the PNP's base by  $-0.6\text{ V}$ . Then with a  $0\text{-V}$  input, both transistors are just at the edge of conduction, but in opposition so their currents cancel. This wastes a bit of power at idle, but largely removes the problem with crossover distortion. The uncorrected amplifier is called a **class-B amplifier**, while the bias-corrected amplifier is called a **class-AB amplifier** (the less-efficient emitter follower is a **class-A amplifier**).

One method of biasing the push-pull pair is to use a voltage divider to generate the  $\pm 0.6\text{-V}$  bias points. Variable resistances can allow the bias points to be fine-tuned. Another common technique is to generate the bias voltages that match the base-emitter bias junctions using another pair of junctions, such as diodes. The basic scheme is shown below.



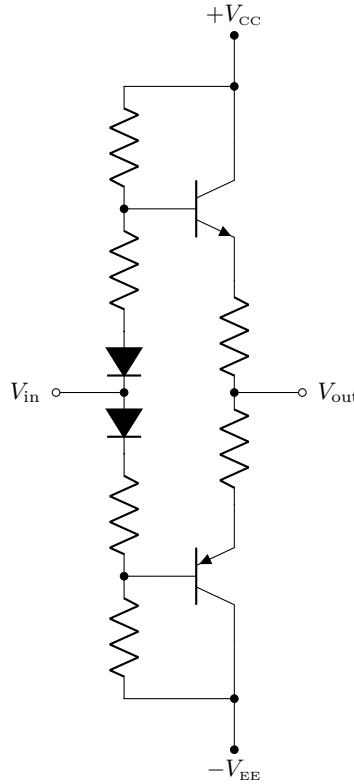
The forward-biased diodes act as part of a stiff voltage divider that set the appropriate bias points; the choice of resistors will affect the impedance seen by the signal source at the input (these could be replaced

by current sources to have high effective input impedance while selecting the proper current to match the transistor  $V_{BE}$ 's). A common variation on the above idea is to use matching biasing transistors in place of the diodes, as shown below.



This better matches the bias voltages to the base-emitter drops, but is probably only practical for signal-level circuits (this is a common trick in the internals of an op-amp, for example).

The one remaining problem is the temperature stability of the bias in high-current applications. The output power transistors heat up as they deliver power to the load; this causes their  $V_{BE}$ 's to drop. If the low-power bias diodes are thermally decoupled from the power transistors, the fixed bias voltage causes the quiescent current to increase, just as we saw in the grounded-emitter amplifier in Section 4.12. This can lead to **thermal runaway**, where the increased quiescent current leads to more heating, which leads to yet more quiescent current, and so on. The solution, as in the grounded-emitter amplifier, is to add emitter-degeneration resistors, as shown below.



This increases the voltage drop from the transistor bases to  $V_{out}$ —the idea is that the voltage drop across the (relatively small) emitter resistors should be large compared to the anticipated change in  $V_{BE}$  due to temperature changes (i.e., larger than  $\sim 100$  mV, but not *too* much larger). Extra (small) resistors have also been added to correspondingly boost the bias voltages in the diode chain; adjustable resistors here are useful in setting the quiescent current (which should be something bigger than zero so that both transistors are in slight conduction, ready to “take over” the signal when it crosses through zero).

## 4.14 Mathematical Modeling of DC BJT Behavior

The basic design equations that we have used so far are the basic Ebers–Moll model (4.60), along with the current-gain equation:

$$\begin{aligned} I_C &= I_S \left( e^{V_{BE}/nV_T} - 1 \right) \\ I_B &= I_C / \beta. \end{aligned} \tag{4.83}$$

We have also inserted some rough voltage drops as appropriate, such as  $V_{BE} \approx 0.6$  V when the base is forward-biased with respect to the emitter, and  $V_{CE} \approx 0.2$  V in saturation. Here we will briefly discuss more detailed mathematical modeling of the transistors, and as a by-product discuss the BJT “characteristic curves” that are used in diagnosis and characterization of transistors, and were in the past used for graphical design of transistor circuits.

### 4.14.1 Ebers–Moll Model

The main issue with the model equations (4.83) are that they are too simplistic for certain circumstances. For example, they predict that  $I_C$  is completely independent of  $V_{CE}$ , so they cannot correctly handle any saturation behavior (without inserting a saturation voltage by hand). However, the same diode-law treatment

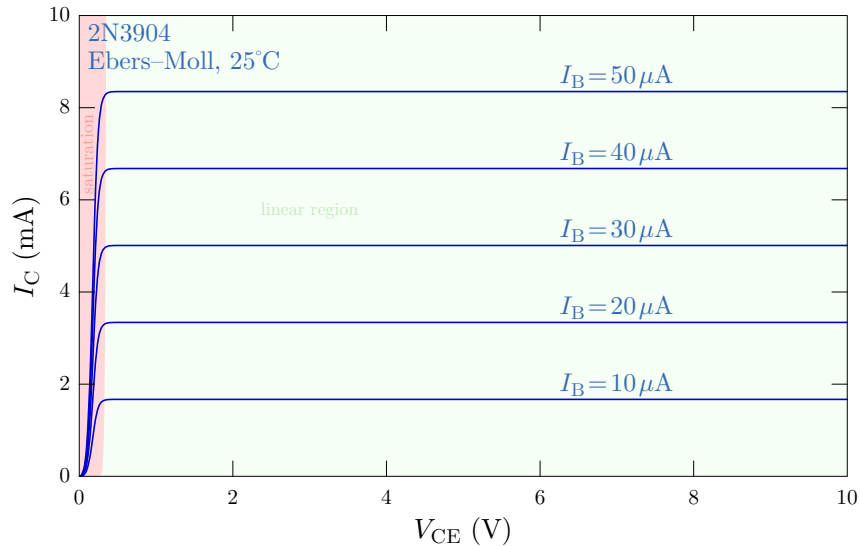
of the base–emitter junction applies to the other junctions, so we can write the more detailed Ebers–Moll equations

$$\begin{aligned} I_C &= I_S \left( e^{V_{BE}/nV_T} - 1 \right) - \left( 1 + \frac{1}{\beta_R} \right) I_S \left( e^{V_{BC}/nV_T} - 1 \right) \\ I_B &= \frac{I_S}{\beta_F} \left( e^{V_{BE}/nV_T} - 1 \right) - \frac{I_S}{\beta_R} \left( e^{V_{BC}/nV_T} - 1 \right). \end{aligned} \quad (\text{detailed Ebers–Moll equations}) \quad (4.84)$$

The first term of the first equation is what we had before. The second term of the  $I_C$  equation treats the base–collector junction, including the collector–emitter current (which goes as  $I_S$ ) and collector–base current (which goes as  $I_S/\beta_R$ ). The  $I_B$  equation has the normal base–emitter current that enforces the  $I_B = I_C/\beta$  equation, with a similar term for the base–collector current. Note that there are separate “forward” and “reverse”  $\beta$  factors, because the two junctions are not equivalent in geometry or semiconductor doping. Also, the idealities are assumed to be the same here for simplicity—otherwise the coefficients will involve idealities as well, and the  $\beta$  factors will not be constant (which is more realistic, but unnecessarily complex within this model).

This model can properly account for saturation, because at saturation the base–collector junction becomes weakly forward-biased, and the forward conduction contributes to the currents. Note that in linear operation (away from saturation)  $V_{BC}$  is typically large and negative, so the second terms in both of Eqs. (4.84) are negligible, justifying the simpler model (4.83) from before.

To visualize where we are with the model, the plot shows typical characteristic curves, showing the “output” current ( $I_C$ ) variation with the applied  $V_{CE}$  for different “input” currents  $I_B$ . Note that  $V_{CE}$  enters the Ebers–Moll equations indirectly via  $V_{BC} = V_{BE} - V_{CE}$ ; the  $I_B$  equation can be solved to yield  $V_{BE}$  in terms of the input  $I_B$ , which then determines  $I_C$  via the first equation. The parameters are appropriate for the popular 2N3904 small-signal transistor; the specific values are  $I_S = 2 \text{ fA}$  (irrelevant for the specific calculation here),  $n = 1.33$ ,  $\beta_F = 167$ , and  $\beta_R = 1$ . The right-hand side (shaded in light green) of the plot corresponds to linear operation, where  $I_C = \beta_F I_B$  holds (hence the even spacing of the curves, where the spacing is proportional to  $\beta_F$ ). Saturation occurs for small  $V_{CE}$  at the left-hand side of the plot (shaded in pink). More specifically, saturation occurs *visually* when the curve acquires a large slope. The exact  $V_{CE}$  where this occurs is not completely well-defined, and depends on  $I_B$ , but it should be clear that 0.2 V is a reasonable start for the saturation voltage. In more quantitative terms, saturation occurs below  $V_{CE} \sim V_{BE}$ , so that the base–collector junction is no longer reverse-biased. The actual saturation boundary shown in the plot is where  $V_{CE}$  is 35% of  $V_{BE}$ .



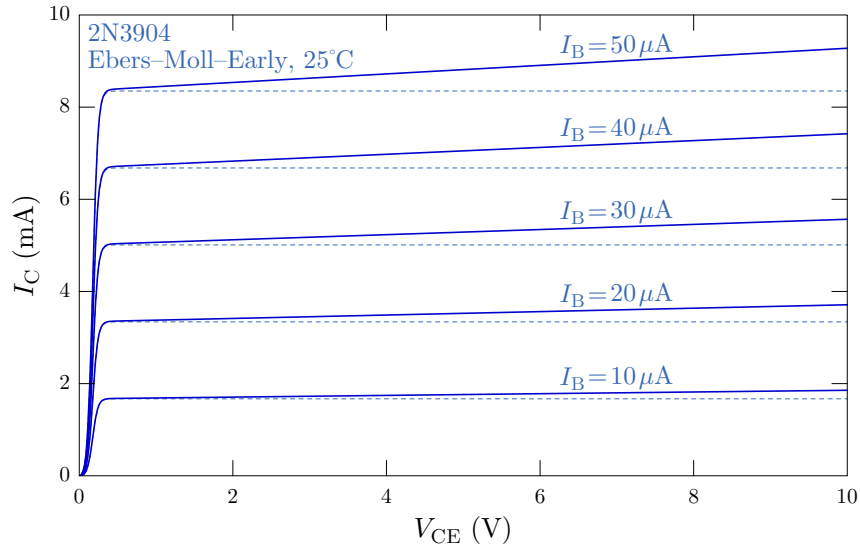
The major remaining defect of this plot is that the curves are essentially horizontal in the linear regime, which indicates a divergent collector resistance  $r_c$ ; this is the next correction to consider below.

### 4.14.2 Ebers–Moll–Early Model

In the above discussion of the Early effect, the narrowing of the base region is modeled by modifying the basic Ebers–Moll equation to have a factor  $(1 + V_{CE}/V_A)$ , as in Eq. (4.77). No such factor should influence the base current, however. With this modification, Eqs. (4.84) read as follows:

$$\begin{aligned} I_C &= I_S \left( e^{V_{BE}/nV_T} - 1 \right) \left( 1 + \frac{V_{CE}}{V_{AF}} \right) - I_S \left( e^{V_{BC}/nV_T} - 1 \right) \left( 1 + \frac{V_{CE}}{V_{AR}} \right) - \frac{I_S}{\beta_R} \left( e^{V_{BC}/nV_T} - 1 \right) \\ I_B &= \frac{I_S}{\beta_F} \left( e^{V_{BE}/nV_T} - 1 \right) - \frac{I_S}{\beta_R} \left( e^{V_{BC}/nV_T} - 1 \right). \end{aligned} \quad (\text{Ebers–Moll equations with Early correction}) \quad (4.85)$$

Note that the correction factor applies only to the two terms that correspond to the collector–emitter current; the base–collector current is unaffected. Since the two terms refer to different junctions, there are in principle distinct forward and reverse Early voltages  $V_{AF}$  and  $V_{AR}$ . Note, however, that  $V_{AR}$  will have negligible effect under normal (forward) operation. Also, these terms can be referenced to  $V_{BC}$  and  $V_{BE}$  in the first and second terms, respectively, instead of  $V_{CE}$  (this will be the case in the next model below). The effect of incorporating an Early voltage of  $V_{AF} = 90$  V for the 2N3904 is shown in the plot below (with  $V_{AR} = 0$  V). Note that the effect of the Early correction is that the curves acquire a slope in the linear regime. Due to the form  $(1 + V_{CE}/V_A)$  of the correction factor, the slopes extrapolate back to a common intercept of the voltage axis at  $-V_A$ . (In a real transistor, the slopes for larger  $I_B$  and  $V_{CE}$  will not quite intercept at the same point as higher-order effects become important.)



The previous Ebers–Moll curves are superimposed as dashed lines to make the difference between models more clear. The curves are no longer horizontal in the linear regime, and their changing slope means that  $\beta$  is effectively no longer independent of  $V_{CE}$ . Also, the Early voltage of 90 V leads to a collector resistance of  $r_c \approx V_{AF}/I_C \approx 20$  k $\Omega$  for the  $I_B = 30$   $\mu$ A curve, where  $I_C \approx 5$  mA.

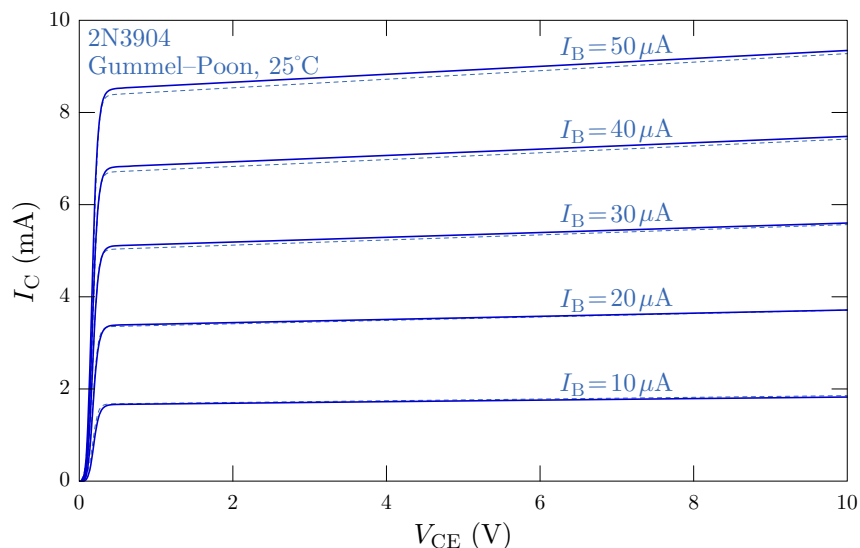
### 4.14.3 Gummel–Poon Model

A substantially more sophisticated transistor is the **Gummel–Poon model**,<sup>10</sup> which is commonly used in computer circuit-modeling software. This model is essentially an extension of the Ebers–Moll–Early model above, but bases the collector–emitter current flow on the base *charge*, rather than voltage or current. It also incorporates modifications to transistor behavior at higher currents, and includes internal resistances in

<sup>10</sup>H. K. Gummel and H. C. Poon, “An Integral Charge Control Model of Bipolar Transistors,” *Bell System Technical Journal* **49**, 827 (1970) (archive: bstj49-5-827) (doi: 10.1002/j.1538-7305.1970.tb01803.x)

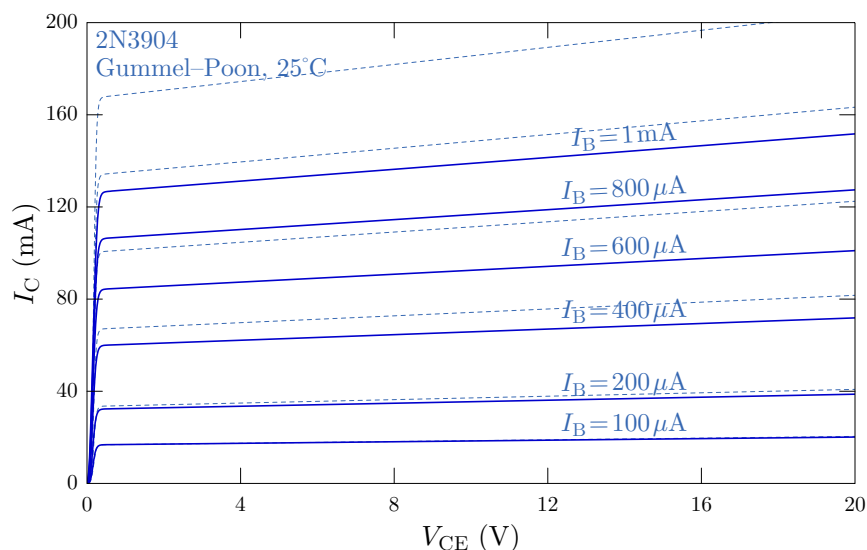


series with all three terminals (with a variable base resistance). The equations in the form commonly used in software are not especially enlightening, so we won't list them here, but they may be found elsewhere.<sup>11</sup> The characteristic curves from this Gummel–Poon model are shown below for the 2N3904.<sup>12</sup>



Again, the curves from the Ebers–Moll–Early model are included as dashed lines to make the differences between models clear. There are clear differences (for example, more variation in  $\beta$  with  $I_B$ , and thus nonlinearity), but the Ebers–Moll–Early model is not far off when the parameters are well-chosen.

The match is good enough in the above plot that you might even wonder if it's worth bothering with this more complicated model. You just have to push harder on the transistor to see more significant deviations, however.



<sup>11</sup>“The Spice Gummel-Poon Model,” Section 5.5 in in the *DynLAB Course on Modeling and Control of Multidisciplinary Systems, Part II: Modelling of Semiconductor Devices*, <https://web.archive.org/web/20090130135801/http://virtual.cvut.cz:80/dynlabcourse/>.

<sup>12</sup>The parameters are taken from the SPICE model published by Central Semiconductor. A number of other models are available for this same device in libraries and by other manufacturers, and they show a surprisingly wide variation in their output. The characteristic curves shown here match measurements reasonably well; the parameters of the Ebers–Moll–Early model were tuned to make those curves match this model. One final note: the SPICE model parameters should be viewed as *fitting* parameters. That is, a value for  $\beta_F$  listed in a SPICE model may not be a good representation for the  $\beta$  of the model, because other parameters can enhance or counteract the effect of  $\beta_F$ .

This plot is getting closer to the limits of this transistor, which is rated for a maximum  $I_C$  of 200 mA, and a maximum  $V_{CE}$  of 40 V (at  $I_B = 0$ ). One physical feature *still* typically missing from the more complete models is breakdown at high  $V_{CE}$  (and reverse breakdown), but this is not a big issue provided you respect the device limits published in the data sheet!

## 4.15 Little-“h” Notation

In the development of the BJT here, we have stuck to a few relatively simple parameters, such as  $\beta$ , maybe an Early voltage, and of course the various input and output currents and voltages. We have also briefly mentioned a couple of other parameters, like  $h_{FE}$  and  $h_{fe}$  (p. 93), and  $h_{oe}$  (Section 4.11.5.2). This notation is widespread in transistor data sheets, and so it’s worth explaining this briefly.

The idea is that for any of the basic transistor-amplifier circuits that we have studied, there is an input and output voltages, as well as input and output currents. Mostly we have stuck to relating input and output voltages, but the currents can also influence or be caused by the voltages. Suppose we consider very generally how  $V_{in}$ ,  $I_{in}$ ,  $V_{out}$ , and  $I_{out}$  are interrelated. For not-very-obvious reasons it is conventional to take  $I_{in}$  and  $V_{out}$  as given, and the other variables to be functions of these:

$$\begin{aligned} V_{in} &= V_{in}(I_{in}, V_{out}) \\ I_{out} &= I_{out}(I_{in}, V_{out}). \end{aligned} \quad (4.86)$$

Considering small variations  $\Delta V_{in}$  and  $\Delta I_{out}$  of the “outputs” due to small variations  $\delta I_{in}$  and  $\delta V_{out}$  of the “inputs” gives, by the chain rule,

$$\begin{aligned} \Delta V_{in} &= \frac{\partial V_{in}}{\partial I_{in}} \Delta I_{in} + \frac{\partial V_{in}}{\partial V_{out}} \Delta V_{out} \\ \Delta I_{out} &= \frac{\partial I_{out}}{\partial I_{in}} \Delta I_{in} + \frac{\partial I_{out}}{\partial V_{out}} \Delta V_{out}. \end{aligned} \quad (4.87)$$

It is then conventional to represent the partial derivatives with “constants” with labels of the form  $h_{xy}$ . The “h” comes from this being the “hybrid model” of the transistor. The different possible values for the  $y$  subscript are:

- e: for the common-emitter amplifier (Section 4.8); this is the option most commonly found in data sheets.
- c: for the common-collector amplifier (i.e., the emitter follower, Section 4.6)
- b: for the common-base amplifier (see Problem 4.3)

The basic idea is that a BJT has three terminals; one will act as an input, another an output, and the third a common voltage reference for the other two (i.e., ground). This leaves three sensible options for transistor configurations, corresponding to grounding one of each of three terminals (the input and output are fixed by requiring that the transistor have some amplifying action).

The first subscript  $x$  then takes on four different values to denote one of the partial derivatives in Eqs. (4.87);

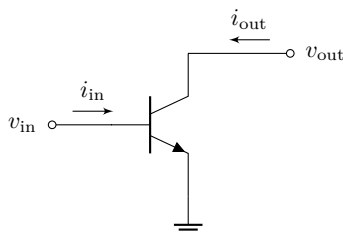
- i: for the **input impedance**  $\partial V_{in}/\partial I_{in}$
- r: for the **reverse voltage ratio** or **voltage feedback ratio**  $\partial V_{in}/\partial V_{out}$
- f: for the **forward current ratio** or **small-signal current gain**  $\partial I_{out}/\partial I_{in}$
- o: for the **output admittance**  $\partial I_{out}/\partial V_{out}$

Thus, for a common-emitter amplifier, Eqs. (4.87) become

$$\begin{aligned} v_{\text{in}} &= h_{\text{ie}} i_{\text{in}} + h_{\text{re}} v_{\text{out}} \\ i_{\text{out}} &= h_{\text{fe}} i_{\text{in}} + h_{\text{oe}} v_{\text{out}}, \end{aligned}$$

(hybrid-model equations, common-emitter configuration) (4.88)

if we use ac notation for the changes in voltage and current instead of the “delta” notation of Eqs. (4.87). For the common-emitter amplifier,  $i_{\text{in}} = i_{\text{B}}$  and  $i_{\text{out}} = i_{\text{C}}$ : the “output current” of the amplifier is the current flowing into the “output,” which is the collector, as shown below. (Note that since these parameters characterize only the transistor, the reference circuit omits any external resistances, but we assume it to be biased properly to work as an amplifier.)



Thus, the second equation is just the transistor law if we identify  $\beta = h_{\text{fe}}$  (recall that the second term is an Early-effect correction due to  $v_{\text{CE}}$ , which is small because the collector’s intrinsic resistance  $r_{\text{c}}$  is large).

The lower-case subscripts indicate that these are ac quantities (i.e., they refer to partial derivatives at some nonzero bias values of  $I_{\text{in}}$  and  $V_{\text{out}}$ ). These are sometimes reported as dc quantities, in which case the subscripts are written in upper cases. The analogous relations to Eqs. (4.88) for the dc case would then be

$$\begin{aligned} V_{\text{in}} &= h_{\text{IE}} I_{\text{in}} + h_{\text{RE}} V_{\text{out}} \\ I_{\text{out}} &= h_{\text{FE}} I_{\text{in}} + h_{\text{OE}} V_{\text{out}}. \end{aligned} \quad (4.89)$$

For our purposes, we have been referring to both  $h_{\text{FE}}$  and  $h_{\text{fe}}$  interchangeably as  $\beta$ , although they are not the same.

To give some sense of typical values, for the 2N4401,<sup>13</sup> at a (dc) bias current  $I_{\text{C}} = 1 \text{ mA}$ ,  $V_{\text{CE}} = 10 \text{ V}$ , and a signal frequency of 1 kHz, the guaranteed range of each parameter is listed below.

parameter	min value	max value
$h_{\text{ie}}$ (input impedance)	1 k $\Omega$	15 k $\Omega$
$h_{\text{re}}$ (reverse voltage ratio)	$10^{-5}$	$8 \times 10^{-4}$
$h_{\text{fe}}$ (forward current gain)	40	500
$h_{\text{oe}}$ (output admittance)	1 $\mu\text{S}$	30 $\mu\text{S}$

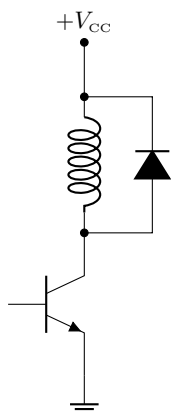
Note that the intrinsic emitter resistance  $r_{\text{e}}$  is a useful parameter that doesn’t appear among these parameters. However, in the common-emitter configuration, you can view the input impedance as the emitter resistance, but amplified by roughly a factor of  $\beta$  [Eq. (4.12)]. Thus, we can estimate  $r_{\text{e}}$  based on the above parameters via  $h_{\text{ie}} \approx h_{\text{fe}} r_{\text{e}}$ . Taking typical values (“unit 1” from the data-sheet plots) of  $h_{\text{fe}} \approx 200$  and  $h_{\text{ie}} \approx 5.5 \text{ k}\Omega$ , this gives  $r_{\text{e}} \approx 28 \Omega$  at  $I_{\text{C}} = 1 \text{ mA}$ . This agrees with the prediction of  $r_{\text{e}} = nV_{\text{T}}/I_{\text{C}} \approx 28.3 \Omega$  from the Ebers–Moll model (at 25°C and assuming  $n = 1.1$ ).

## 4.16 Circuit Practice

### 4.16.1 Transistor Switching an Inductive Load

Consider the circuit below, where an NPN transistor switches an inductive load (electromagnet, motor, etc.).

<sup>13</sup>as specified by ON Semiconductor, <http://www.onsemi.com/pub/Collateral/2N4401-D.PDF>



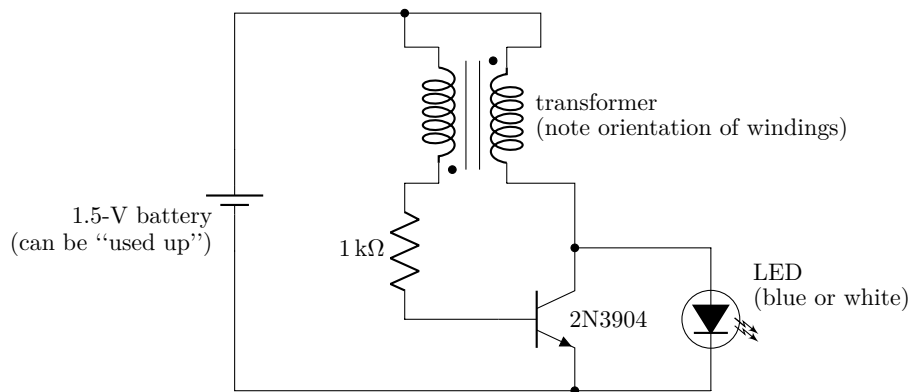
Explain what you need to do to switch the inductive load on and off. Also, why the diode, which is necessary for large inductances, to protect the transistor when switching off the current?

**Solution.** Without any connection to the transistor base, the inductor is switched off. To run current through the load (say a current  $I$ , limited by the resistive part of the load impedance), we need to inject about  $1/100$  ( $1/\beta$ ) of  $I$  into the base, e.g., using a voltage and a resistor as a voltage-to-current convertor.

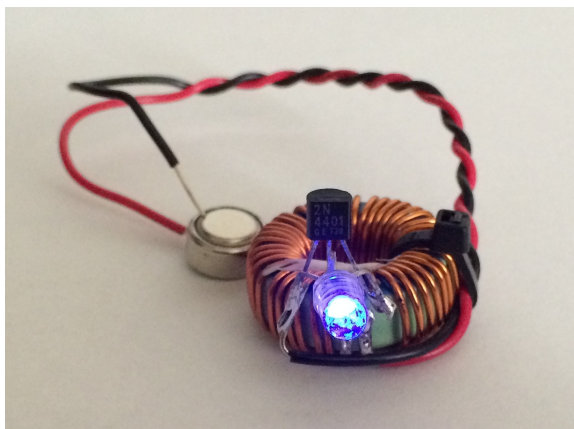
The problem with this setup is that when switching off the current, the inductor will develop a large EMF to try to sustain the current. This will pull the collector voltage far above  $+V_{CC}$  if the inductance is large and the switching is rapid [remember the EMF is  $L(dI/dt)$ ]. This can exceed the transistor breakdown voltage, destroying the transistor. When  $V_C > V_{CC}$ , the diode shorts the collector to  $V_{CC}$ , clamping the voltage at  $V_{CC}$  and protecting the transistor.

### 4.16.2 Joule Thief

Let's take a break from quantitative analysis of transistor circuits, and look at a fun and elegant circuit. The circuit below is called the **joule thief**.<sup>14</sup> To understand the name, consider that the circuit is powered by a 1.5 V battery, but as we discussed before in Section 4.5, turning on a blue or white LED takes about 3.3 V. The “thievery” comes from noting that this circuit works even with a sagging battery with an even lower voltage—hence, we can get steal every last joule from the battery (here, for the purposes of lighting up the LED).



<sup>14</sup>For more on this circuit, including instructions on how to build this, see <http://www.evilmadscientist.com/2007/weekend-projects-with-bre-pettis-make-a-joule-thief/>. See the rest of the site for a lot more cool stuff on electronics and other subjects.

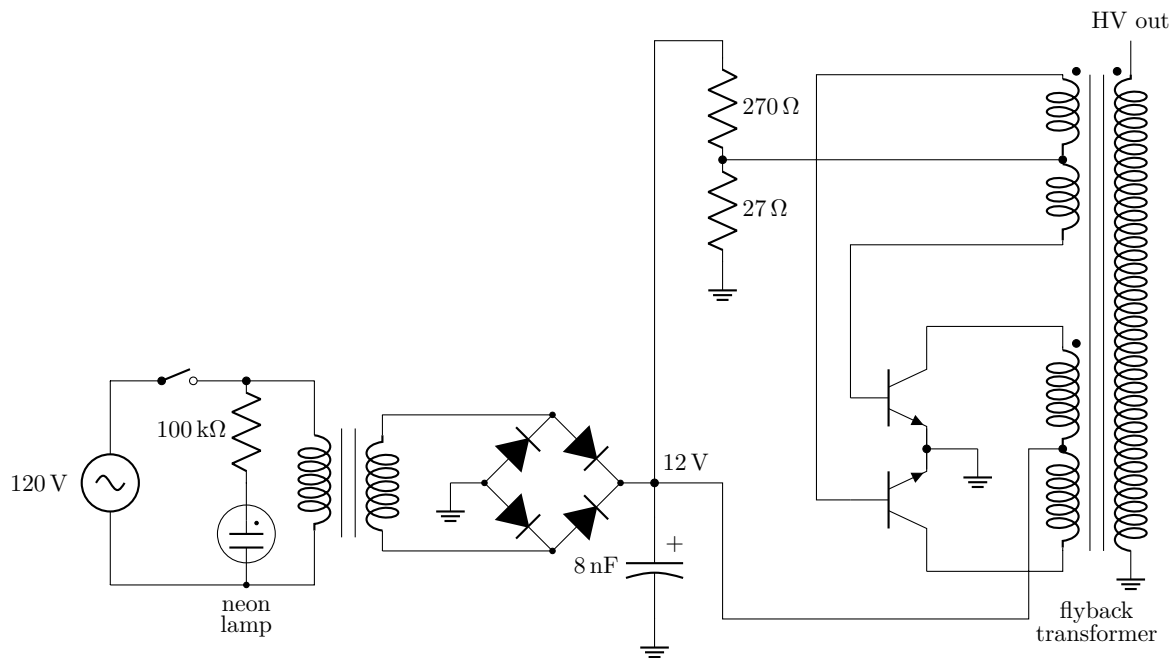


Your exercise here is to trace through the circuit and explain how it works. Just treat the transistor as a switch; no need to do any calculations. Also, as a hint, note that inductors are fairly rare in circuits; usually they show up either as filters (as in the outputs of switching power supplies for computers), or—as is the case here—they are around to give an inductive kick when they are interrupted (like the spark-plug coil in gasoline-powered car engines).

**Solution.** Suppose we start with the transistor off. Current starts flowing from the battery, but it can't go into the collector. The only path it has is into the base; the resistor is there to make sure we don't overdo the current. This base current turns on the transistor, so a much larger current can start to flow into the collector. As the current flows through the secondary of the transformer to the collector, the *opposite* current is induced in the primary (note the primary and secondary are wound/oriented oppositely). This induced current opposes and thus interrupts the base current, switching off the transistor. This interrupts the collector current, interrupting the relatively large current in the transformer secondary. This leads to an inductive kick, and the collector voltage builds up potentially far above the battery voltage (recall the same thing happens for a transistor switching off an inductive load, as in Section 4.16.1). The collector thus builds up to a sufficiently high voltage that the LED turns on, dumping the inductive kick. Then we are back to the initial state, and the process repeats.

### 4.16.3 Solid-State Tesla Coil

Here is another good circuit to practice *qualitative* reading of transistor operation.



The goal of this circuit<sup>15</sup> is to develop high voltages (of order 25 kV) by using a flyback transformer from a television (the old, cathode-ray-type, which needs large voltages to accelerate the electron beam). The left-hand side of the circuit is basically a dc power supply to drive the right-hand side. The neon lamp (a tiny, neon-filled glass discharge tube) is a handy way of indicating 120-V power; the lamp needs about 90 V to “fire,” and then once the discharge starts, little current is needed to sustain it, so the voltage is regulated by the 100-kΩ resistor.

The flyback transformer is driven through a center-tapped primary coil by an alternating pair of power transistors, which are driven by a center-tapped “feedback” coil. Trace through the circuit to understand how the transistor pair switches between the “on-off” and “off-on” pair states. To get started, assume that both transistors are initially off, and pick an arbitrary path for the current from the voltage divider driving the feedback coil.

#### 4.16.4 Eric Clapton Signature Stratocaster Preamplifier

Look at the linked schematic diagram<sup>16</sup> for the Eric Clapton Signature Stratocaster from Fender Musical Instruments Corp.<sup>17</sup> The guitar emulates Clapton’s beloved old “Blackie” guitar, but the preamp has an unusual feature of enabling a boost in the midrange audio band. This allows electronic control of the guitar’s tone from the traditional “Strat” sound to something more akin to a Gibson Les Paul (which featured “humbucking” pickups, which have more midrange gain).

This schematic is somewhat awkward to read, but you can still look through it and try to spot a few elements. In particular, look at transistors Q1–Q4 and identify what amplifiers each one makes up. Also try to identify the bias network in each case.

#### 4.16.5 Op-Amp Internals

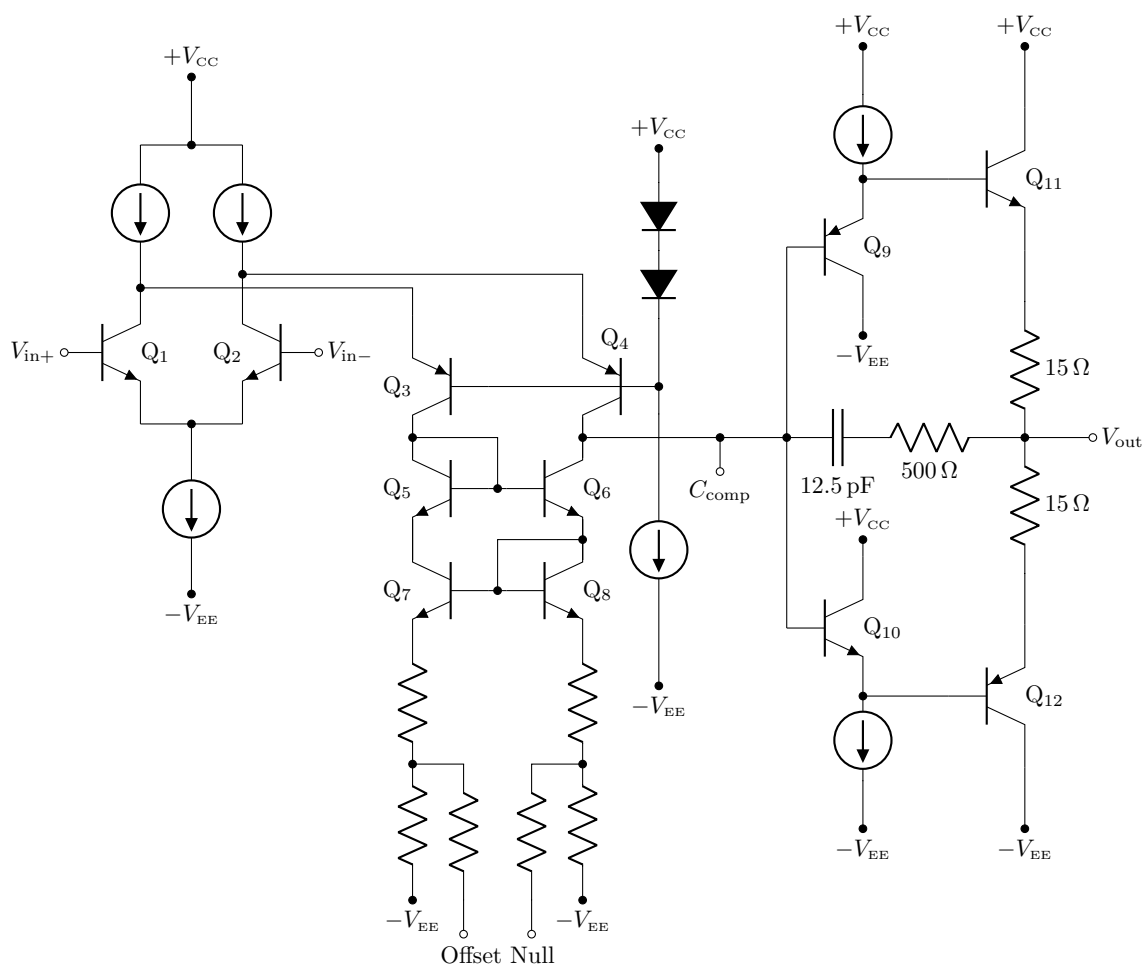
The schematic below shows the simplified schematic diagram of the internals of the AD829 high-speed, low-noise op-amp from the data sheet.<sup>18</sup> Isn’t it pretty?

<sup>15</sup>Robert E Iannini, *Build Your Own Laser, Phaser, Ion Ray Gun and Other Working Space Age Projects* (McGraw-Hill, 1983) (ISBN: 0830606041).

<sup>16</sup>[http://www.blueguitar.org/new/schem/\\_gtr/ec\\_schem\\_fact.jpg](http://www.blueguitar.org/new/schem/_gtr/ec_schem_fact.jpg)

<sup>17</sup><http://www.fender.com/guitars/stratocaster/eric-clapton-stratocaster/product-011760.html>

<sup>18</sup><http://www.analog.com/media/en/technical-documentation/data-sheets/AD829.pdf>; see also Analog Devices’ *Op-Amp Basics* by James Bryant, Walt Jung, and Walt Kester, <http://www.analog.com/media/en/training-seminars/>

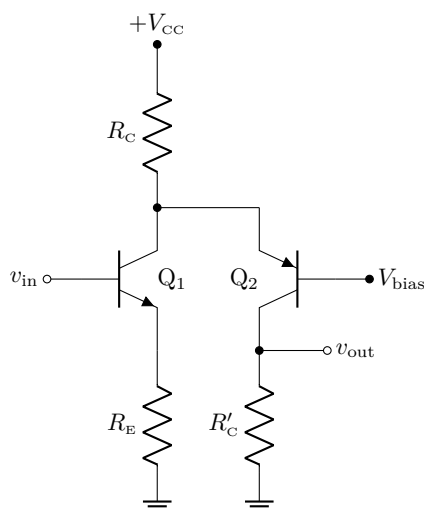


The point of this exercise is to go through and understand the functions of all parts of this circuit, based on what you've learned in this chapter.

Before you start, note that one transistor circuit here that we haven't yet studied is the **folded cascode amplifier**, shown below, which has ac voltage gain (Problem 4.14)

$$G = -\alpha_1 \alpha_2 \frac{R'_C}{R_E} \approx -\frac{R'_C}{R_E}, \quad (4.90)$$

where  $\alpha_1 = \beta_1/(\beta_1 + 1)$  for  $Q_1$  and  $\alpha_2 = \beta_2/(\beta_2 + 1)$  for  $Q_2$ .



That is, this circuit behaves fairly similarly to the common-emitter amplifier, but with two transistors instead of one (functioning something like a common-emitter amp cascaded with a common-base amp). The main advantages of this circuit are that it suppresses the Miller and Early effects.

**Solution.**  $Q_1$  and  $Q_2$  form a transistor differential amplifier; the current source to  $-V_{EE}$  acts as the “long tail,” or an effectively large  $R_{EE}$ , while the other two current sources appear to act as large  $R_C$  resistances (the  $R_E$  resistors are omitted for high gain based on the transistor  $r_e$ ). But this is a little more complicated, because  $Q_3$  and  $Q_4$  are part of the differential amplifier two; the “transistors” involved in the differential amplifier are actually folded-cascode amplifiers instead of common-emitter amplifiers as we originally learned the differential amplifier. The main benefit is probably the reduced Miller effect for high speed of the input stage. The two-diode network in series with the current source is the bias circuit for the two cascode circuits (biased close to  $V_{CC}$  to allow a wide swing voltage for the cascode outputs), and the “loads” of the cascode circuits is the pair of current mirrors formed by  $Q_5$  with  $Q_6$  and  $Q_7$  with  $Q_8$ . The current mirror is a good load to replace  $R_C$  resistors to obtain high gain in a differential amplifier, but here the mirror is the load at the outputs of the of the folded-cascode circuits rather than the directly at the collectors of the input transistors. Why the stacked mirrors? Recall that there is some error in the current mirror associated with the small base currents flowing through only one of the emitters; this trick of stacking two current mirrors in opposite senses cancels this base-current error, and is called an **improved Wilson current mirror**. The resistors to  $-V_{EE}$  provide some emitter degeneration increase the effective collector impedance at the output (the collector of  $Q_6$ ). Of course it also gives a method of trimming the balance of the current mirror using the offset-null pins (where a potentiometer with wiper connected to  $-V_{EE}$  would be connected).

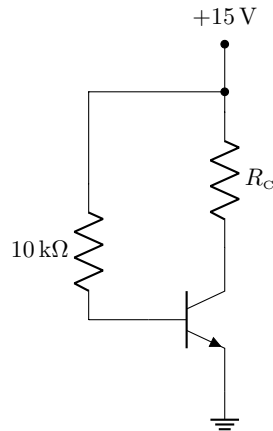
Now the output of the current mirror drives a push-pull output stage, with complementary output pair  $Q_{11}$  and  $Q_{12}$ . The output transistors have emitter-degeneration resistors for thermal stability.  $Q_9$  and  $Q_{10}$  function to bias the bases of the output transistors, but also serve as first-stage emitter followers (with high effective  $R_E$ ), so the output current is boosted in two stages (for good output current while minimizing loading effects on the current mirror; the data sheet specifies a current gain of 40 000 for the output stage). The  $C_{comp}$  is a terminal for an external compensation capacitor (to ground), which can roll off the amplifier gain at high frequencies for stability. The remaining  $R$ - $C$  pair is a little more complicated. Note that at low frequencies, this capacitor is “bootstrapped”, meaning that the voltage on either side of the capacitor should be the same, in which case the capacitor is irrelevant. At higher frequencies, when the output of the push-pull stage lags the input (especially with a capacitive load connected), this  $R$ - $C$  network provides a (negative) feedback path to the push-pull input, which effectively also rolls off the gain of this output stage at high frequencies (again, to help ensure stability).



## 4.17 Exercises

### Problem 4.1

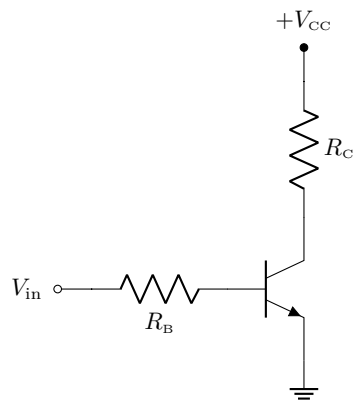
In the circuit below, for what range of resistances  $R_C$  will the transistor be saturated?



Assume  $\beta = 100$ .

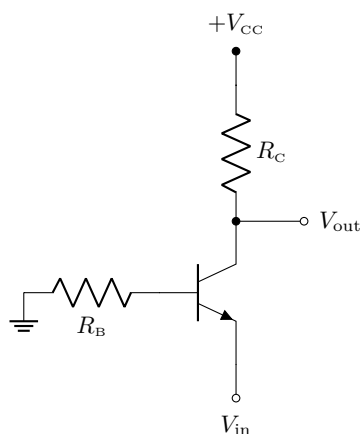
### Problem 4.2

Consider the following transistor circuit. Consider  $V_{CC}$ ,  $R_B$ , and  $R_C$  to be fixed. For what range of input voltage  $V_{in}$  is the transistor saturated?



### Problem 4.3

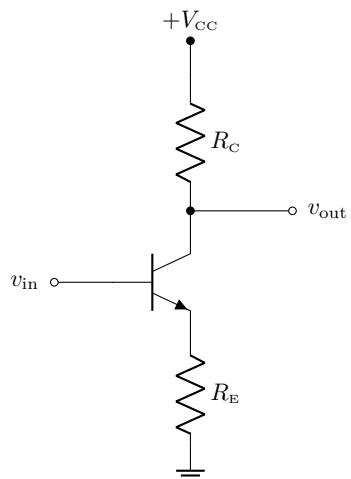
Consider the transistor circuit below. Assume the input is biased such that the transistor works normally, and ignore any internal resistances in the transistor. For the calculations in this problem, you may assume  $\beta \gg 1$ .



- Consider small voltage changes  $v_{in}$  and  $v_{out}$  at the input and output (i.e., ignore bias voltages). Compute the voltage gain  $G = v_{out}/v_{in}$ .
- Compute the input impedance.
- Compute the output impedance.

#### Problem 4.4

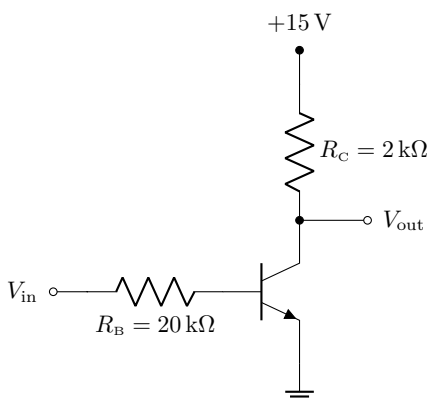
Consider the common-emitter circuit below.



In Section 4.8, we showed the ac voltage gain of the circuit is  $G = -R_C/R_E$ , assuming  $\beta \gg 1$ . Rederive the ac gain of the circuit, this time **without** making any approximations involving  $\beta \gg 1$ .

#### Problem 4.5

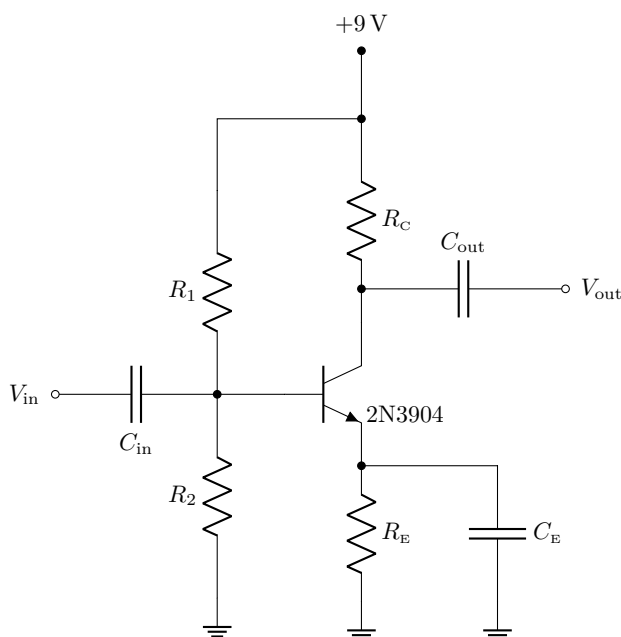
Consider the amplifier shown below.



- (a) What is the minimum value of  $V_{in}$  that saturates the transistor? Assume  $\beta = 100$ .
- (b) What is the (ac) voltage gain of the amplifier, assuming the input is biased to something above 0.6 V and below your result from (a)? (Why should you *not* think of this circuit as a grounded-emitter amplifier in the sense of a common-emitter amp with  $R_E = 0$ ?)

#### Problem 4.6

Consider the common-emitter amplifier below. This is a fairly involved problem, but the idea is to give you a template for how to design a more realistic transistor amplifier.



The design criteria for this amplifier are: a large gain over a usable bandwidth of 100 Hz–20 kHz, an input impedance of 10 kΩ, a quiescent current (dc current) of 100  $\mu$ A (dominated by the collector current). The dynamic range of the output should also be reasonably close to the maximum possible, so we will fix the collector voltage at  $+V_{CC}/2$ . Also assume an ideality value of  $n = 1.09$  for the 2N3904.

- (a) Now we must set the gain of the circuit, and so we need to explain the function of  $C_E$ . The idea is that at frequencies above 20 Hz, the capacitor bypasses the resistor, so that over the amplifier bandwidth the intrinsic emitter resistance sets the voltage gain, not  $R_E$ . Derive an expression for the

collector current in terms of  $R_C$ , which is yet to be determined. (Start with the expression for the ac voltage gain of the amp.)

(b) Use your result from (a) to show that the (ac) voltage gain in this circuit can be written

$$G = -\frac{V_{R_C}}{nV_T}, \quad (4.91)$$

where  $V_{R_C}$  is the voltage drop across  $R_C$ .

(c) Using the specified collector voltage, what is the gain at 20°C?

(d) What are  $R_C$ ,  $r_e$ , and  $I_B$ ? Use  $\beta$  ( $h_{FE}$ ) from the data sheet.<sup>19</sup>

(e) For the emitter resistor, a good choice for  $R_E$  is often about 10% of  $R_C$ . What is this resistance, and what is  $V_E$ ? The reason for this choice is as follows. Recall that  $V_{BE}$  varies with slightly with temperature. If we fix  $V_B$ , then this variation leads to a temperature dependence of  $V_E$  and thus the bias current. To minimize this effect,  $V_E$  should be much larger than the variation in  $V_{BE}$ . Verify explicitly that this is the case for this circuit, over, say, variations of 10°C (i.e., estimate the effect on the quiescent current for this temperature change).

(f) What is  $V_B$ ? (Use the data sheet for any values you need, don't assume standard values.) Now choose  $R_1$  and  $R_2$ , accounting for the design specs and not forgetting to account for the effect of  $R_E$  on the input impedance. A computer here is probably helpful in solving the resulting equations; ask for help with this if you need it.

(g) Set  $C_E$  by guaranteeing that its impedance is negligible compared to  $r_e$  (*not*  $R_E$ ) over the amplifier's ac bandwidth.

(h) Choose  $C_{in}$  and  $C_{out}$ .

(i) As a sanity check, note that transistors have a parameter  $h_{oe}$ , called the **output admittance**, defined by

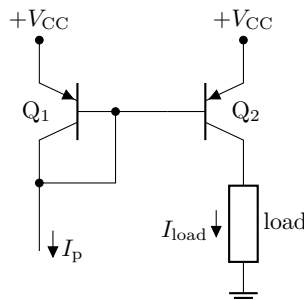
$$h_{oe} := \frac{\partial I_C}{\partial V_{CE}}. \quad (4.92)$$

Note that this has units of conductance, and in fact  $1/h_{oe}$  acts as an effective collector resistance that appears in parallel with  $R_C$ . Verify from the data sheet that we are justified in ignoring this.

(j) Give an order-of-magnitude estimate for the life of a 9-V battery powering this circuit. (Look at some battery data sheets to find some useful information; cite any sources you use.)

### Problem 4.7

Consider the current mirror below.

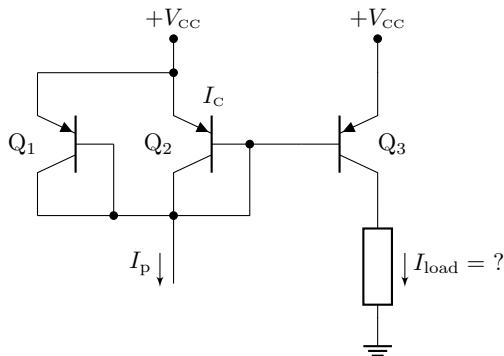


Suppose that  $Q_1$  is at temperature  $T_1$  and  $Q_2$  is at temperature  $T_2$  (in class we assumed these were the same). Derive an expression for the current  $I_{load}$  in terms of the program current  $I_p$  and the temperatures. Assume that the transistors are otherwise identical.

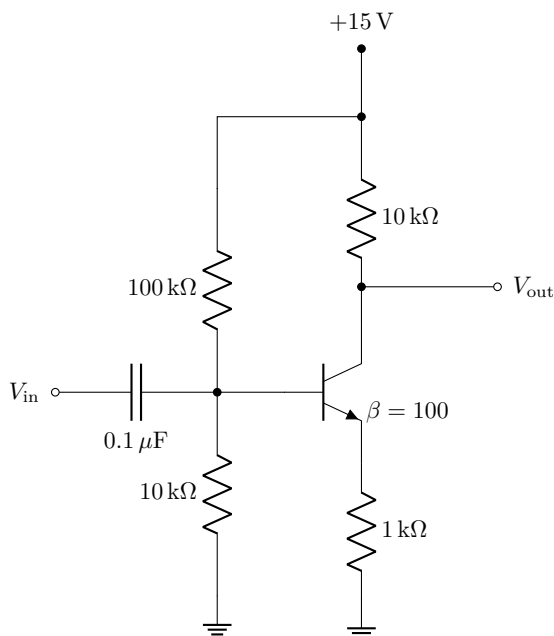
<sup>19</sup><http://www.fairchildsemi.com/ds/2N/2N3904.pdf>

**Problem 4.8**

The circuit below is a variation on the transistor current mirror. Compute  $I_{\text{load}}$ , assuming that the three transistors are identical.

**Problem 4.9**

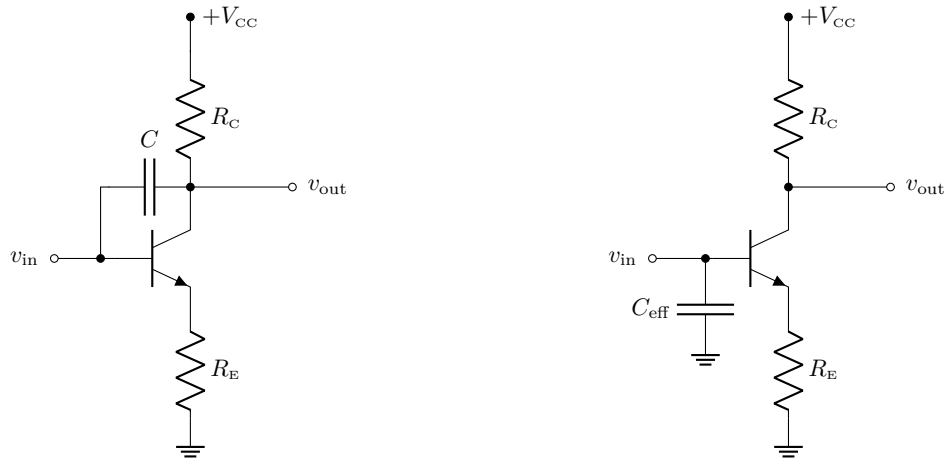
Consider the circuit below.



- What is the input impedance at high frequencies?
- What is the dc bias level of  $V_{\text{out}}$ ?
- How does the intrinsic emitter resistance enter into the analysis of this circuit? For the specific component values here, how much does it affect the quantities you calculated in (a) and (b)?

**Problem 4.10**

Consider the common-emitter amplifier shown below, assumed to be biased as needed for the circuit to work normally as an amplifier. Recall that the **Miller effect** refers to an intrinsic capacitance that appears across the reverse-biased base-collector junction, shown as capacitance  $C$  in the left-hand diagram.



The reason that the Miller capacitance is a serious problem in a high-gain circuit is that seen as an effective input capacitance  $C_{\text{eff}}$  as shown in the right-hand diagram, the effective capacitance is

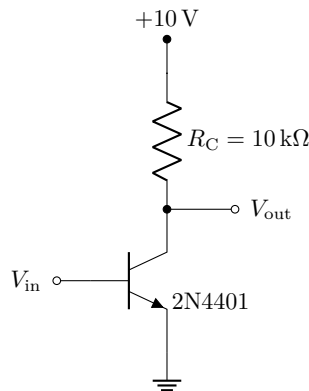
$$C_{\text{eff}} = (1 + |G|)C, \quad (4.93)$$

where  $G \approx -R_C/R_E$  is the ac voltage gain of the amplifier.

Show that this statement is true by considering the change  $q$  in the charge  $Q$  on the capacitor due to a change  $v_{\text{in}}$  in the input voltage.

#### Problem 4.11

Consider the **grounded-emitter amplifier** (i.e., a common-emitter amplifier with  $R_E = 0$ ) shown below.



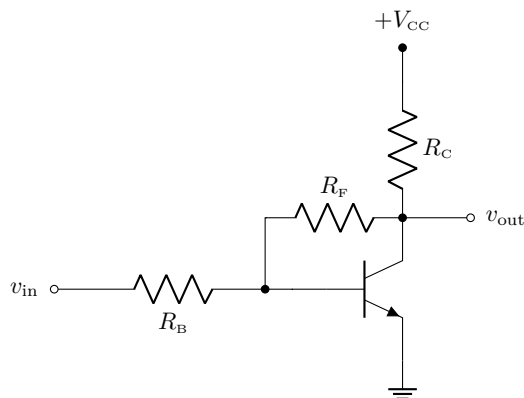
- What is the minimum value of  $V_{\text{in}}$  that saturates the transistor? Assume a transistor temperature of  $20^\circ\text{C}$ , with parameters  $I_{\text{CS}} = 1.75 \text{ pA}$  and  $n = 1.1$ .<sup>20</sup>
- Suppose you set the dc bias of  $V_{\text{in}}$  to  $0.54 \text{ V}$ . What is the (ac) voltage gain of the amplifier? (Recall that in this configuration, the gain is limited by the intrinsic emitter resistance.)
- Note that, in practice, this circuit has some drawbacks. The emitter resistance varies with  $I_C$ , and thus with  $V_{\text{in}}$ . That means that both the gain and the input impedance of the transistor depend on  $V_{\text{in}}$ , leading to distortion in the output signal. To see this distortion explicitly, make a plot of the (dc) output voltage  $V_{\text{out}}$  as a function of the (dc) input  $V_{\text{in}}$ , assuming  $V_{\text{in}}$  to be a  $36\text{-mV}$  triangle wave

<sup>20</sup>The parameters in this problem are from the 2N4401 model in LTspice (<http://www.linear.com/solutions/ltspace>).

(peak-to-peak). Make sure to include the input dc bias from (b), which should prevent any clipping of the output signal.

### Problem 4.12

Compute the ac voltage gain for the transistor amplifier shown below, which is a grounded-emitter amplifier, but with a negative-feedback path via resistor  $R_F$  from the collector to the base. Assume the input is biased as needed to make the transistor work normally.

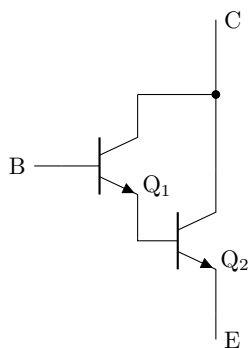


*Hint:* first solve the circuit *without* resistor  $R_F$  in the circuit, then adapt your solution to the case where  $R_F$  is present.

Why is the feedback via  $R_F$  *negative*? (Explain **briefly**.)

### Problem 4.13

A **Darlington pair** is a pair of transistors arranged as in the diagram below; the idea is that it acts as a single, “super” transistor. Darlington pairs integrated to work like a regular BJT are available, such as the TIP120 power transistor. In this application, this circuit helps address the issue that power transistors tend to have lower  $\beta$ ’s than small-signal transistors.



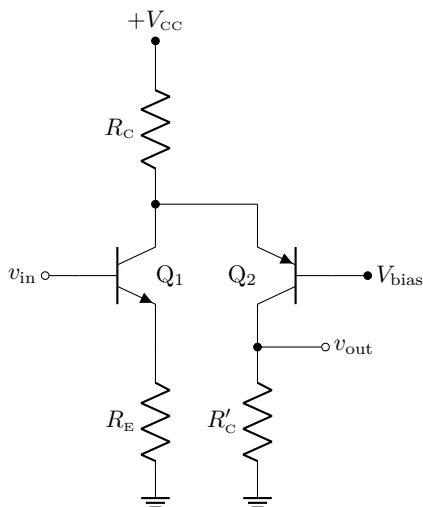
Assume transistor  $Q_1$  has parameters  $\beta_1$ ,  $V_{BE1}$ ,  $r_{e1}$ ,  $h_{oe1}$ , etc., while  $Q_2$  has parameters  $\beta_2$ ,  $V_{BE2}$ ,  $r_{e2}$ ,  $h_{oe2}$ , etc.

- Compute the effective  $\beta$  of the Darlington pair (in terms of  $\beta_1$  and  $\beta_2$ ).
- Compute the effective intrinsic emitter resistance  $r_e$  (in terms of the individual  $r_e$ ’s and  $\beta$ ’s).
- Recall that  $1/h_{oe}$  acts as an intrinsic resistance as seen at the collector of a BJT. Compute the effective  $h_{oe}$  of the pair. (You may take  $h_{oe} := i_C/v_{CE}$  as the defining relation of the output admittance.)

(d) Derive an effective Ebers–Moll equation for the pair, noting that some of the effective “constants” may not be completely constant. Assume the same temperature, ideality, and forward saturation current for both transistors.

#### Problem 4.14

Consider the **folded cascode amplifier** shown below. Here,  $V_{\text{bias}}$  is an externally imposed, constant bias voltage.



- Derive an expression for the ac voltage gain.
- Explain why this circuit does not suffer from the Miller effect (Problem 4.10). Be *brief*.
- Explain why this circuit is not subject to the Early effect (Section 4.11.5.2). Be *brief*.





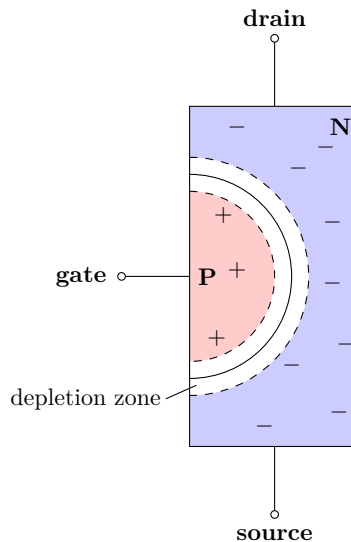
## Chapter 5

# Field-Effect Transistors and Semiconductor Switching Devices

**Field-effect transistors** are three-terminal devices, like bipolar junction transistors. Although they operate somewhat differently from BJTs, they can also be used as amplifiers and switches. By considering how they work in basic circuits, we'll see some of the advantages and disadvantages of FETs relative to BJTs. In this chapter we will discuss FETs as well as some related semiconductor devices that are useful as switches, particularly in power-switching applications.

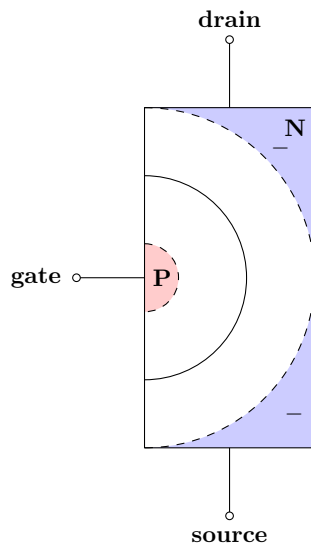
### 5.1 JFET (Depletion-Mode FET)

A **junction FET (JFET)** is basically a p-n junction with a special geometry and three terminals. JFETs always work as **depletion-mode** devices, which refers to the depletion zone at the junction, which is responsible for the switching action of the JFET. The basic scheme for an **n-channel JFET** is shown below (the **p-channel** counterpart is basically the same, under the exchange of p- and n-type semiconductors).



As in the semiconductor diode, a depletion zone forms at the p-n junction. In normal operation, the voltage at the gate terminal is kept below that of the drain and source, so the junction is reverse-biased. Only a small leakage current flows via the gate terminal. However, in the configuration shown, there is a low-impedance path between the drain and source terminals, through the n-type region (hence, the “n-channel”).

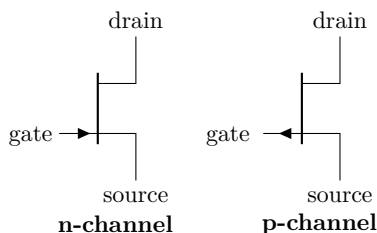
However, when the gate is brought to a negative voltage with respect to drain and source, the depletion zone expands as in the reverse-biased diode. At a sufficiently negative voltage, the depletion zone “pinches off” the n-channel, preventing current flow from drain to source.



For intermediate gate voltages, the gate voltage acts as a control that modulates the resistance of the drain-source path. This acts something like modulating the flow of water in a garden hose by changing the clamping force of a pair of pliers on the hose.

Note that the convention in the n-channel JFET is that **current flows from drain to source** (i.e., the n-type carriers are flowing from the source and out the drain). The current flows from source to drain in the p-channel JFET. These two terminals appear to be interchangeable according to the above diagrams, and for some devices they are so in practice. However, due to geometric differences, the drain and source are not always equivalent. For example, the drain and source often have different capacitances, so reversing them can affect the speed of a JFET amplifier.

The symbols for n-channel and p-channel JFETS is shown below. The difference between the two transistors is only the direction of the gate arrow, which indicates the orientation of the p-n junction (like the arrow in the diode symbol).



The asymmetric placement of the gate differentiates the drain and source.

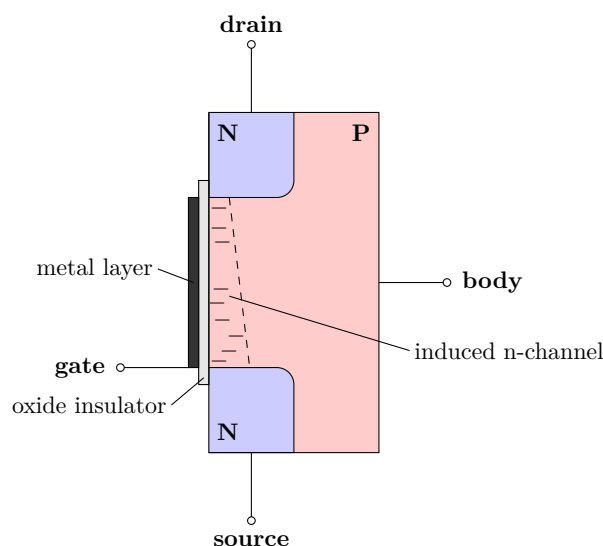
To summarize the operation of the n-channel (depletion-mode) JFET:

- If  $V_{GS} = 0$ , then current can flow through the n channel, typically from drain to source (i.e.,  $I_D \approx I_{DS} > 0$ ).
- If  $V_{GS} < 0$ , the junction is reverse-biased, and the expanding depletion zone restricts current flow.
- There is some **threshold voltage**  $V_{\Theta}$ : if  $V_{GS} < V_{\Theta}$ , then the n-channel is “pinched off,” and no current flows ( $I_D \approx I_{DS} = 0$ ). Typically  $V_{\Theta}$  ranges from  $-2$  to  $-15$  V. *Note:* it’s common for a symbol like  $V_T$  to be used for both the threshold voltage  $V_{\Theta}$  and the thermal voltage  $V_T$  of diode-law and Ebers–Moll fame; make sure not to confuse these!

- The forward-biased case  $V_{GS} > 0$  doesn't normally happen. You may as well just use a diode.
- The JFET is thus a **transconductance device**, where a voltage controls a current (like the BJT, from the standpoint of the Ebers–Moll equation).

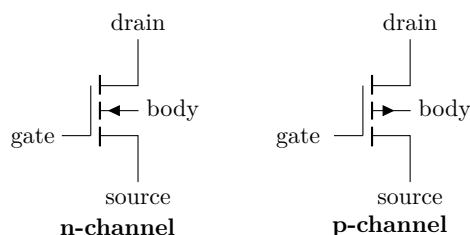
## 5.2 MOSFET (Enhancement-Mode FET)

A second major class of FETs is the **metal-oxide-semiconductor FET (MOSFET)**, also called the **insulated-gate FET (IGFET)**. The names refer to the gate terminal, which is connected to a metal conducting layer, and insulated by an oxide layer from the semiconductor regions of the MOSFET. Some MOSFETs behave as depletion-mode devices, but these are relatively rare. Here we will focus on the much more common case of **enhancement-mode MOSFETs**. The basic scheme for an n-channel MOSFET is shown below (note that the body is a p-type semiconductor, unlike the n-channel JFET; the operating principle is quite different).



The drain and source terminals connect to n-type regions, which are embedded in the p-type substrate or body. In principle MOSFETs have a separate body connection, but these are not always explicitly available (often the body and source terminals are combined into one terminal). The idea here is that if no voltage is applied to the gate, no drain–source current can flow because it will be blocked by one reverse-biased p–n junction. However, when a *positive* control voltage (with respect to drain and source) the gate's  $E$ -field pulls *n-type* carriers out of the p-type substrate (the p-type carriers are the **majority carrier** in the p-type semiconductor, but n-type carriers are also present as the **minority carrier**). The n-type carriers bunched against the gate insulator form an effective n-channel, or **induced n-channel**, that bridges the drain and source. The diagram shows the induced n-channel as being somewhat asymmetric, as appropriate if  $V_{DS} > 0$ .

The symbols for n-channel and p-channel enhancement-mode MOSFETs are shown below. Again, there are extra body connections, but often these are internally shorted to the source.



Also, the three short, vertical lines that represent the three semiconductor regions are sometimes drawn as a single line. Again, the asymmetry of the gate distinguishes drain from source, and the orientation of the arrow distinguishes n- and p-channel types, by indicating the orientation of the p-n body-drain and body-source junctions.

To summarize the operation of the n-channel (enhancement-mode) MOSFET:

- If  $V_{GS} = 0$ , no current  $I_D$  can flow, because of a reverse-biased junction as in the NPN transistor.
- If  $V_{GS} > 0$ , the gate field induces an n-channel, and current  $I_D$  can flow, typically with  $I_D > 0$ .
- Due to the oxide insulator, there is very little gate-current leakage. However, the insulator layer is easily damaged by static discharges.
- The switching operation is similar to the n-channel JFET, but the threshold voltage  $V_{\Theta} > 0$ . That is, conduction occurs when  $V_{GS} > V_{\Theta}$ , so the gate voltage must be *positive* for the MOSFET. Remember the gate voltage is generally *negative* in the JFET, with conduction turning off for sufficiently negative voltages.

MOSFETs are *surface-effect* devices, because the “action” happens at the surface of the oxide layer, whereas the action in JFETs happens in the bulk of the semiconductor. Thus MOSFETs tend to have lower gain compared to JFETs, because mobility is higher for bulk carriers than for surface carriers. (And p-channel MOSFETs tend to have lower gain still because surface hole mobility is lower than surface electron mobility.)

MOSFETs are very common in digital circuits, in the form of **complementary MOS (CMOS)** circuits, where n- and p-channel MOSFETs are paired together, because CMOS pairs can be fabricated very compactly (this is by far the dominant transistor type in modern microprocessors).

### 5.3 Quantitative FET Behavior

Having established that JFETs and MOSFETs have similar behavior except for the particular value (i.e., sign) of the threshold  $V_{\Theta}$ , we can treat both cases together, so long as we track the control voltage  $V_{GS}$  *relative to*  $V_{\Theta}$ . For small input signals (small ac signals on a dc bias), the transconductance nature of a FET means we can write

$$i_D = g_m v_{GS}, \quad (5.1)$$

(FET transconductance relation)

where  $g_m$  is the **transconductance**. This has dimensions of  $\Omega^{-1}$ , which is often written  $\mathcal{U}$  and called a **mho** (or **siemens**, abbreviated “S,” if you want to be all SI about it). The transconductance depends, however, on the bias levels  $V_{GS}$  and  $V_{DS}$ . This relation is the analogue of  $I_C = \beta I_B$  for BJTs.

The FET current formulas can describe the current-voltage characteristics, like the Ebers–Moll equation for BJTs, but in a piecewise way.

1. **Linear region:** when  $V_{DS} < (V_{GS} - V_{\Theta})$ , we have

$$I_D = 2k \left[ (V_{GS} - V_{\Theta})V_{DS} - \frac{V_{DS}^2}{2} \right] \quad (5.2)$$

(FET in linear region)

where  $k$  is a conductance parameter, which is device-dependent and scales with temperature as  $T^{-3/2}$ . (The threshold voltage  $V_{\Theta}$  also depends on temperature.) Note that due to the quadratic term in  $V_{DS}$ , this relation is not really “linear.” However, this term is negligible if  $V_{DS} \ll (V_{GS} - V_{\Theta})$ , or if more linear behavior is desirable, there are tricks to compensate for this term. In the “really” linear case where we can ignore the quadratic term, we have a resistance

$$R = \frac{1}{2k(V_{GS} - V_{\Theta})}. \quad (5.3)$$

(FET resistance in linear region)

This says that the FET in this regime is useful as a voltage-controlled resistor, for example to control variable gain or attenuation in a circuit.

Also in the linear region, the transconductance is

$$g_m = 2kV_{DS}, \quad (\text{FET transconductance in linear region}) \quad (5.4)$$

using Eq. (5.1) in the form  $g_m = \partial I_D / \partial V_{GS}$ .

2. **Saturation or quadratic region:** when  $V_{DS} > (V_{GS} - V_{\Theta})$ . In this region,  $I_D$  is independent of  $V_{DS}$ , in some sense like the saturated BJT where the transistor no longer directly modulates the collector current. In this region,

$$I_D = k(V_{GS} - V_{\Theta})^2, \quad (\text{FET in saturation region}) \quad (5.5)$$

and differentiating again to find  $g_m$ ,

$$g_m = 2k(V_{GS} - V_{\Theta}) = 2\sqrt{kI_D}. \quad (\text{FET transconductance in saturation region}) \quad (5.6)$$

As useful observation in this region is thus to use  $g_m \propto \sqrt{I_D}$ , so that  $g_m$  can be found at different  $I_D$  even when it is only known at one point.

3. **Subthreshold region:** when  $I_D$  is small, because the gate voltage is close to or below the  $V_{\Theta}$ . The simplest model for this region is that the current decreases linearly until  $V_{GS}$  goes below  $V_{\Theta}$ , and at that point the current ceases. A more careful model in this region gives diode-law-like behavior

$$I_D = I_0 e^{(V_{GS} - V_{\Theta})/nV_T}, \quad (\text{FET in subthreshold region}) \quad (5.7)$$

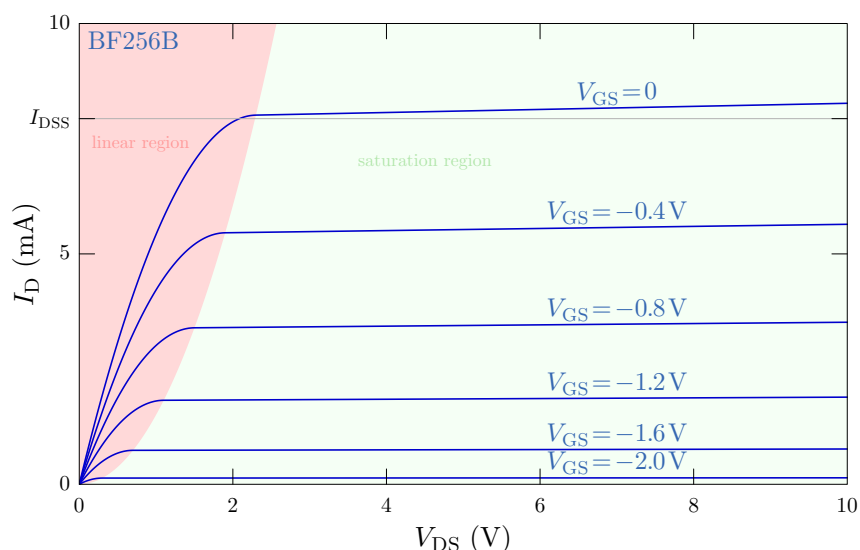
where  $V_T$  is the thermal voltage from the diode law and Ebers–Moll equations. Note that in a JFET, drain current via the gate can become significant. Then  $g_m$  is given by

$$g_m = \frac{I_0}{nV_T} e^{(V_{GS} - V_{\Theta})/nV_T} = \frac{I_D}{nV_T}, \quad (\text{FET transconductance in subthreshold region}) \quad (5.8)$$

according to the exponential law.

### 5.3.1 Visualization

It's much easier to understand the modeling relations as given above by visualizing them. As with the characteristic curves for BJTs (Section 4.14), there are standard characteristic curves for FETs. One example shown below is for the BF256B, a radio-frequency, small-signal, n-channel (depletion-mode) JFET. Again the control signal  $V_{GS}$  labels the different curves, and the results give the drain current  $I_D$  as it varies with applied drain–source voltage  $V_{DS}$ .



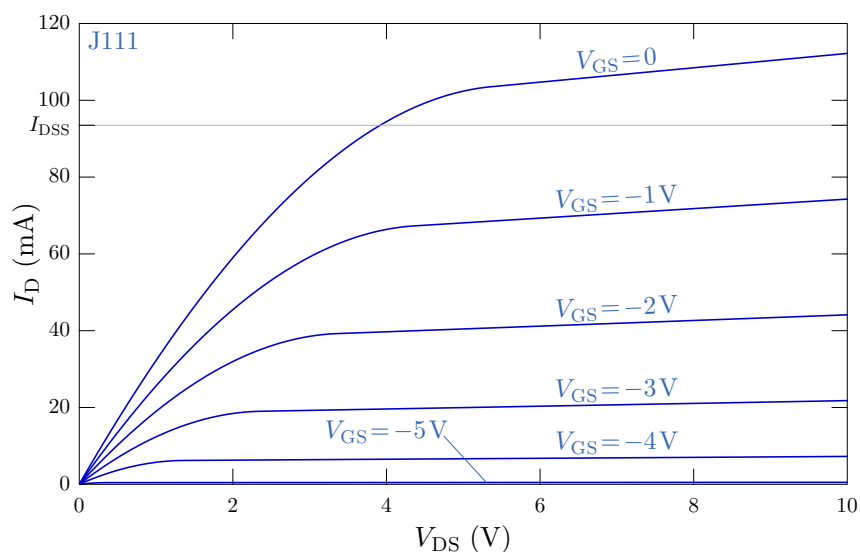
For this transistor model<sup>1</sup> the threshold is  $V_{\Theta} = -2.3$  V (but ranges from  $-0.5$  V to  $-8$  V in the data sheet!); also  $I_{DSS}$  ranges from 6 to 13 mA in the data sheet.<sup>2</sup> Recall that the boundary between the linear and saturation regions is given by the condition  $V_{GS} - V_{\Theta} = V_{DS}$ , although these labels are unfortunately reversed from what you might expect after looking at an analogous diagram of BJT characteristic curves (where recall that saturation occurs to the left and linearity occurs to the right). In the FET case, the “linearity” of the linear region refers to the linear behavior for small  $V_{DS}$  (and small  $I_D$ , where the device acts as a variable resistor across the drain–source connections (controlled by  $V_{GS}$ ). In the saturation region, the drain current is insensitive to the applied  $V_{DS}$ , but is modulated quadratically with  $V_{GS}$  (hence this is also the quadratic region). Note that Eq. (5.5) predicts *no* dependence on  $V_{DS}$ , but the SPICE-model equations used in the plot allow some linear variation in this region [otherwise, the equations are basically (5.2) and (5.5)]. The subthreshold region is not shown (and not commonly included in SPICE models), but would occur along the bottom edge of the plot. Also, for large enough  $V_{DS}$ , the JFET will of course eventually break down (past the right-hand edge of the plot), causing the characteristic curves to turn sharply upward; this effect is also not illustrated here.

As a somewhat different example, the characteristic curves for the J111 n-channel JFET (designed for switching larger currents) is shown below. In this model,<sup>3</sup>  $V_{\Theta} = -5.36$  V, and  $I_{DSS}$  is much larger, almost 100 mA. The behavior is generally similar to that of the 2N3819, but on a larger scale of  $I_D$  and  $V_{GS}$ . Note also the larger slopes in the saturation region, particularly as  $V_{GS}$  approaches zero.

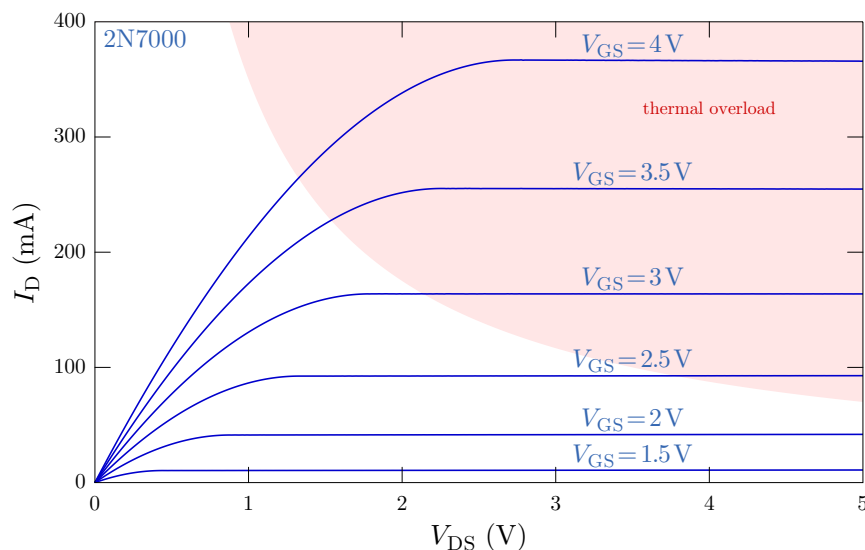
<sup>1</sup>The parameters are from the SPICE model from the MicroSim library, [http://ppd.fnal.gov/experiments/cdms/old\\_files/electronics/FLIP/3U/QampDiscrete/proto/schematics\\_layouts/JFET.LIB](http://ppd.fnal.gov/experiments/cdms/old_files/electronics/FLIP/3U/QampDiscrete/proto/schematics_layouts/JFET.LIB).

<sup>2</sup><https://www.fairchildsemi.com/datasheets/BF/BF256B.pdf>

<sup>3</sup>The parameters here are from the Philips SPICE model; see <http://web.rfoe.net:8000/ziliaoxiazai/PHILIPS/models/spicespar/data/j111.html>.



To show just one more example, the characteristic curves of a popular small-signal (n-channel, enhancement-mode) MOSFET are shown in the plot below.



The basic feature of enhancement mode is apparent, because now all the curves correspond to *positive*  $V_{GS}$  (the model used to make this plot assumed  $V_{\Theta} = 1$  V. The modeling for a MOSFET is considerably more complex than that of a JFET,<sup>4</sup> although the basic shape of the curves is the same. Note the much larger “gain” (i.e., transconductance, or amount of current modulated by a gate-voltage change) compared to the JFET in the plots above; recall that this is *not* very characteristic of an enhancement-mode device (which relies on lower surface mobility vs. larger bulk mobility). In fact, this small-signal transistor is only rated for 350 mW, so much of the plot corresponds to a potentially fried MOSFET (shaded “thermal overload” region). One final curious feature is that in the saturation regime, the curves look quite level, but actually *decrease* slightly with increasing  $V_{DS}$ —this indicates a large, but *negative* effective saturation impedance as seen at the drain (in the grounded-source amplifier).

<sup>4</sup>See, for example, the *HSPICE MOSFET Models Manual* ([http://www.ece.tamu.edu/~spalermo/ecen474/hspice\\_mosfet.pdf](http://www.ece.tamu.edu/~spalermo/ecen474/hspice_mosfet.pdf)), which goes into many possibilities. The plot here assumes a “LEVEL 3” model, with parameters published by Supertex Inc.

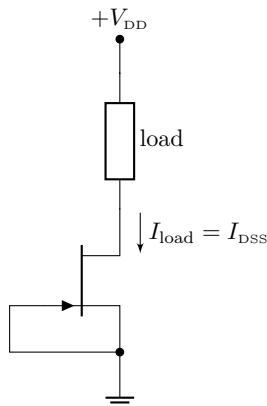


## 5.4 Basic FET Circuits

Now let's see how to apply FETs in some basic circuits. Many of the JFET circuits here have analogues to BJT circuits, and analyzing these will both introduce analysis techniques for FET circuits and highlight the differences as compared to BJTs.

### 5.4.1 JFET Current Source

A simple JFET circuit is the JFET current source, shown below.



The idea is to operate the JFET in the saturation region, where  $I_D$  is independent of  $V_{DS}$ , allowing the JFET to adjust the load voltage to regulate its current. Since the gate and source are shorted (and grounded),  $V_{GS} = 0$ . Then using Eq. (5.5), we find the constant current

$$I_D = kV_{GS}^2 =: I_{DSS}, \quad (5.9)$$

(JFET regulated current)

where  $I_{DSS}$  is the drain current at zero gate–source voltage (i.e., effectively the maximum JFET current).

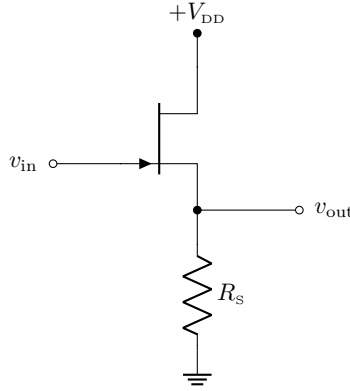
Looking at the characteristic curves above for the 2N3819, you can indeed see that  $I_D$  is relatively insensitive to  $V_{DS}$  above about 2 V in this application, so it makes a fairly good current source (the effective resistance with the slope in the plot as shown is 38 kΩ). From looking at the characteristic curves of the J111 JFET above, it is apparent that it is not especially good as a current source (with an effective resistance of about 7 kΩ). From looking at these same plots, it is apparent that it is also possible to vary the regulated current by reducing  $V_{GS}$  (say, by inserting a resistor, especially a variable resistor, between the JFET source and ground).

The advantage of this circuit, compared to the BJT current source, is its simplicity. The disadvantage is more serious:  $I_{DSS}$  varies significantly between devices. However, it is possible to get hand-picked FETs with particular values of  $I_{DSS}$ , called **current-regulator diodes** (something like Zener diodes, but for current instead of voltage). A good example is the 1N5283–1N5314 series, which covers a nominal current range of 0.22–4.70 mA.<sup>5</sup>

### 5.4.2 JFET Source Follower

The next circuit is the JFET analogue of the BJT emitter follower.

<sup>5</sup>[https://www.centralsemi.com/get\\_document.php?cmp=1&mergetype=pd&mergepath=pd&pdf\\_id=1N5283-5314.PDF](https://www.centralsemi.com/get_document.php?cmp=1&mergetype=pd&mergepath=pd&pdf_id=1N5283-5314.PDF)



Ignoring offsets, Ohm's law for the source resistance gives

$$v_S = i_{DS} R_S, \quad (5.10)$$

while the transconductance relation (5.1) gives

$$i_{DS} = g_m v_{GS} = g_m (v_G - v_S). \quad (5.11)$$

[Note that typically  $V_{DD}$  is large enough to put the JFET into the saturation region—this only takes a few V. Then  $g_m$  is given by Eq. (5.6).] The latter equation becomes

$$v_G = \frac{i_{DS}}{g_m} + v_S = \frac{i_{DS}}{g_m} + i_{DS} R_S. \quad (5.12)$$

after solving for  $v_G$  and using Eq. (5.10). Then the ac voltage gain is

$$G := \frac{v_{out}}{v_{in}} = \frac{v_S}{v_G} = \frac{i_{DS} R_S}{i_{DS}/g_m + i_{DS} R_S}, \quad (5.13)$$

after using Eqs. (5.10) and (5.12). This simplifies to

$$G = \frac{g_m R_S}{1 + g_m R_S}. \quad (5.14)$$

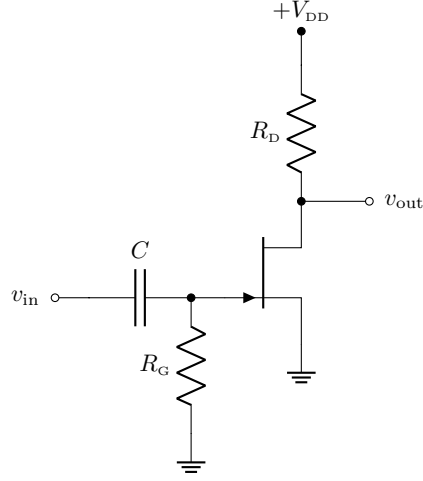
(voltage gain, JFET source follower)

Note that  $G \approx 1$  (hence, a follower) for large resistance  $R_S \gg 1/g_m$ . For a typical signal JFET like the 2N5485,  $g_m \approx 5000 \mu\text{S}$ , so  $1/g_m \approx 200 \Omega$ , so have a larger resistance than the internal FET resistance is not difficult.

This circuit has high input impedance, because the input is essentially a reverse-biased diode. However, the output impedance is basically  $1/g_m$ , which could be somewhat large compared to the emitter follower. Another disadvantage of this circuit is the unpredictable dc offset, since  $V_{GS}$  for a certain current is not well-controlled in the fabrication process.

### 5.4.3 JFET Common-Source Amplifier

Next is a JFET common-source amplifier, the FET analogue of the common-emitter amplifier.



Here,  $R_G$  maintains the dc bias of  $V_{GS}$  at ground, but allows ac modulation of the gate via the capacitor  $C$ .

If  $V_{GS} = 0$ , the quiescent current is  $I_D = I_{DSS}$ , as in the JFET current source (thus assuming that  $V_{DD}$  is large enough to put the JFET into saturation). Then the output is biased at

$$V_{out} = V_{DD} - I_D R_D. \quad (5.15)$$

This bias level may be adjusted by introducing a source resistor, thus lowering  $V_{GS}$  and decreasing  $I_{DS}$ , though biasing this amplifier is more difficult than its BJT counterpart, due to large variation among particular devices in the JFET's  $I_D$  at a given  $V_{GS}$ .

For small ac signals,

$$v_{out} = v_D = -i_D R_D. \quad (5.16)$$

Then using the transconductance relation (5.1),

$$i_D = g_m v_{GS} = g_m v_{in}, \quad (5.17)$$

we have

$$v_{out} = -g_m R_D v_{in}, \quad (5.18)$$

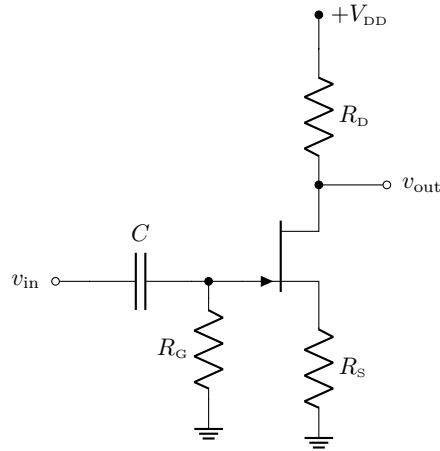
which implies a voltage gain

$$G = -g_m R_D. \quad (5.19)$$

(voltage gain, grounded-source amplifier)

So for example, if  $R_D = 1 \text{ k}\Omega$  and for the 2N5485,  $1/g_m = 200 \Omega$ , then  $G = -5$ , which is not a huge gain. By contrast, in the common-emitter amplifier, using the intrinsic emitter resistance  $r_e = 25 \Omega$  at  $I_C = 1 \text{ mA}$ , the gain is about  $8\times$  larger (because  $r_e$  is eight times smaller than  $1/g_m$ ). So the advantage of the JFET is high input impedance (good for input stages of amplifiers, especially op-amps), the advantage of the BJT is better speed (due to lower capacitance) and amplification.

The idea is basically the same if we also add a source resistor (which could be useful for setting the dc bias, as we noted above), as shown below.



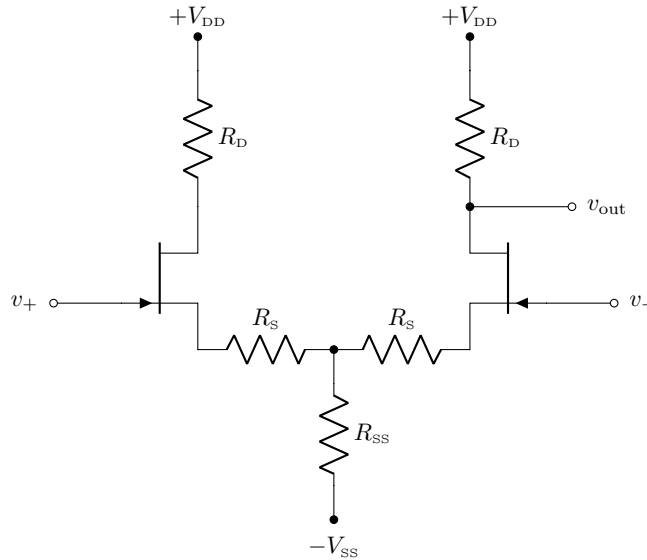
The analysis is a combination of the analysis of the grounded-source amplifier above and that of the source follower in the previous section. The result is

$$G = -\frac{g_m R_D}{1 + g_m R_S} = -\frac{R_D}{R_S + g_m^{-1}}. \quad (\text{voltage gain, common-source amplifier}) \quad (5.20)$$

The last expression here makes the parallel to the BJT more clear—this amplifier is closely analogous to the BJT common-emitter amplifier [see Eq. (4.82)], where  $1/g_m$  plays the role of the intrinsic emitter resistance  $r_e$ . Like the common-emitter case,  $R_S$  acts as a “source-degeneration” resistor [see Section 4.12], reducing the gain, and it dominates the effect of  $g_m$  if it is large enough, so that the gain simplifies to  $G = -R_D/R_S$ .

#### 5.4.4 JFET Differential Amplifier

A straightforward generalization of the BJT differential amplifier of Section 4.10 is its JFET counterpart. The advantage is high input impedance, while the main disadvantage is difficulty in designing and tuning the circuit due to the intrinsic uncertainty in the JFET device parameters.



The analysis here is closely analogous to that of the BJT case, so we will only sketch the idea. Assuming a purely differential (ac) input  $\Delta v := v_+ - v_-$ , the voltage drop across  $R_{SS}$  is fixed, so the two transistors act

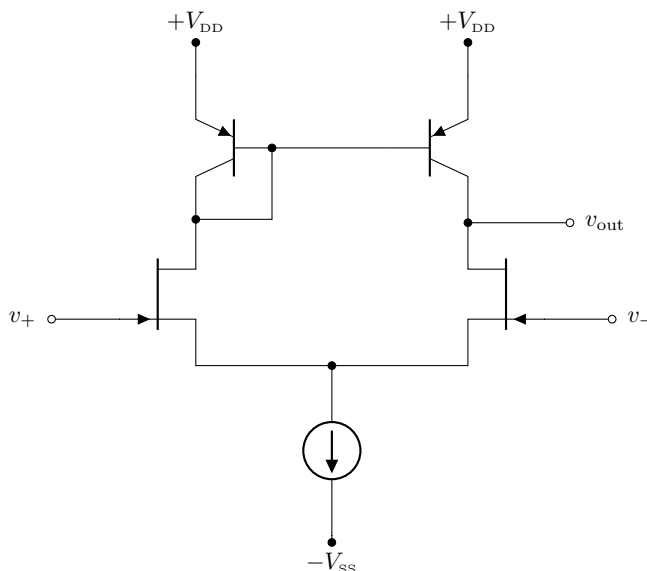
as uncoupled common-source amplifiers. Thus, Eq. (5.20) implies a differential gain

$$G_{\text{diff}} = -\frac{R_D}{2(R_S + g_m^{-1})}. \quad (\text{voltage gain, common-source amplifier}) \quad (5.21)$$

Similarly, a common-mode input  $\bar{v} := (v_+ - v_-)/2$  implies a voltage change across  $R_{SS}$  by the same amount, leading to a current change  $i_{SS} = \bar{v}/R_{SS}$ , which is split between the two transistors. This causes the current  $R_D$  to rise, dropping the output, but also tempers this by reducing  $v_{GS}$ . Accounting for all this leads to a common-mode gain of  $G_{CM} = -R_D/(R_S + g_m^{-1} + 2R_{SS})$ , and thus a CMRR

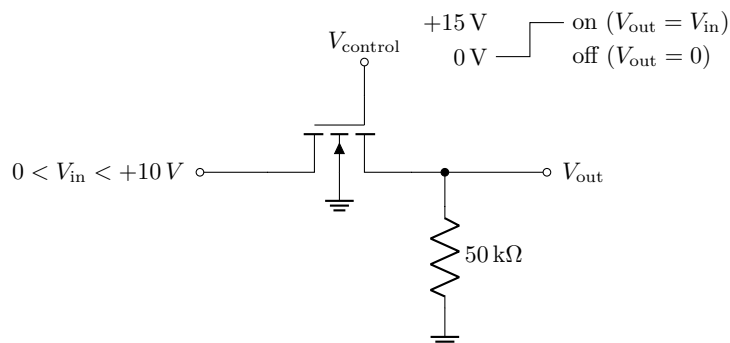
$$\text{CMRR} = \frac{R_S + g_m^{-1} + 2R_{SS}}{2(R_S + g_m^{-1})} \approx g_m R_{SS}. \quad (\text{common-mode rejection ratio, JFET differential amplifier}) \quad (5.22)$$

The last expression here applies to a “good” differential amplifier, which has large  $R_{SS}$  to maximize the CMRR, and  $R_S = 0$  to maximize the differential gain. Thus, the same improvements as in Sections 4.10.4 and 4.11.4.1 of replacing  $R_{SS}$  by a current source and replacing the  $R_D$  pair by a (BJT) current mirror apply here. See the schematic below for a generic example of an “improved” configuration with large differential gain and CMRR.



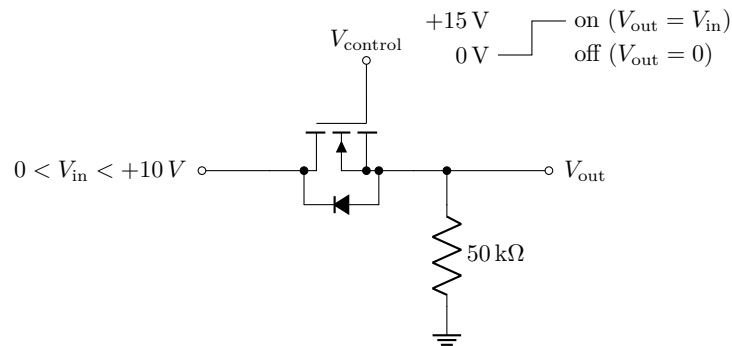
### 5.4.5 MOSFET Analog Switch

One nice example of a MOSFET application is as an **analog switch**, which passes or blocks an analog signal based on a control voltage.



Here, the base terminal is tied to ground; the MOSFET acts as a short if  $V_{\text{control}}$  is well above any input, but acts as an open circuit if the gate drops to zero. The output resistor ensures a zero (not floating) output when the MOSFET is off, and it suppresses any tendency of the signal to cross the transistor in the off state due to drain–source capacitance. This circuit can be adapted to positive/negative signals by dropping the base and the “off” voltage to  $-15\text{ V}$ .

MOSFETs with a separate base connection (like the obsolete 2N4351) are actually the exception, rather than the rule. More common are MOSFETs with the base and source shorted together. The net effect of this connection is similar to connecting a diode from source to drain. The circuit still works with this kind of MOSFET (e.g., 2N7000), if the orientation is correct.



However, for a positive/negative signal, *two* 2N7000's, back-to-back, are needed for the switch to function properly. Note also that, in the case of the 2N7000, the characteristic curves on p. 145 show that having the 0V control signal at least a few V over the maximum possible input voltage is enough to ensure that the MOSFET stays saturated.

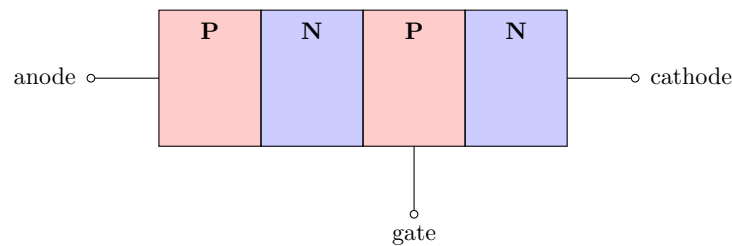
Of course, more sophisticated (and prepackaged) versions of switches are available. These are briefly discussed later in Section 12.5.

## 5.5 Thyristors

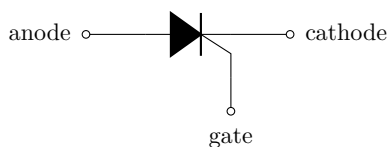
Now we will move on to a few semiconductor components that are somewhat more complicated than the basic transistors (BJT and FET). **Thyristors** refer to a family of devices based on a three semiconductor junctions, though the name “thyristor” can also be used synonymously with the basic SCR component that we will treat first. These are devices that utilize p-n junctions to produce hysteresis, which is useful in switching circuits.

### 5.5.1 SCRs

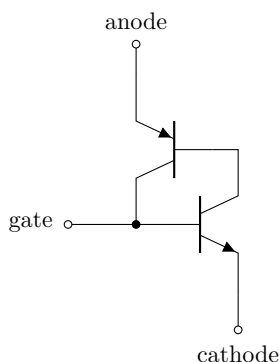
The basic form of a thyristor is the **semiconductor-controlled rectifier**, or **SCR**. An SCR is constructed from alternating four regions. The result is something like a diode, but with two p-n junctions and one n-p junction, all in series.



The basic operation will be something like a diode, but there is one additional “gate” connection to the internal p-type region. This is reflected in the schematic symbol, shown below.



Thinking of the p-n junctions as diode-type junctions is only useful in understanding the “OFF” state; in either direction there is a reverse-biased diode, so no current can flow. However, to understand the switching behavior of the SCR, it helps to think of the p-n-p stack to the left as a PNP transistor, and the n-p-n stack to the right as an NPN transistor. Since the inner two regions are common to each transistor, the transistors are interconnected as shown in the equivalent circuit below.



Again, if no current is initially flowing in the circuit, then both transistors are OFF, and no current can flow from anode to cathode, because the current would have to pass via the collector of one or the other transistor.

However, if a (sufficiently large) anode–cathode current is already flowing in the circuit, then the collector current of the NPN transistor ensures the PNP transistor is ON (since the current comes via the emitter–base junction), and the collector current of the PNP transistor flows through the base–emitter junction of the NPN transistor, ensuring that it is ON.

As an aside, note that although this transistor model is a useful way to understand how an SCR works, an SCR can’t necessarily be replaced by two transistors. For example, the collector leakage current of one of the transistors could supply sufficient base current to turn the other ON, in which case this SCR OFF state wouldn’t work as described above.

So now we see that there are two possible states, where current is flowing or not, and the states are self-sustaining (the current-flowing state of course requiring a sufficient forward voltage). Then the question is, how can we transition between the states? In particular, if no current is flowing, how can we get the current started? There are two useful answers to this in practice:

1. Inject some gate current (flowing from gate to cathode). This turns on the NPN transistor, which turns on the PNP transistor, and the SCR transitions to the conducting state.
2. Bring the forward voltage sufficiently high that the reverse-biased junction breaks down. Just before the breakdown, the forward voltage can be quite high (hundreds of volts), but after breakdown, the forward voltage will become much smaller, on the order of a volt. Note that a gate current will reduce the breakdown voltage, and a gate current above a switching threshold brings the breakdown voltage to essentially zero.

To make a transition in the other direction, or that is, to stop the conduction, the only way is to bring the forward current to zero (or equivalently, drop the forward voltage to zero). This “resets” the SCR into the non-conducting state, assuming there is no gate current is off.

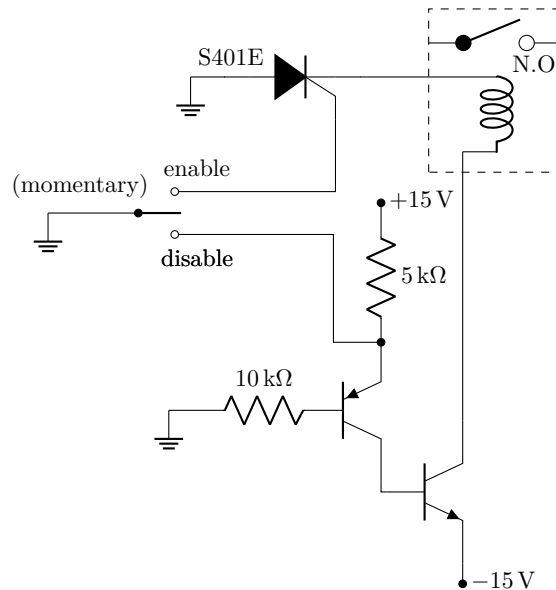
As an example, the small-signal S401E SCR has 1-A maximum (rms) current, 400-V off voltage (in either direction), a gate trigger current guaranteed to be in the range of 1–10 mA, a 1.6-V maximum forward-conduction voltage, and a maximum gate forward voltage of 1.5 V.

SCRs are useful in switching circuits where some “memory” (latching behavior) is required, and they are also useful in switching large currents and voltages. The key to the utility in switching high-power loads is that the forward voltage drop is small in the conducting state, so large currents do not necessarily result in large power dissipation. It is also possible to break down the SCR in the reverse direction (just like a diode) with a sufficiently high voltage, but then the voltage gets clamped to the (high) breakdown voltage, as in a zener diode. In that case, a large current could result in a large power dissipation, because both  $V$  and  $I$  are large.

The name “thyristor” derives from the tube analogue of the SCR, the **thyatron**, which is a gas-filled tube where a control voltage can initiate an ionization breakdown of the gas, putting the tube into a conducting state (like the neon lamp, discussed below with DIACs). The thyristor name is then a combination of thyatron + transistor.

### 5.5.1.1 Example: Latching Switch with Power-Supply Fault Protection

One example of a useful switching circuit is shown below. The idea here is that sometimes one has to switch a delicate and/or expensive load, and the driving circuitry must include some protection against certain faults, such as a failed power supply. (Power supplies can fail due to faulty capacitors, and capacitors elsewhere in the circuit, such as the bypass capacitors discussed in Section 7.6.2.2, can fail to a short and effectively cause a power-supply failure.) The circuit below is intended to switch a delicate load via a relay (magnet-controlled switch). The diagram shows the relay, but not the load, and the relay is normally open (N.O.) when there is no magnet current.

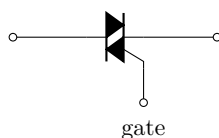


A momentary switch allows the user to toggle the relay on or off. In the “enable” position, the switch activates the SCR, but only if the  $\pm 15$  V power supplies are working, so that there is a path to  $-15$  V via the transistor pair. Once the SCR is activated, current flows through the relay coil, connecting the load. The “disable” switch position forces the PNP emitter to ground, which turns off the PNP and thus the NPN transistor opening the relay and resetting the SCR. (The disable operation also emulates the failure of the  $+15$ -V power supply.)

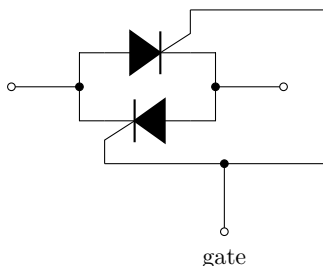
## 5.5.2 DIACs and TRIACs

There are two related devices in the thyristor family: **DIACs** and **TRIACs**. Since SCRs are diode-like in their asymmetric operation, it is useful to have analogous devices to work in ac circuits. The TRIAC is basically an SCR that can conduct (with low forward voltage) in either direction. The TRIAC symbol is shown below.





This component is basically equivalent to two SCRs, oppositely oriented and wired in parallel, with a common gate, as shown below.



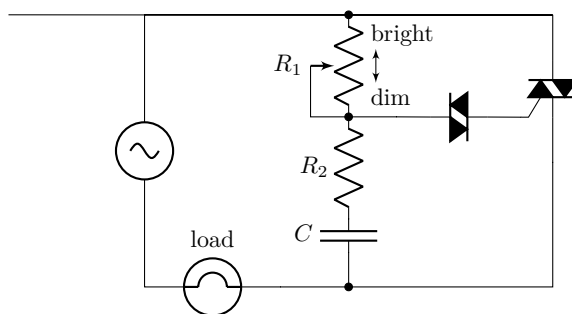
The idea is similar with a DIAC, but this component has no gate; the intent is for this device to conduct by bringing the voltage above threshold (in either direction), to break down the device into conduction.



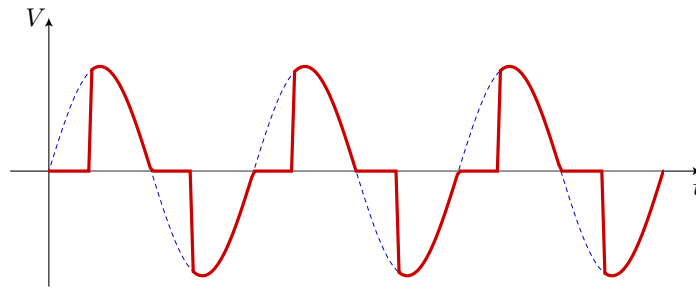
While it's useful to think of both devices as a combination of two SCRs, they are fabricated in a simpler way as a single stack of junctions. For example, the DIAC, because it doesn't need a gate connection, can be fabricated as a p-n-p stack. The symmetry of the stack allows for symmetric-breakdown operation. DIACs are typically fabricated with breakdown voltages in the range of 32–40 V. Gas-discharge lamps can also work in a similar way to DIACs. For example, the small NE-2 neon lamp has a breakdown voltage of 90 V (i.e., it does not conduct until the voltage exceeds 90 V in either direction), and the discharge-state (i.e., lit) voltage drops to around 60 V when the lamp conducts. The nominal conduction current is about 0.5 mA, so if used as a lamp across 120 V ac, a dropping resistor of around 100 k $\Omega$  is necessary to limit the lamp current. However, for switching applications, a neon lamp can function in place of a DIAC, for example in the dimmer circuit to follow.

### 5.5.2.1 Light Dimmer

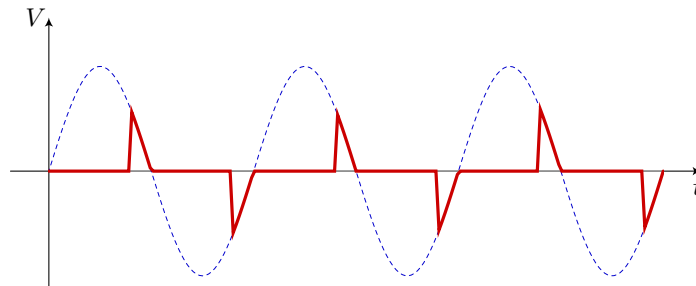
One of the most important applications of thyristor devices is in dimmers for ac lighting circuits. There are a number of designs, but a common DIAC–TRIAC design is shown below.



As each “pulse” (of either sign) from the ac source begins, the TRIAC is initially off, and the voltage begins to charge the capacitor via  $R_1 + R_2$ . Once the capacitor charges to the breakdown voltage of the DIAC (around 30 V), it conducts, and discharges the capacitor via  $R_2$  through the TRIAC gate, causing it to conduct. The TRIAC continues to conduct until the end of the pulse, supplying current to the load during the last part of the cycle, and the capacitor also discharges and resets before the next voltage pulse. The net result is that the load is only on for part of each cycle; the voltage presented to the load is illustrated below.



The brightness can be adjusted by changing  $R_1$ : a larger  $R_1$  leads to a longer charging time, which leads to a longer delay until the load turns on each cycle. A resulting dimmer signal is illustrated in the plot below.



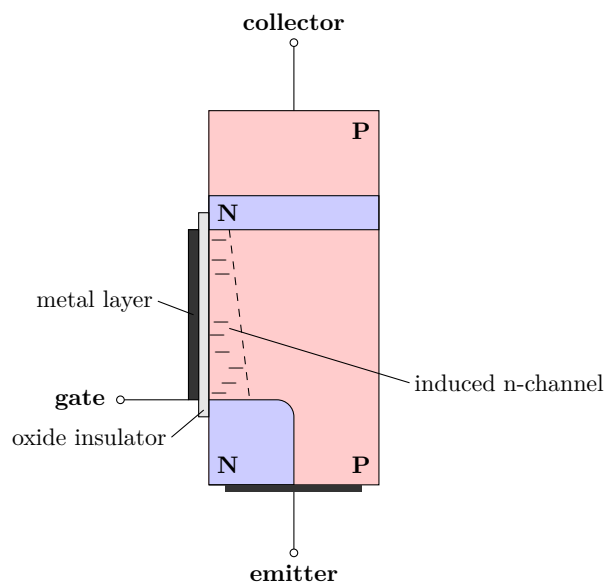
The  $R_2$  resistor protects against excessive capacitor charging or discharging currents.

One disadvantage of this circuit is the sudden turn-on edge in each cycle (in the voltage as well as the current), which is largest at 50% brightness. This can be a source of electromagnetic interference, and it can also cause buzzing in the load (for example, causing mechanical motion of the filaments in incandescent bulbs). Good dimmer designs include inductors in series with the load to suppress the interference and smooth the hard edges.

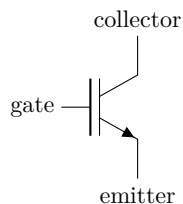
This circuit must also be designed with a particular load in mind (or range of loads), because the load must draw sufficient current to keep the TRIAC on until it resets at the end of the cycle. Thus, for example, older thyristor dimmers designed for incandescent loads can perform poorly when the lights are replaced by lower-current LED equivalents, which draw much less current. The resulting problems typically come in the form of flicker in the new lighting. The other problem that can crop up in this situation is that during the charging cycle, a small current flows through the load from the voltage source. Generally  $R_2$  limits this to something small enough that an incandescent light is effectively off. But with higher-efficiency LED lights in an older dimming fixture, the LEDs can stay dimly lit, even in the dimmest setting (or even the “OFF” state of the fixture).

## 5.6 IGBTs: Switching Very Large Voltages and Currents

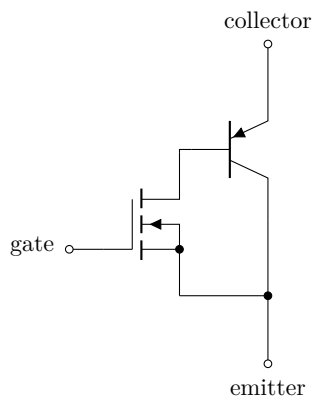
The **insulated-gate bipolar transistor (IGBT)** is a device fabricated as shown below. It is something like the n-channel (enhancement-mode) MOSFET (see the diagram on p. 141), but with an extra p-type region, and different names for the conduction terminals.



The schematic symbol is shown below: it is basically the symbol for a BJT, but with a gate separation to emphasize the insulated gate. Note that while MOSFETs and BJTs come in n- and p-flavors, the IGBT comes in only the polarity shown (or, that is, the opposite polarity is so rare as to be effectively nonexistent).



Ignoring the upper p-type region, this is just a MOSFET with source and base shorted together. The extra p-type region then forms something like an PNP BJT with the upper n-type and lower p-type regions. The MOSFET drain and BJT base are connected (because they are formed by the same n-type region), and the MOSFET body/source and BJT emitter are connected. Thus, an equivalent circuit is the transistor pair as connected below.



This circuit functions more or less as a BJT, except the C–E conduction is controlled by a voltage ( $V_{GE}$ ), rather than a base–emitter current. When  $V_{GE}$  is sufficiently positive, the MOSFET turns ON, which switches on the PNP transistor and allows current to flow from collector to emitter. Note that in principle this current may flow via either the PNP collector (which is *not* the IGBT collector!) or via the MOSFET drain–source

path, but the device is fabricated such that most current flows through the collector. For large currents, this arrangement can be efficient compared to a power BJT, which may require substantial base current to control the transistor (around 10% of the collector current for good saturation, since power BJTs are on the low end of the  $\beta$  spectrum).

For high-power switching applications, power MOSFETs are also available. However, for switching loads with hundreds of amps and hundreds of volts, IGBTs tend to be the best choice in terms of power dissipation. Recall that if the gate voltage is sufficiently large, the MOSFET behaves like a variable resistor. For example, the Infineon IRF1405 power MOSFET is specified to have a typical ON resistance of  $4.6\text{ m}\Omega$ , a maximum continuous drain current of  $169\text{ A}$  (at  $25^\circ\text{C}$ ), and a drain–source breakdown voltage of at least  $55\text{ V}$ .<sup>6</sup> In this case, when the MOSFET switches the power to a load, the power dissipated by the MOSFET is  $P = I^2 R_{\text{DS}}$  while the transistor is ON. Since the conduction path for an IGBT behaves like a BJT, the collector–emitter voltage is relatively small when the transistor is saturated, and the IGBT ON power dissipated is  $P = IV_{\text{CE}}$ . A comparably sized example of an IGBT (in the sense of being available in a package of similar size) is the Fairchild FGA60N65SMD,<sup>7</sup> which is rated to switch  $120\text{ A}$  (at  $25^\circ\text{C}$ ) at  $650\text{ V}$ , with  $V_{\text{CE}} = 1.9\text{ V}$  typical.

Since the MOSFET power scales quadratically with current while the IGBT scales only linearly, the IGBT can have a big advantage for high-power loads. For the two example devices above, we can estimate a power dissipation of  $228\text{ W}$  for the IGBT vs.  $66\text{ W}$  for the MOSFET at  $120\text{ A}$ . For these “small” currents, the MOSFET is the clear winner (though the IGBT can of course switch much higher voltages). However, the crossover to where the IGBT has the advantage occurs at  $I^2 R_{\text{DS}} = IV_{\text{CE}}$ , or  $I = V_{\text{CE}}/R_{\text{DS}}$ , which comes out to about  $400\text{ A}$ , if we extrapolate assuming the same device parameters carry over to larger devices. This is a typical figure, above which IGBTs are preferred, especially to stand off large voltages. Of course, if MOSFETs are preferable for some reason, they can be used in a parallel transistor bank to reduce the effective resistance. However, in MOSFETs capable of standing off larger voltages, the conduction path must be longer, typically leading to larger ON  $R_{\text{DS}}$ , and thus worse power dissipation at a particular current.

Such large IGBTs certainly exist, in “mini-brick” and “brick” packages, with screw terminals. For example, the somewhat larger IXYN80N90C3H1 IGBT by IXYS<sup>8</sup> comes in a SOT-227B “miniBLOC” package, with 4 screw terminals. (Two terminals are emitter terminals, so one can function as a “Kelvin emitter.” The second terminal can connect to a voltage-sensing circuit, such as a gate driver, to accurately sample the emitter voltage without a spurious drop across a high-current contact resistance at the other emitter connection to the current load.) This is designed to switch a pulsed current of  $340\text{ A}$  (or  $115\text{ A}$  continuous) at  $900\text{ V}$ , with  $V_{\text{CE}} = 2.3\text{ V}$ . The more expensive ( $\sim \$200$ ) CM600HA-24A by Powerex<sup>9</sup> is rated for  $600\text{ A}$  (dc,  $1200\text{ A}$  pulsed) at  $1200\text{ V}$ , with  $V_{\text{CE}} = 3\text{ V}$ .

IGBTs have a reputation for being somewhat slower than MOSFETs for switching, because the main conduction path is governed by the carriers injected into the “base” n-type region. However, the carriers enter and leave this region only indirectly, via the neighboring p-type regions. This is in contrast to a standard BJT, where carriers can enter or leave directly via the base connection. Especially when switching off, excess carriers can be “stuck” in the base region, only escaping by causing collector–emitter current to flow.<sup>10</sup> On the other hand, in the FET–IGBT comparison above, the FGA60N65SMD IGBT happens to have a rise/fall time of around  $50\text{ ns}$ , while the IRF1405 MOSFET has a rise/fall time of  $190/110\text{ ns}$ , respectively.

### 5.6.1 Driver Circuitry

In both power MOSFET and IGBT switching circuits, the transistor is only thermally well-behaved if the device is fully OFF (zero current, large voltage) or fully ON (large current, small resistance/voltage). In between, the transistor may drop a substantial voltage while conducting a large current, which can be a

<sup>6</sup><http://www.infineon.com/dgdl/irf1405pbf.pdf?fileId=5546d462533600a4015355db084a18bb>

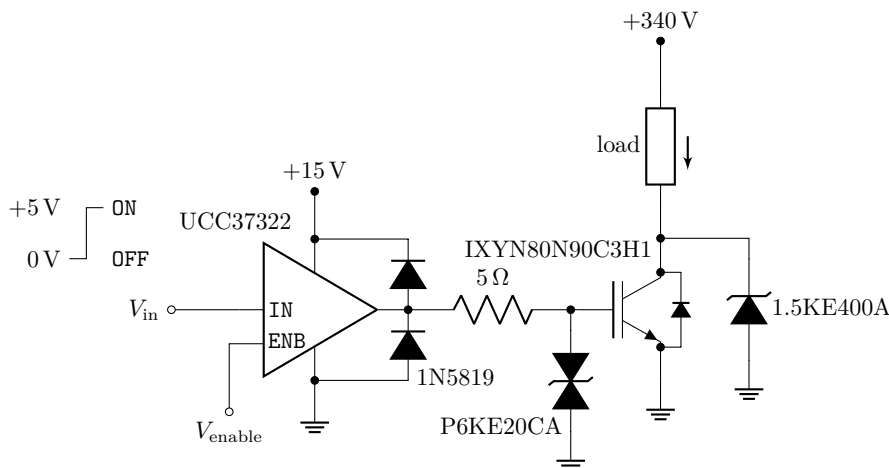
<sup>7</sup><https://www.fairchildsemi.com/datasheets/FG/FGA60N65SMD.pdf>

<sup>8</sup>[http://ixapps.ixys.com/DataSheet/DS100522A\(IXYN80N90C3H1\).pdf](http://ixapps.ixys.com/DataSheet/DS100522A(IXYN80N90C3H1).pdf)

<sup>9</sup>[http://www.pwr.com/pwr/docs/cm600ha\\_24a.pdf](http://www.pwr.com/pwr/docs/cm600ha_24a.pdf)

<sup>10</sup>International Rectifier, *Application Note AN-983: IGBT Characteristics*, available online at <http://www.infineon.com/dgdl/an-983.pdf?fileId=5546d462533600a40153559f8d921224>

disaster if it carries on for long. For this reason, it is a huge advantage to switch between the ON and OFF states as quickly as possible. However, the gate acts as a capacitive load to the input (5.5 nF for the IRF1405, 2.9 nF for the FGA60N65SMD, 4.6 nF for the IXYN80N90C3H1), and to change the gate voltage rapidly requires a large, transient current. Thus, such high-power devices typically require special gate-driver circuitry, typically consisting of a pair of smaller FETs or BJTs acting as a fast push-pull buffer. These parts are conveniently packaged in some integrated gate-driver IC's; one example with a relatively large current capability is the UCC37321/UCC37322 by TI,<sup>11</sup> which can drive a 9-A peak current, with the output switching between ground and a maximum power supply voltage of +15 V (the UCC37321 inverts the input signal in the sense of a logic inversion, as in Section 9.3.1, while the UCC37322 performs no inversion). An driver example driving up to 340 A (pulsed) current through a load at high voltage via a power IGBT is shown below.



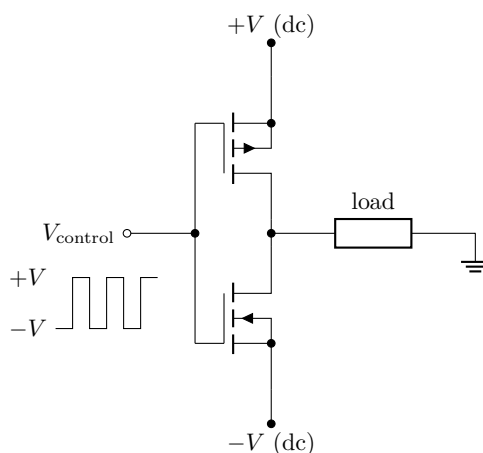
The driver takes a TTL-compatible input (see Section 11.4), but the ON level can also be up to +15 V. The driver also has an ENABLE input, which forces the output LOW, close to 0 V (this is also true on the inverting UCC37321), independent of the input voltage.

A number of other components ensure proper operation and protection of the switching transistor. First, the 5-Ω gate resistor limits the gate current to  $15\text{ V}/5\text{ Ω} = 3\text{ A}$  in the worst case of a short to ground, ignoring the effects of capacitance, etc. The data sheet implies a minimum value of 2 Ω for this resistance, and gives specifications for up to 15 Ω. This resistance limits the current into the gate capacitance, easing demands on the driver and helping to control against voltage overshoot. The 1N5819 Schottky diodes on the driver output clamp the output voltage to the power-supply range, preventing output overshoot; the Schottky diodes are fast and have a (relatively) small forward voltage. There is also a transient-voltage suppressor (TVS) on the IGBT gate, which is basically two back-to-back Zener diodes. The P6KE20CA TVS break down if the voltage exceeds 20 V in either direction (and is rated for 600 W), protecting against gate overvoltages (which can come from capacitive coupling via the IGBT as well as from the driver). Finally, as with many IGBTs, the IXYN80N90C3H1 incorporates an internal reverse C-E diode, so that it only stands off voltage in one direction—the device is not designed to stand off high reverse voltages, so the diode offers protection in the event of (intentional or unintentional) reverse transient voltages. However, it is good practice to include another (unidirectional) TVS to protect against reverse and overvoltage transients. The 1.5KE400A TVS acts as a fast Zener diode that breaks down at 400 V and can handle 1.5 kW.

## 5.6.2 Inverter Circuits

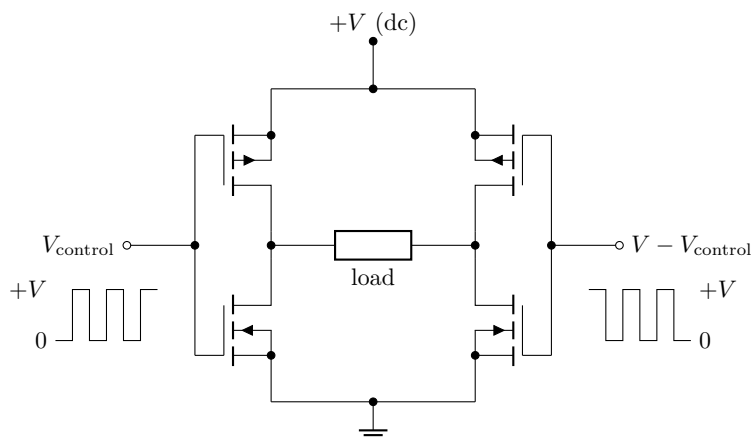
One important application of power MOSFET and IGBT switching circuits is the **inverter**, which is used to convert dc to ac signals. An example of a **half-bridge inverter** is shown below, where a push-pull pair of MOSFETs switches one side of a load between  $\pm V$ , giving a square-wave ac signal of amplitude  $V$ .

<sup>11</sup><http://www.ti.com/lit/ds/symlink/ucc27322.pdf>



Circuits of this type are particularly common in battery-powered applications. For example, an uninterruptible power supply (UPS) creates an ac waveform from dc battery power, after a dc–dc converter boosts the voltage to 120 V (for U.S. supplies). Cheaper UPSs make a crude approximation of a sine wave by switching to +120 V for a quarter-cycle and –120 V for a quarter-cycle, and to OFF for the other quarter cycles in between; this gives both the correct peak and rms voltages, compared to the sine wave. More sophisticated UPSs use rapid switching, pulse-width modulation, and filtering to create a clean sine wave. Other important applications of inverter circuits is in creating ac power from solar cells, motor-driving circuits in battery-powered electric cars, induction “burners” in kitchen stoves, and arc welders. Switching power supplies also use inverters to convert rectified mains voltage into a high-frequency (tens of kHz) signal that can be changed to other dc voltages by a transformer and another rectifier and filter circuit; because of the high frequency, the transformer can be compact and inexpensive compared to traditional 60-Hz transformers. Compared to linear power supplies, they are noisier but cheaper and more compact, and are preferred, for example, in computer power supplies requiring high current.

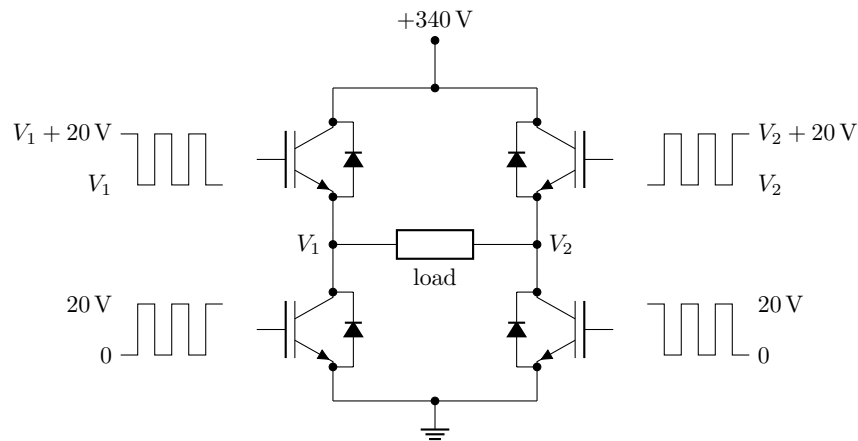
A circuit that makes better use of the available voltage, at the expense of more parts and a possibly ungrounded load is the **full bridge** or **H-bridge** (after the shape made by the switching transistors and the load), as shown below.



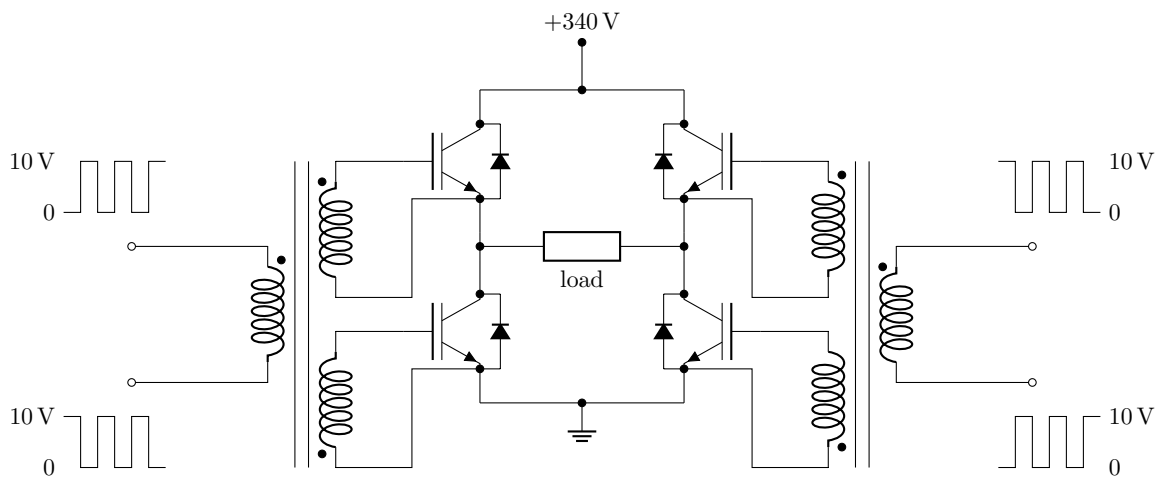
This circuit is basically two MOSFET half-bridges switched out of phase. As in the half-bridge, the control voltage swings between the power-supply rails; the pairing of complementary n-channel and p-channel devices makes the control relatively simple. In each state of  $V_{\text{control}}$ , a diagonally opposed pair of transistors conducts, while the other pair is open. The result is that the voltage across the load alternates between  $\pm V$ , but the  $-V$  power supply is no longer necessary, as in the half-bridge example above.

For high voltages and high currents, again, IGBTs are preferable. However, IGBTs are generally not available in complementary pairs, which makes the gate control considerably more complicated. An example

H-bridge is shown below.



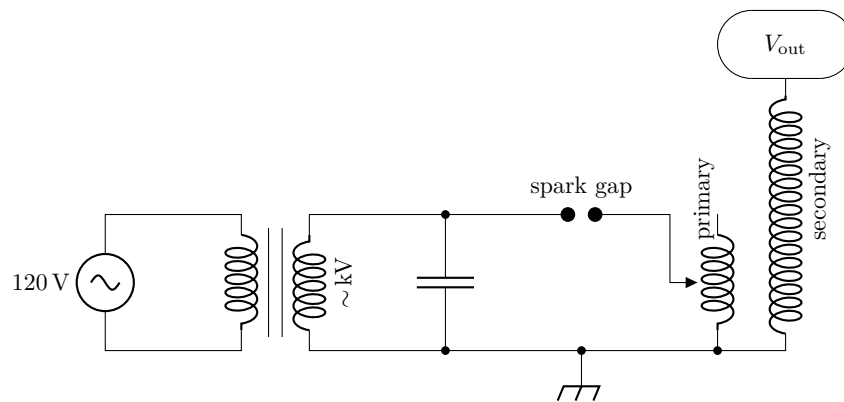
The voltages shown are appropriate for the IXYN80N90C3H1 IGBT, for example. Note that the IGBT gate voltage is always measured with respect to the collector. The lower two IGBT's can thus be switched normally, except the control voltage need not swing all the way to the upper supply rail. The tricky part, however, is that the *upper* two IGBTs require similar control voltages, but referenced to their emitter voltages, which change according to the phase of the inverter. Thus, the control circuitry for these two IGBTs must float along with the appropriate load voltage. One solution to this is to use special driver ICs with optocouplers that isolate the input and allow the output voltage to float with respect to the input. Another solution is to couple the gates via gate-drive transformers, as illustrated below.



In this case, two gate driver ICs drive the ends of the primary coil of a transformer with out-of phase signals, making the primary voltage oscillate between  $\pm 20$  V. The same drivers can drive the primaries of both transformers, but then the primaries should be connected in antiparallel. The two antiparallel secondaries (with a 1:1 primary:secondary ratio here) drive the gates, and there is no problem floating the gate voltages. Note that this only works for relatively rapid switching and a symmetric drive waveform, as the transformers cannot support a dc signal.

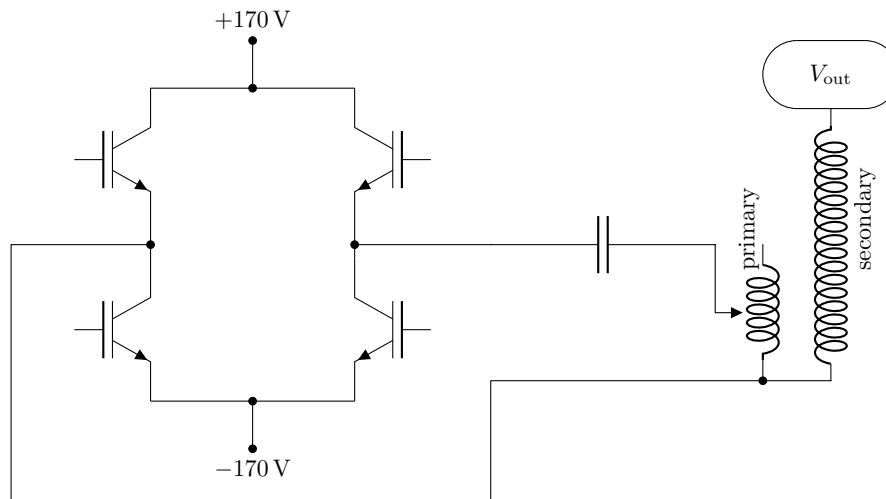
### 5.6.2.1 Inverter-Based Tesla Coil

Let's return to the Tesla-coil example of Section 2.8.1, shown again below.



Again, the LC tank circuit resonates at the same frequency as the secondary, and it acts as a harmonic oscillator. The tank circuit is ultimately powered by the 60-Hz transformer output, but because the resonant frequency is much higher ( $\sim 100$  kHz), we need a way to “whack” the tank circuit to get it oscillating. This is a useful way to look at one function of the spark gap: it holds off the power-supply voltage until it builds up, and then it suddenly lets it into the LC circuit, exciting it with a step voltage. Once the spark gap breaks down, it acts as a short circuit, connecting the LC loop; since the resonance frequency of the LC circuit is high, the oscillating signal bypasses the secondary of the high-inductance, 60-Hz transformer.

A more recent variation on the classic Tesla coil involves replacing the spark gap with an inverter circuit.<sup>12</sup>



In this example, a dc power supply of  $\pm 170$  V (derived by rectifying and filtering 120-V mains voltage) is switched by a full-wave bridge to give the same, step-function excitation to the tank circuit. This design can be more consistent in its performance, because it does not rely on a discharge to excite the resonator. However, the overall circuitry is generally much more complicated, involving the driver and protection circuitry discussed above, as well as yet more control circuitry. For example, to not overload the IGBTs, the bridge should only switch when the current in the tank circuit crosses through zero; otherwise, the sudden change in current would lead to a large inductive-kick voltage, which can strain the IGBTs. So a current detector is needed to feed back to the inverter inputs so that switching occurs at the correct time. Inverter-type coils are also typically modulated, being active only a small fraction of the time, so that the high currents involved do not thermally overload the IGBTs. This requires more control circuitry which must also turn off the inverter bridge only at the proper times.

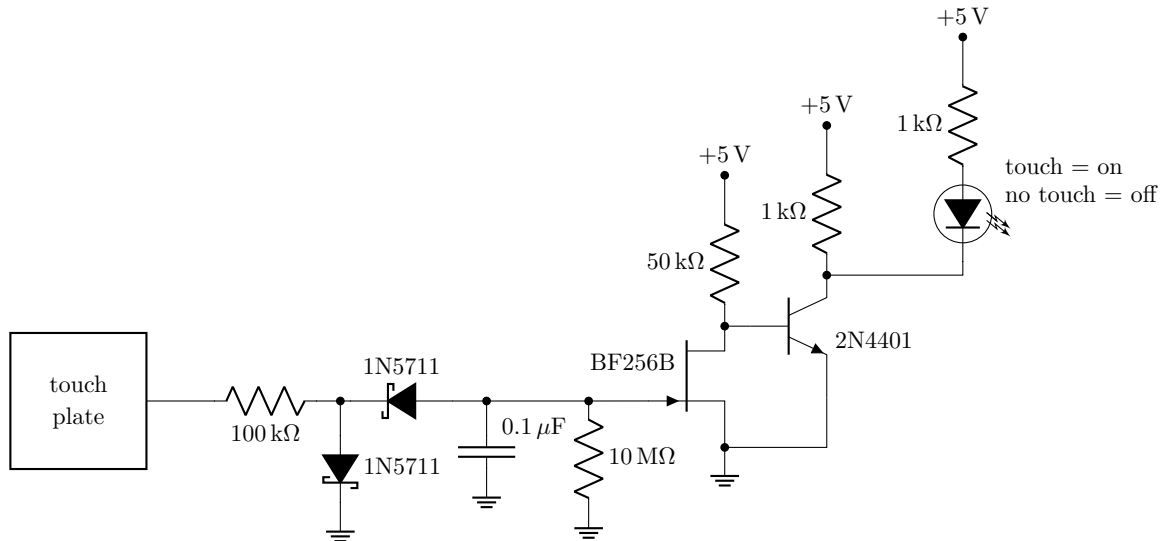
<sup>12</sup>See, for example, <http://www.stevehv.4hv.org>.



## 5.7 Circuit Practice

### 5.7.1 Touch Switch

Starting in the 1980's, table lamps with a “touch base” became popular. The idea is, when you touch a metallic base, the lamp toggles between on and off (or between off, on, and two dim states). The key to the operation of this circuit is that the human body acts as an efficient antenna for 60-Hz radiation, which can be coupled into a metallic plate (or even a “dangling” wire) and into a circuit. The circuit below shows how the coupled radiation can be rectified and used as a switch.



- How does this circuit work?
- What is the function of the 10-M $\Omega$  resistor?
- Note that Schottky diodes have a small forward voltage (about 0.2 V), to make the circuit more sensitive.
- How would you have to modify the circuit to use a MOSFET in place of the JFET, and can you think of an advantage of this circuit over that one?
- To add in some digital electronics, what can you add to the circuit to toggle the LED between ON and OFF on each touch?
- How should this circuit be modified to control a 120-V incandescent lamp?

**Solution.** When nobody is touching the plate, there is essentially no signal at the input; then the FET is ON (the default state), which means the npn base is at low voltage. Thus, the npn transistor is OFF, so the connection to the LED is pulled to +5 V, so the LED is off. (Strictly, the 1-k $\Omega$  resistor at the collector is not necessary, but it allows the circuit to work if the LED is replaced by something else, such as a relay, discussed below.

When touching the plate, 60-Hz radiation is coupled in, and rectified by the diode network (Schottky diodes with small forward voltages) to pass only the negative parts of the signal. The first diode (to ground) discards the positive half of the signal, which shouldn't make it through the second diode anyway; so it seems like this first diode might not really be necessary. However, this first diode gives a path to ground for the positive part of the input signal (so the diode network is basically a charge pump, transferring charge from the capacitor to ground via the diodes). This process charges the capacitor to negative voltage; when this

crosses the threshold ( $-4\text{ V}$  max for this JFET), the JFET turns OFF, which turns ON the npn transistor, and lights the LED.

When the touching stops, the  $10\text{-M}\Omega$  resistor slowly ( $1\text{-s } 1/e$  time) discharges the capacitor to turn the JFET back ON. Empirically, the LED turns off in a fraction of a second with a brief touch of the antenna, but can take about 2 seconds to turn off with a long touch of the antenna.

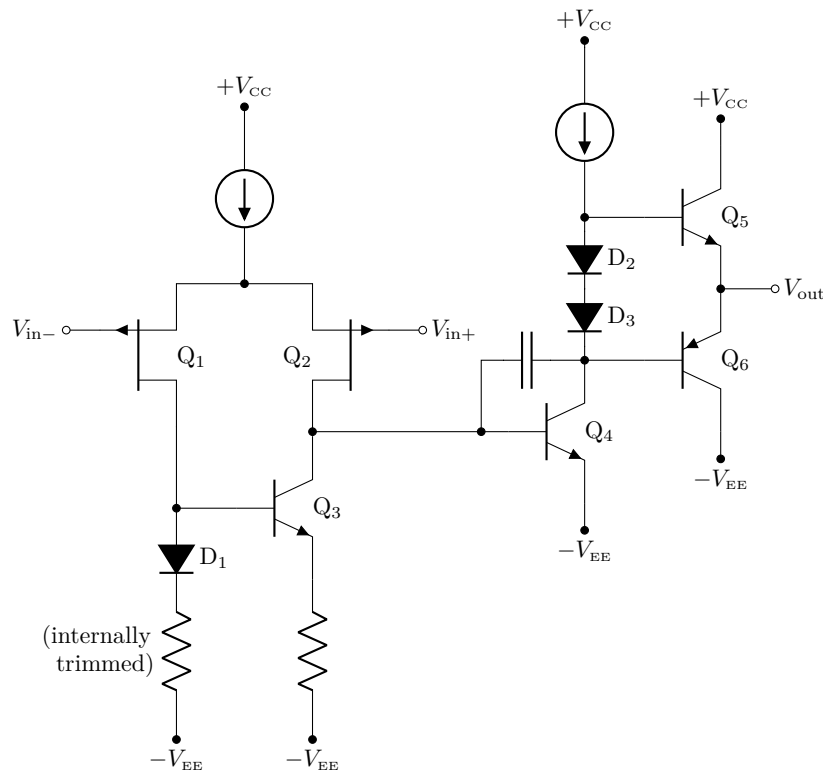
To use a MOSFET, the direction of the diode would have to be reversed, because we'd want to keep the positive parts of the signal to change the MOSFET state. This would also reverse the sense of the LED being on/off, so the LED could be placed between the collector and ground to restore previous operation. Note that MOSFETs are more sensitive to destruction by static discharge and overvoltage. While the  $100\text{-k}\Omega$  input resistor should help in this regard, it would be safer to use a clamping diode at the gate to make sure the MOSFET is protected. Overall, the JFET makes for a simpler design here.

To get the LED to toggle, a divide-by-2 counter (flip-flop) would help here.

To drive a heavy load like a light bulb, the output of this circuit should drive a relay (or even a small relay that drives a large relay, which in turn drives the light bulb). It could also drive a solid-state version of a relay, like a power-MOSFET switch.

### 5.7.2 JFET Op-Amp Internals

The schematic below shows a simplified diagram of the internals of several classic JFET-input op-amps, as given in the data sheets for the TL082,<sup>13</sup> LF411,<sup>14</sup> and the LF353.<sup>15</sup>



As before in Section 4.16.5, the exercise here is to identify what is going on with all the components in the schematic, based on what you know about FETs and BJTs.

**Solution.**  $Q_1$  and  $Q_2$  form a JFET differential amplifier, but with p-channel JFETs. It has the “improvement” of a current source to  $+V_{CC}$ , acting as an effectively large  $R_{SS}$ . The load of the differential amplifier

<sup>13</sup><http://www.ti.com/lit/ds/symlink/tl082-n.pdf>

<sup>14</sup><http://www.ti.com/lit/ds/symlink/lf411-n.pdf>

<sup>15</sup><http://www.ti.com/lit/ds/symlink/lf353-n.pdf>

is a current mirror, here shown as  $Q_3$  and  $D_1$  (which sets  $V_{BE}$  for  $Q_3$ , and is implemented internally as a transistor as in the usual current mirror in Section 4.11.4. Note that all of the *rest* of the transistors are BJT's, to get the high gain desired in an op-amp (recall that JFETs have inherently lower gain compared to BJT's because  $g_m^{-1}$  tends to be a fair bit larger than  $r_e$  in a comparable device).

The two resistors (of the order of a couple  $k\Omega$ ) that appear as a load in addition to the current mirror allows the voltage offset (i.e., the symmetry of the differential amplifier) to be trimmed. But they also improve the behavior of the current mirror. For example, if they drop something on the order of a couple of tenths of a volt, they greatly reduce the sensitivity of the current mirror to imbalances in the  $V_{BE}$ 's of the transistors (due to temperature difference or manufacturing variation). Also according to the Early effect, where  $V_{BE}$  varies with  $V_{CE}$ , the resistors temper this effect, which would otherwise cause a mismatch of the mirrored current due to varying voltage at the output of this stage. This effectively increases the (differential) output resistance of the current mirror.

The differential amplifier drives a grounded-emitter gain stage ( $Q_4$ , with a current source as a large effective collector resistance  $R_C$ ). A compensation capacitor limits the bandwidth of the op-amp at this stage. The diodes bias the push-pull complementary emitter-follower pair ( $Q_5$  and  $Q_6$ ), which boost the current drive at the output.

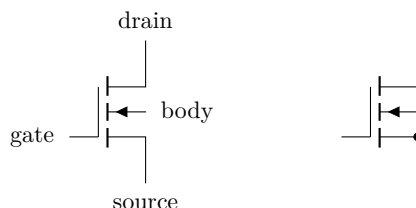
Note that because the JFET inputs act as reverse-biased diodes, they can withstand a wide range of differential input voltages. However, they can't be taken below  $-V_{BE}$ , because this would forward-bias the JFET "diodes," causing them to conduct.

The actual innards of the op-amps are considerably more complicated; detailed diagrams appear in the data sheets referenced above.

## 5.8 Exercises

### Problem 5.1

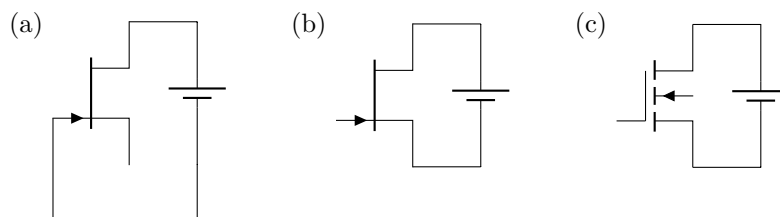
Consider an n-channel (enhancement-mode) MOSFET. The standard connections are shown on the left in the diagram below. Recall that if the G–S voltage is zero, the MOSFET is in the “pinched-off” state, so no D–S current flows in either direction.



Many MOSFETs have the body connected internally to the source (often to fit the MOSFET into a cheaper, 3-pin package), as shown to the right in the diagram. In this case, when  $V_{GS} = 0$ , the MOSFET will block current *in one direction only*, while conducting freely in the other direction. Which direction is which? Explain *why* in terms of the underlying n-type and p-type semiconductor regions.

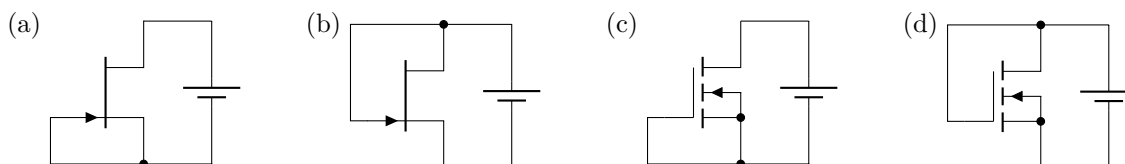
### Problem 5.2

In each circuit, does current flow due to the applied EMF? (Ignore leakage currents, e.g., in a reverse-biased diode, and assume the applied EMF is a few volts.) **Briefly** explain your answer.



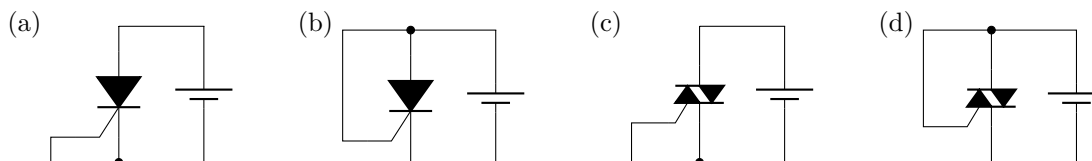
### Problem 5.3

In each circuit, does current flow due to the applied EMF? (Ignore leakage currents, e.g., in a reverse-biased diode, and assume the applied EMF is a few volts.) **Briefly** explain your answer.



### Problem 5.4

In each circuit, does current flow due to the applied EMF? Assume the applied EMF is enough to cause a forward-biased p–n junction to conduct, but **not** enough to break down a reverse-biased p–n junction. Also ignore any leakage currents, and assume the device was **OFF** when the battery was connected. **Briefly** explain your answer, and indicate the paths along which current flows.





## Chapter 6

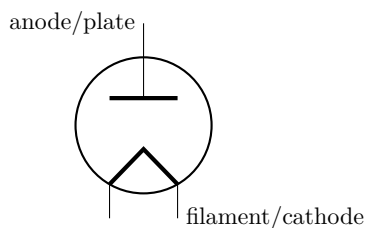
# Vacuum Tubes

Nowadays we take it for granted that there can be billions of transistors in a modern computer. But before transistors there were vacuum tubes, which are much larger, hotter, higher-voltage, and more power-consuming than their silicon successors. Despite the low cost and ubiquity of transistors, vacuum tubes still have a few important applications, from audio amplifiers to circuits with very high power (e.g., for radio-frequency and microwave applications).

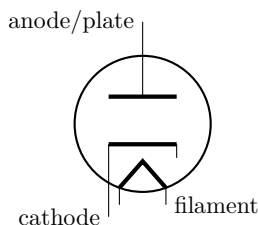
This chapter is located just after FETs because they fit in logically here. However, vacuum tubes remain something of a niche subject. They're worth studying for fun and more circuit practice, but ultimately the op-amps in the next chapter are much more useful: make sure to learn about them first and come back to tubes later if you're still interested.

### 6.1 Vacuum Diodes

The simplest vacuum tube is the **vacuum diode**, so named because of the two fundamental elements, the anode and the cathode. The main function of the cathode is to act as a source of emitted electrons. There are two basic types. The simpler **filamentary cathode** or **directly heated cathode** is simply a filament of wire, heated to high temperature (typically dull glowing red, typically 700°C). The wire is coated so that electrons “boil” easily from the surface at operating temperature. The anode, often called the **plate** is a conductor that attracts or repels the electrons for a positive or negative voltage, respectively, relative to the cathode. Current only flows through the tube in the positive-voltage case where the electrons are drawn to the anode. The symbol for the diode with directly heated cathode is shown below.



The other type of cathode is the **heated cathode indirectly heated cathode**. In this case the cathode is physically and electrically separate from the heating filament, but the cathode still becomes hot and acts as an electron source. In this arrangement the cathode acts as a shield for any ac fields produced by the filament, which is typically heated with a low-voltage ac current. This is an advantage in amplifier vacuum tubes to avoid a filament-induced “hum” in the output, and the vast majority of more sophisticated tubes below have indirectly heated cathodes. However, vacuum diodes in power-supply circuits don't benefit greatly in this respect, so filamentary cathodes are common. The symbol for a diode with an indirectly heated cathode is shown below.



### 6.1.1 Child–Langmuir Law

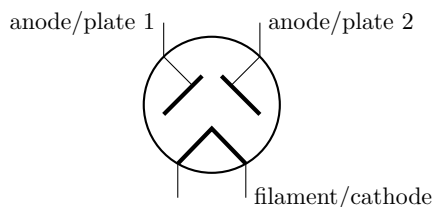
For either cathode type, as the cathode temperature increases, the density of emitted electrons increases, leading to an increased current in forward conduction mode. However this increase does not continue indefinitely: the electrons themselves will modify the electric field between the cathode and anode, and the net effect is that the current is limited to

$$I = \frac{4\epsilon_0 A}{9} \sqrt{\frac{2e}{m_e}} \frac{V^{3/2}}{d^2}. \quad (6.1)$$

This expression, called the **Child–Langmuir Law**,<sup>1</sup> assumes parallel, identical plates for the anode and cathode, where  $A$  is the anode (cathode) area,  $e$  is the electron-charge magnitude,  $m_e$  is the electron mass,  $V$  is the diode forward voltage, and  $d$  is the anode–cathode distance. This scaling of  $I \propto V^{3/2}$  holds in other geometries, but the proportionality factor is geometry-dependent.<sup>2</sup> This law can act as a rough guide to the voltage–current characteristics of a vacuum tube, although it is not accurate enough for quantitative design purposes. However, it is clear from this law that the turn-on of a vacuum tube with voltage is much less “sudden” than the semiconductor diode, which turns on exponentially with forward voltage.

### 6.1.2 Vacuum Full-Wave Rectifier

In semiconductor full-wave rectifier circuits, it is typical to use a full-wave bridge consisting of four diodes (see Section 3.6.2). However, vacuum diodes are more complicated and expensive, so such an arrangement is rare. However, to economize on the number of vacuum tubes, a common variation on the vacuum diode is the **vacuum full-wave rectifier**. This is the same as a vacuum diode, but with two anodes, forming a bridge of two diodes. The schematic symbol is shown below.



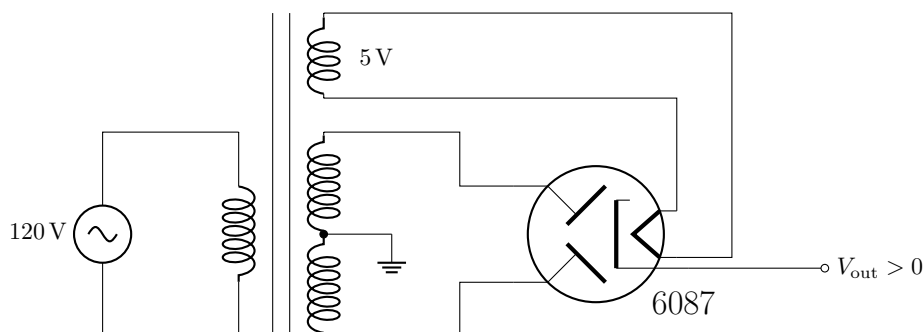
A common, low-power rectifier tube is the 5Y3, a directly heated full-wave rectifier. The part number in this case has a typical numbering that is happily informative: the “5” says the filament operates at 5 V, and the trailing “3” says that there are three elements (two anodes and one cathode). The “Y” is a symbol to differentiate the tube. The part is often listed as, for example, 5Y3GT, where the “GT” indicates a “glass tube” envelope. The 5Y3 operates with 2 A of filament current, can stand off 1400 V of reverse voltage, and can handle 400 mA max forward current (in steady state, rising to 2.5 A max transient current). In typical

<sup>1</sup>C. D. Child, “Discharge From Hot CaO,” *Physical Review (Series I)*, **32**, 492 (1911) (doi: 10.1103/PhysRevSeriesI.32.492); Irving Langmuir, “The Effect of Space Charge and Residual Gases on Thermionic Currents in High Vacuum,” *Physical Review* **2**, 450 (1913) (doi: 10.1103/PhysRev.2.450).

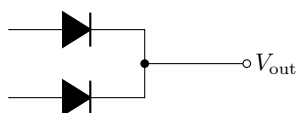
<sup>2</sup>For example, for long, coaxial cylinders, the relation becomes  $I = (8\pi\epsilon_0 L/9) \sqrt{2e/m_e} V^{3/2} / r\beta^2$ , where  $L$  is the cylinder length,  $V$  is the potential at the radius  $r$ , and  $\beta$  is a dimensionless factor of order unity. See Irving Langmuir and Katharine B. Blodgett, “Currents Limited by Space Charge between Coaxial Cylinders,” *Physical Review* **22**, 347 (1923) (doi: 10.1103/PhysRev.22.347).

operation, it would operate at a 360-V (dc) plate voltage, and 125 mA of output current, with a forward voltage drop of 50 V (well above the  $\sim 0.7$  V for a semiconductor diode!). To characterize the reverse leakage, the tube is specified with a reverse resistance of around  $10\text{ M}\Omega$ .

A typical full-wave rectifier circuit using a 6087 rectifier tube (essentially the same as the 5Y3, but with indirect heat) appears below.



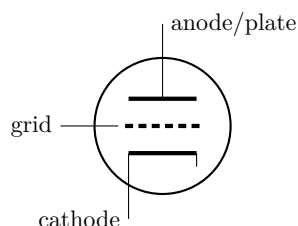
A transformer steps up the mains voltage to several hundred volts (a typical value is 320 V rms, which translates to 453 V peak, not accounting for voltage dropped by the internal resistance of the secondary winding). Again, the tube drops about 50 V, and accounting for the secondary output impedance, the output voltage drops to around 320 V (dc peak). Note that the transformer output is center-tapped, as needed to make a full-wave rectifier with only two diodes. A four-diode full-wave rectifier would only need half the secondary winding, but the extra complexity of the transformer is an acceptable tradeoff to reduce the tube count. The vacuum rectifier again acts as a diode pair, as illustrated below.



The heater in the tube is driven by a separate, low-voltage winding on the same transformer.

## 6.2 Vacuum Triodes

The next tube in the hierarchy, the simplest tube intended for use as an amplifier, is the **vacuum triode**. The name of course indicates that there is now a third element, the **grid** or **control grid**. The grid is something like a fine-mesh electrode between the anode and cathode, but is often made from a fine wire wrapped around supporting posts. The schematic symbol for a triode appears below.



The idea is that the grid can modulate the plate current by attracting or repelling electrons, in the same way as the anode voltage controls current. The grid area is small (i.e., the wire mesh is mostly open), so it is intended only to modulate the plate current, not to conduct any of the electrons itself. In typical operation, the plate-cathode voltage  $V_{PC} > 0$  is fixed, while the grid-cathode voltage  $V_{GC} < 0$  and is relatively small in magnitude. Thus the grid tends to oppose the action of the anode, and a more negative  $V_{GC}$  reduces the plate current. In this sense, the triode grid is analogous to the gate of an n-channel (depletion-mode) JFET, where an increasingly negative voltage pinches off the drain-source current (see Section 5.1). However, the



scale of the voltages is different, and, unlike for the JFET, the “pinch-off” voltage for the triode grid depends on the plate voltage.

In the case of a diode, the Child–Langmuir Law stated that the plate–cathode current and voltage are related by  $I_{PC} \propto V_{PC}^{3/2}$ , where the proportionality constant is geometry-dependent. We can adapt this law to the triode by writing

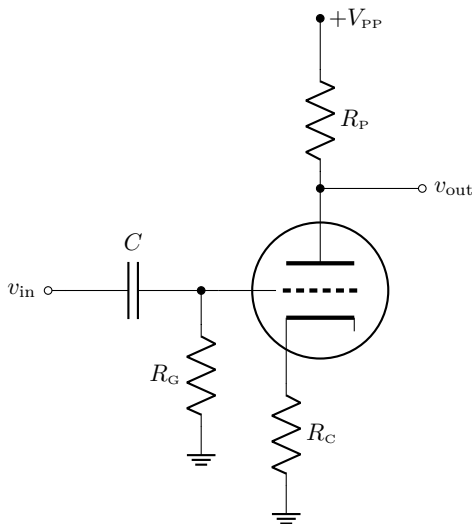
$$I_{PC} \propto (\mu V_{GC} + V_{PC})^{3/2} \quad (\mu V_{GC} + V_{PC} > 0), \quad (6.2)$$

where  $\mu$  is the triode **magnification factor**, and represents the ratio of the two proportionality constants for the two geometries (grid–cathode to plate–cathode). The magnification factor typically falls in the range of 10 to 100 or more. In this simple model, the current is zero if the tube is “reverse-biased” ( $\mu V_{GC} + V_{PC} < 0$ ). Again, the above triode law only acts as a rough guide to triode operation; triode circuits are usually analyzed graphically or on the computer for quantitatively accurate results.

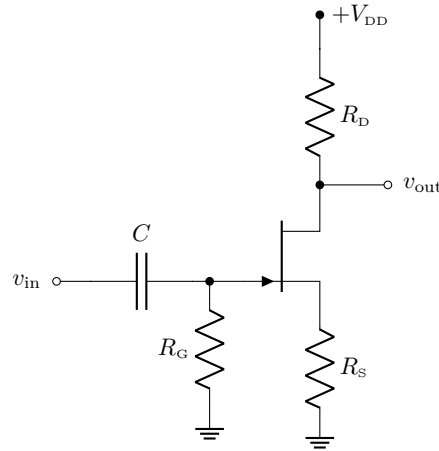
An example of a popular triode tube is the 12AX7 dual triode (i.e., 2 triodes in one package; “12” = 12.6-V filament, center-tapped to also work with 6.3 V; “7” = 7 elements: 2 anodes, 2 grids, 2 cathodes, 1 heater). This tube is especially popular in audio preamplifiers. It typically operates with a plate voltage  $V_{PC}$  of 100 to 250 V. The typical grid voltage  $V_{GC}$  is  $-1$  to  $-2$  V; the magnification factor is  $\mu = 100$ ; the grid transconductance  $g_m$  is in the range of 1250 to 1600  $\mu S$ ; and the typical plate current is 0.5 to 1.2 mA.

### 6.2.1 Triode Voltage Amplifier (Common-Cathode Amplifier)

A typical circuit for using a triode as an amplifier, the **common-cathode amplifier**, is shown below.



Again, due to the analogy with the JFET, the tube circuit is very similar to the JFET (common-source) amplifier shown below, which we analyzed before (Section 5.4.3).



In both cases, the grid/gate is biased at 0 V, and the cathode/source resistor functions to elevate the cathode/source voltage so the grid/gate is biased slightly negative with respect to the cathode/source. In the vacuum-tube case, this is called a **cathode-biased** amplifier.

The common-cathode amplifier will generally also drive a load, which we can regard as having a Thévenin resistance  $R_L$ . The net effect is that  $R_P$  should be replaced by  $R_P \parallel R_L$ , and the supply voltage  $V_{PP}$  should be replaced by the appropriate linear combination of the supply voltage and the Thévenin voltage of the load. By using  $R_P$  and  $V_{PP}$ , the following analysis also handles a loaded amplifier if these modifications are kept in mind.

### 6.2.1.1 DC Bias

As an example of tube-amplifier analysis, we will work through the analysis of the triode amplifier here. The plate power supply  $+V_{PP}$ , powers the arrangement, supplies plate-cathode current, and as a result it again sets up the bias voltage for the grid. In traditional lingo, the plate power supply is called the “B+ voltage,” in reference to the three elements A, B, and C of the triode (“A” being the cathode and “C” being the grid), and the A, B, and C batteries that were used to power the early circuits.

To begin, the dc current is determined by the power-supply voltage, which must be dropped across the plate and cathode resistors as well as the tube itself along the plate-cathode path:

$$V_{PP} - V_{PC} = I_{PC}(R_P + R_C). \quad (6.3)$$

This is called a “load-line equation,” because when written as a relation for the plate current  $I_{PC}$  as a function of the plate-cathode voltage  $V_{PC}$ , this becomes the equation for a line with slope  $-1/(R_P + R_C)$  and intercept  $V_{PP}/(R_P + R_C)$ :

$$I_{PC} = \frac{V_{PP} - V_{PC}}{R_P + R_C}. \quad (6.4)$$

This current then determines the grid bias point via

$$V_{GC} = -V_C = -I_{PC}R_C. \quad (6.5)$$

Finally, the tube itself has a response current  $I_{PC}(V_{GC}, V_{PC})$  to the grid and plate voltages. In general, this response function does not have a simple form, so the design is often done graphically on a plot of  $I_{PC}$  vs.  $V_{PC}$  for various values of  $V_{GC}$ . The load line and grid-bias relations are superimposed on this plot, and the intersection of the two determines the operating point. However, we can do this using the Child–Langmuir Law (6.2) to illustrate the process algebraically (at the cost of accuracy):

$$I_{PC} = a(\mu V_{GC} + V_{PC})^{3/2}. \quad (6.6)$$

Using Eqs. (6.3) and (6.5) to eliminate the voltages, the result is

$$I_{PC} = a \left[ V_{PP} - [(\mu + 1)R_C + R_P] I_{PC} \right]^{3/2}. \quad (6.7)$$

This is effectively a cubic equation in  $I_{PC}^{1/3}$ , and thus has an analytic (but messy) solution. The point is that it can in principle be solved for  $I_{PC}$ , given the resistances and the power-supply voltage. Once this dc current is known, this fixes the grid bias via Eq. (6.5).

### 6.2.1.2 Naïve AC Analysis

Once the grid bias is known, the grid (ac) transconductance

$$g_m := \frac{\partial I_{PC}}{\partial V_{GC}} \quad (6.8)$$

(tube transconductance: definition)

is then known. In the Child–Langmuir model the transconductance is

$$g_m = \frac{\partial I_{PC}}{\partial V_{GC}} = \frac{3a\mu}{2}(\mu V_{GC} + V_{PC})^{1/2}. \quad (6.9)$$

Note that even though the Child–Langmuir model is a rough approximation, the derivative still serves to define the transconductance, and the concept hold in general (though with a different quantitative dependence on the plate and grid voltages).

Now we will proceed with an analysis that traces along the lines of the JFET voltage amplifier. However, it yields expressions that in general do not work well for triode amplifiers (although for pentodes below this analysis works well). Afterwards, we will show how to improve it to handle triode amplifiers.

Given a small input ac voltage  $v_{in}$  (on top of a dc bias, as in the previous section), assuming a high enough frequency that the input-capacitor impedance is negligibly small, the input voltage causes an ac plate current (i.e., a small change in  $I_{PC}$ )

$$i_{PC} = g_m(v_G - v_C) = g_m(v_{in} - i_{PC}R_C). \quad (6.10)$$

Solving for  $i_{PC}$ ,

$$i_{PC} = \frac{g_m}{(1 + g_m R_C)} v_{in}. \quad (6.11)$$

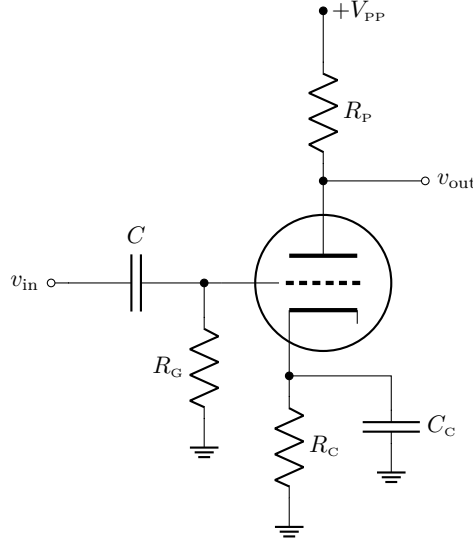
Then this current determines the ac output via

$$v_{out} = -i_{PC}R_P = -\frac{g_m R_P}{(1 + g_m R_C)} v_{in}, \quad (6.12)$$

for an ac voltage gain

$$G = -\frac{g_m R_P}{(1 + g_m R_C)}. \quad (\text{naïve!}) \quad (6.13)$$

Note that in Eq. (6.10) leading up to this, the effect of the cathode resistor was to take away some of the gain, because an increase  $v_{in}$  causes increasing  $i_{PC}$ , but this somewhat counteracts the input voltage via an increase in  $v_C$ . If desired, this effect may be counteracted by including a bypass capacitor across  $R_C$ , which removes its effect for ac signals, as shown below.



In this case, the gain is simply

$$G = -g_m R_P. \quad (\text{still naïve!}) \quad (6.14)$$

As in analogous transistor circuits, the input impedance is controlled mainly by  $R_G$ , and the output impedance is set mainly by  $R_P$ .

### 6.2.1.3 Proper AC Analysis: Plate Resistance

The important complication that we left out of the above analysis is the **plate resistance**: the plate current depends not only on the gate voltage  $V_{GC}$ , but also on the plate voltage  $V_{PC}$ . In the JFET-amplifier analysis of Section 5.1, we explicitly assumed the saturation regime, where the drain–source current is (approximately) independent of the drain–source voltage. But this same assumption does not carry over in the same way in many tube-amplifier circuits.

The plate resistance  $r_p$  is defined by the current–voltage slope

$$\frac{1}{r_p} := \frac{\partial I_{PC}}{\partial V_{PC}}, \quad (\text{plate resistance: definition}) \quad (6.15)$$

in analogy to the intrinsic emitter resistance of the BJT [see Eq. (4.69)]. For a 12AX7 tube,  $r_p$  is typically in the range of  $60 - 80 \text{ k}\Omega$ , which is usually significant on the scale of the other resistances in the circuit.

Now, to account explicitly for the ability of either the gate or plate voltage to influence the plate current, we can write a differential change in the current as

$$\begin{aligned} dI_{PC} &= \frac{\partial I_{PC}}{\partial V_{PC}} dV_{PC} + \frac{\partial I_{PC}}{\partial V_{GC}} dV_{GC} \\ &= \frac{dV_{PC}}{r_p} + g_m dV_{GC}, \end{aligned} \quad (6.16)$$

after using Eqs. (6.8) and (6.15). Suppose that  $I_{PC}$  is held constant; then the right-hand side of the above equation vanishes, leading to the relation

$$g_m r_p = -\frac{\partial V_{PC}}{\partial V_{GC}}. \quad (6.17)$$

To evaluate the derivative, going back to the triode Child–Langmuir relation (6.2), at constant  $I_{PC}$  this relation leads to

$$\mu V_{GC} + V_{PC} = \text{const.} \quad (6.18)$$

if  $V_{GC}$  and  $V_{PC}$  change by infinitesimal amount  $dV_{GC}$  and  $dV_{PC}$ , respectively, this relation implies that they must cancel, and thus

$$\mu = -\frac{\partial V_{PC}}{\partial V_{GC}}. \quad (6.19)$$

This of course holds for more general models where the current depends on the combination  $\mu V_{GC} + V_{PC}$ , but more fundamentally, we can view this partial derivative as the definition of the magnification factor, as it expresses the relative effect of the two voltages on the current via  $\mu = -(\partial I_{PC}/\partial V_{GC})/(\partial I_{PC}/\partial V_{PC})$ . Combining this relation with Eq. (6.17), this leads to the relation

$$\mu = g_m r_p \quad (\text{small-signal-parameter relation}) \quad (6.20)$$

connecting the various small-signal tube parameters

With the above development in hand, we can return to the ac gain of the common-cathode amplifier on p. 170. Using Eqs. (6.16) in terms of biased ac signals, we have

$$\begin{aligned} i_{PC} &= g_m v_{GC} + \frac{v_{PC}}{r_p} \\ &= g_m (v_{in} - v_C) + \frac{v_{PC}}{r_p}. \end{aligned} \quad (6.21)$$

To evaluate  $v_{PC}$ , note that the plate-cathode voltage drops if  $I_{PC}$  rises, because voltage drops across  $R_P$  and  $R_C$  increase:

$$v_{PC} = -i_{PC}(R_P + R_C). \quad (6.22)$$

Putting this and  $v_C = i_{PC}R_C$  into Eq. (6.21) gives

$$i_{PC} = g_m (v_{in} - i_{PC}R_C) - \frac{R_P + R_C}{r_p} i_{PC}, \quad (6.23)$$

or after solving for  $i_{PC}$ , we find

$$i_{PC} = \frac{g_m r_p v_{in}}{r_p + (g_m r_p + 1)R_C + R_P} = \frac{\mu v_{in}}{r_p + (\mu + 1)R_C + R_P} \quad (6.24)$$

after applying Eq. (6.20). Then using  $v_{out} = -i_{PC}R_P$  and  $G = v_{out}/v_{in}$ , the voltage gain becomes

$$G = -\frac{\mu R_P}{r_p + R_P + (\mu + 1)R_C}. \quad (\text{common-cathode amplifier voltage gain}) \quad (6.25)$$

For the amplifier on p. 172 where  $R_C$  is bypassed by a capacitor, we can set  $R_C = 0$  in the gain expression, which reduces to

$$G = -\frac{\mu R_P}{r_p + R_P}. \quad (\text{common-cathode amplifier voltage gain, with cathode bypass}) \quad (6.26)$$

It is only in the limit of *large* plate resistance (so  $I_{PC}$  is weakly affected by  $V_{PC}$ ) that the naïve from the previous section follow. Rewriting Eq. (6.25) as

$$G = -\frac{g_m R_P}{1 + R_P/r_p + (g_m r_p + 1)R_C/r_p}, \quad (6.27)$$

we can see that taking the limit  $r_p \rightarrow \infty$  leads to

$$G = -\frac{g_m R_P}{1 + g_m R_C}, \quad (\text{common-cathode amplifier voltage gain, infinite plate resistance}) \quad (6.28)$$

in agreement with Eq. (6.13). Bypassing the cathode capacitor then leads to the expression  $G = -g_m R_P$ , as in Eq. (6.14). This situation with large  $r_p$  is more likely to occur in pentode and beam power tubes, as described below, when the plate resistor (and load) are relatively small compared to the plate resistance. As far as obtaining a large gain is concerned, a larger  $r_p$  leads to larger gain.

### 6.2.2 Design Example: 12AX7 Preamplifier

As a specific example, consider a preamplifier stage made from (half of) a 12AX7, with

- $R_P = 100 \text{ k}\Omega$ ,
- $R_C = 1.5 \text{ k}\Omega$ ,
- $R_G = 1 \text{ M}\Omega$  (the grid resistor only really matters for the input impedance and to ensure a proper grid bias),
- $V_{PP} = 330 \text{ V}$ , and
- model values of  $a = 1.89 \times 10^{-6} \text{ A/V}^{3/2}$  and  $\mu = 100$  for the 12AX7 (see the next section),

a numerical solution of Eq. (6.7) yields

$$I_{PC} \approx 1.04 \text{ mA}, \quad (6.29)$$

which leads to a grid bias from Eq. (6.5) of

$$V_{GC} = -1.57 \text{ V}, \quad (6.30)$$

a tube plate-cathode voltage from Eq. (6.3) of

$$V_{PC} = 224 \text{ V}, \quad (6.31)$$

and a grid transconductance from Eq. (6.9) of

$$g_m = 2330 \mu\text{S}. \quad (6.32)$$

Note that the transconductance falls above the range quoted for the 12AX7 in the first part of Section 6.2. With  $\mu = 100$  by assumption and

$$r_p = \mu/g_m = 43.0 \text{ k}\Omega \quad (6.33)$$

from Eq. (6.20), the amplifier gain is

$$G = -34, \quad (6.34)$$

from Eq. (6.25) [ $G = -52$  from the naïve expression (6.13)], or if the cathode resistor is bypassed, the gain is

$$G = -70 \quad (6.35)$$

[ $G = 233$  from the naïve expression (6.14)]. Again, these values should not be taken too seriously, given the simple Child–Langmuir model, although we will see below that they aren't too bad in this example.

### 6.2.3 Phenomenological Triode Model

It is possible to model the triode behavior more accurately than via the Child–Langmuir model. Koren, for example, proposed the following phenomenological model for a triode:<sup>3</sup>

$$\begin{aligned} V_1 &= \frac{V_{PC}}{k_P} \log \left\{ 1 + \exp \left[ k_P \left( \mu^{-1} + \frac{V_{GC}}{\sqrt{k_{VB} + V_{PC}^2}} \right) \right] \right\} \\ I_{PC} &= \frac{V_1^x}{k_G} [1 + \text{sgn}(V_1)]. \end{aligned} \quad (6.36)$$

Note that when  $V_{PC}^2 \gg k_{VB}$ ,  $k_P(\mu^{-1} + V_{GC}/V_{PC}) \gg 1$ , and  $x = 3/2$ , these equations reduce to the Child–Langmuir Law (6.6) in the form  $I_{PC} = (V_{PC}/\mu + V_{GC})^{3/2}/k_G$ , so that we can identify  $a = 2/k_G\mu^{3/2}$ .

<sup>3</sup>Norman L. Koren, “Improved vacuum tube models for SPICE simulations. Part 1: Models and example,” [http://www.normankoren.com/Audio/Tubemodspice\\_article.html](http://www.normankoren.com/Audio/Tubemodspice_article.html) (update 2003).

Koren's parameters for the 12AX7 read as follows:

$$\begin{aligned}
 \mu &= 100 \\
 k_G &= 1060 \text{ V}^x / \text{A} \\
 k_P &= 600 \text{ V} \\
 k_{VB} &= 300 \sqrt{\text{V}} \\
 x &= 1.4.
 \end{aligned} \tag{6.37}$$

Then solving Eqs. (6.36) numerically, along with the load-line equation (6.4) and the grid-bias equation (6.5), we find

$$\begin{aligned}
 I_{PC} &= 1.07 \text{ mA} \\
 V_{PC} &= 221 \text{ V} \\
 V_{GC} &= -1.61 \text{ V}.
 \end{aligned} \tag{6.38}$$

To “seed” the root-finding process, the Child–Langmuir results from the previous section may be used as initial guesses. While the previous results are fairly close (they were within a few percent of these values), it is evident that using the better model makes a difference in the design results. Now the transconductance follows from differentiating Eqs. (6.36),

$$g_m = \frac{\partial I_{PC}}{\partial V_{GC}} = 1880 \mu\text{S}, \tag{6.39}$$

and Eq. (6.20) gives

$$r_p = \frac{\mu}{g_m} = 53.3 \text{ k}\Omega, \tag{6.40}$$

which is somewhat low compared to the RCA specification<sup>4</sup> Eq. (6.25) then leads to a circuit gain of

$$G = -33 \tag{6.41}$$

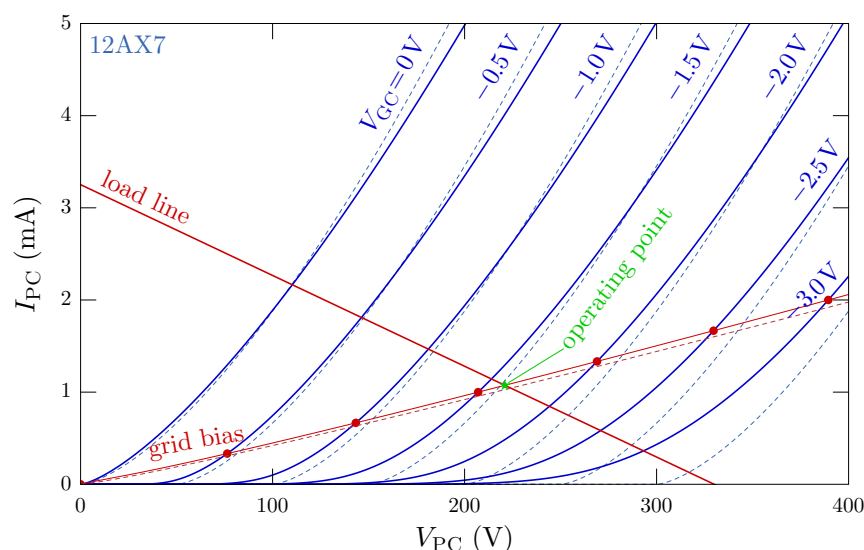
with the feedback from  $R_C$ , or

$$G = -65 \tag{6.42}$$

with  $R_C$  bypassed. The Child–Langmuir versions of the gain values were more than 10% from the correct values.

So is this accuracy typical for the Child–Langmuir model? To see this, and to also see how to do this calculation graphically (in case only graphical tube data are available from a data sheet), the plot below shows the plate  $I_{PC}$ – $V_{PC}$  characteristics for various gate voltages, shown as the solid curves (from the Koren model). The accompanying dashed curves show the Child–Langmuir predictions by comparison; note that these agree in some regions, but not in others.

<sup>4</sup>Available at <http://drtube.com/datasheets/12ax7-rca1962.pdf>, see the plot on the last page.

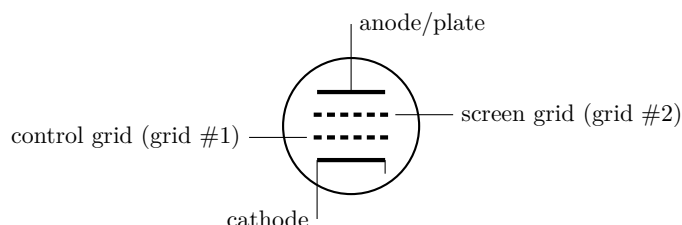


To solve this system, the next step is to draw in the load line (6.3), which is independent of any tube characteristics. Then, draw in the grid-bias curve (6.5) by drawing the point on each curve of constant  $V_{GC}$  that satisfies Eq. (6.5) (i.e., by tracing along the curve to the proper current  $I_{PC} = -V_{GC}/R_C$ ). Then the idea is to fill in the curve; while it is not a straight line, a straight line is a reasonably good approximation near any point. The intersection of the load line and the grid-bias curve determine the (dc) operating point of the circuit, and the slope of the grid-bias curve at the operating point gives the transconductance  $g_m$ . Alternately, it is possible to more directly obtain the ac amplifier gain, in the case of a bypassed  $R_C$ : ac changes in the grid voltage don't affect the bias, so one can follow the load line to find the change in  $I_{PC}$  for a change in  $V_{GC}$ , and convert to a gain via  $R_p$ .

From this diagram it is evident that the circuit happens to operate in the region where the Child–Langmuir model works reasonably well. For tubes in a push–pull circuit, which operate near zero current when the output signal crosses through zero, the simple model makes for a much poorer approximation. Note that the model (6.36) is not perfect either: compared to the RCA data sheet, for example,<sup>5</sup> the  $V_{GC} = 0$  curve has a “kink” with the opposite curvature for small  $V_{PC}$ . However, this model is probably reasonable for estimates in practical designs.

## 6.3 Vacuum Tetrodes

While the triode can work well, it has an issue, the **Miller capacitance** between the plate and grid. This is analogous to the Miller effect in bipolar transistors that we mentioned before in Section 4.11.5.3, but of course the vacuum-tube version came first. The Miller capacitance can cause a loss of bandwidth, but even worse, it can cause instability in high-gain circuits, because it acts as a feedback coupling between the high-current output and the sensitive input of an amplifier stage. One solution to this was to introduce a second grid, called the **screen grid** (also called **grid #2**, with the control grid going by the cheeky alternative **grid #1**). The resulting tube is the **tetrode**, in reference to the fourth element. The schematic symbol for the tetrode is shown below.



<sup>5</sup><http://www.r-type.org/pdfs/12ax7.pdf>



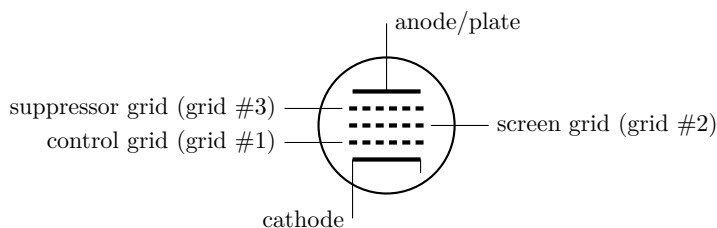
The screen grid sits between the plate and the control grid. The idea is to have a positive, fixed potential of the screen grid,  $V_{G2C} > 0$ , but still negative with respect to the plate ( $V_{G2C} < V_{PC}$ ). Typically this voltage is set via a voltage divider from the plate voltage, and is preferably also bypassed to ground with a capacitor. The control grid is still biased negative ( $V_{G1C} < 0$ ). Since the screen grid is at fixed potential, it acts as a shield (hence the name “screen,” in reference to electrical screening), breaking the capacitive coupling between plate and control grid. The new grid can reduce the control-grid–plate capacitance from several pF to below 0.01 pF.

Since  $V_{G2C} > 0$ , the screen grid also attracts electrons from the cathode, which we will see can be a problem. However, this has an advantage, since the screen grid shields somewhat the effect of the anode on the cathode electrons. The net effect is that the plate current  $I_{PC}$  is relatively insensitive to the plate voltage  $V_{PC}$  (for sufficiently large  $V_{PC}$ ), which can be a useful property.

## 6.4 Vacuum Pentodes and Beam Power Tubes

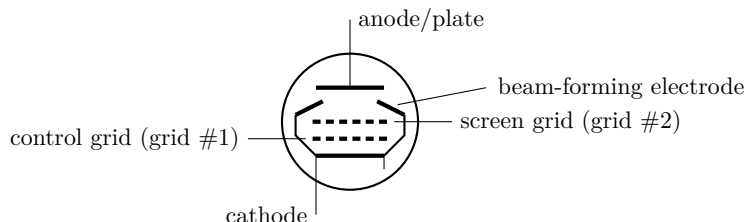
The tetrode, alas, has yet another problem, related to **secondary emission**, which is the emission of electrons from the anode due to cathode electrons striking it. In diodes and triodes, secondary-emission electrons are simply attracted back to the anode, and so don’t affect things much (except through space-charge effects). However, in tetrodes, the secondary electrons are attracted to the screen, especially if the screen grid is biased above the plate. The net effect is to reduce the plate current and overall gain of the circuit, and it reduces the useful range of  $V_{PC}$ .

The solution is to add yet another grid, the **suppressor grid** or **grid #3**, forming the **pentode**. The suppressor grid sits between the screen grid and plate; usually it is shorted to the cathode, and it functions to repel secondary electrons. The pentode diagram is shown below.



The design of tetrode and pentode circuits is essentially the same as in triode circuits, with the extra complication of the additional freedom of choosing extra grid voltages, and of more complicated mathematical models for the tube characteristics.

An alternative solution to the suppressor grid comes in the **beam power tube**. This is generally a pentode, with grid #3 replaced by a **beam-forming electrode**. (Confusingly, data sheets can refer to the same beam power tubes as tetrodes *or* pentodes.) The schematic symbol appears below.



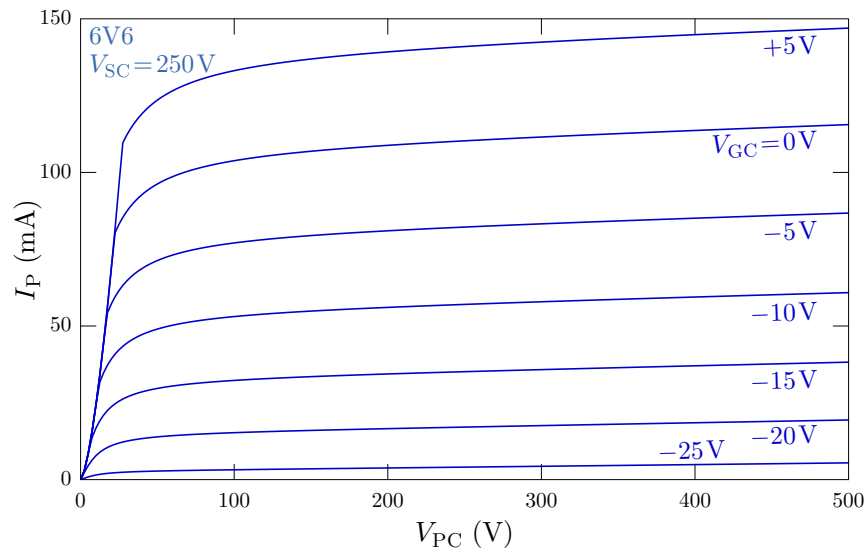
Like the suppressor grid, the beam-forming electrode is typically shorted to the cathode, and serves to direct the cathode electrons in a beam onto a particular part of the plate. Because of the controlled electron impact on the plate, such tubes can handle high currents. It also serves to suppress secondary electrons in the same way as the suppressor grid. Famous examples of beam power tubes are the 6V6 at lower powers and the 6L6 at somewhat higher powers (“6” = 6 elements: cathode, heater, plate, control grid, screen grid, beam-forming electrode). For comparison to the 12AX7 preamplifier tube, the 6V6 in typical operation has a

plate voltage ranging from  $V_{PC} = 180\text{--}315\text{ V}$ , has a screen voltage from  $V_{G_2C} = 180\text{--}225\text{ V}$ , has a control-grid voltage from  $V_{G_1C} = -8.5\text{--}13\text{ V}$ , has a transconductance of  $g_m \sim 3700\ \mu\mathcal{S}$ , an output of  $2\text{--}5.5\text{ W}$ , and a plate current  $I_{PC} = 29\text{--}34\text{ mA}$ .

A naive model of the behavior of pentodes and beam power tubes would simply extend the Child–Langmuir Law in the spirit of Eq. (6.2) for triodes, to include the effect of the screen bias voltage:

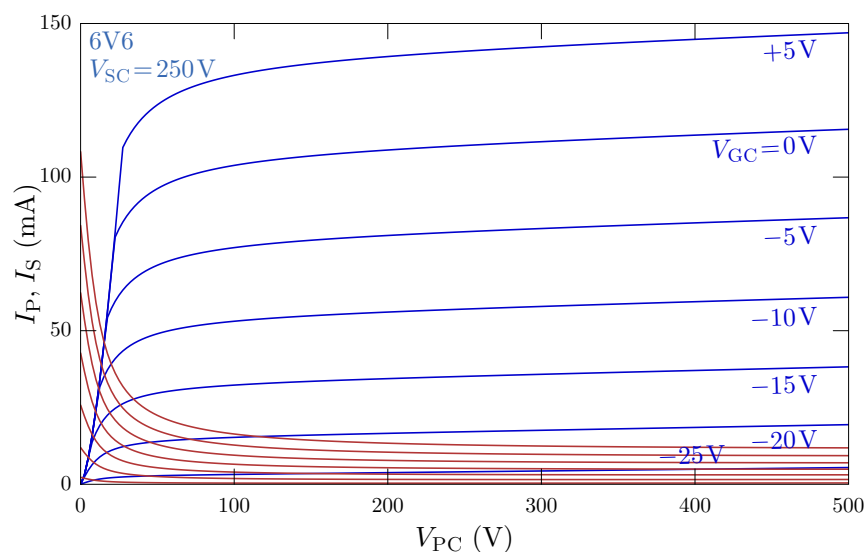
$$I_{PC} \propto (\mu_s V_{SC} + \mu_G V_{GC} + V_{PC})^{3/2} \quad (\mu_s V_{SC} + \mu_G V_{GC} + V_{PC} > 0). \quad (6.43)$$

However, this law has a much smaller range of applicability than the analogous triode law: it does not describe the regime where  $I_{PC}$  becomes insensitive to  $V_{PC}$ , as we mentioned above for tetrodes. To visualize this, we can plot characteristic curves for the typical operation of the 6V6 beam tetrode, with a fixed screen–cathode voltage  $V_{SC} = 250\text{ V}$ . A reasonable mathematical model for this is essentially Eq. (6.43), but with extra factors to enforce the “knee” behavior at larger  $V_{PC}$ .<sup>6</sup> Compare this to the 12AX7 curves in the figure on 176: the plate current here rises rapidly for small  $V_{PC}$  but then levels off for larger plate voltage.



Note that this diagram may be somewhat deceptive, because at small  $V_{PC}$ , the screen voltage is much higher than the plate, and so the screen can attract electrons and develop a substantial current. The figure below superimposes the plate-current curves  $I_s$  corresponding to the same gate voltages  $V_{GC}$  (ordered with increasing current in the same way as the plate-current curves).

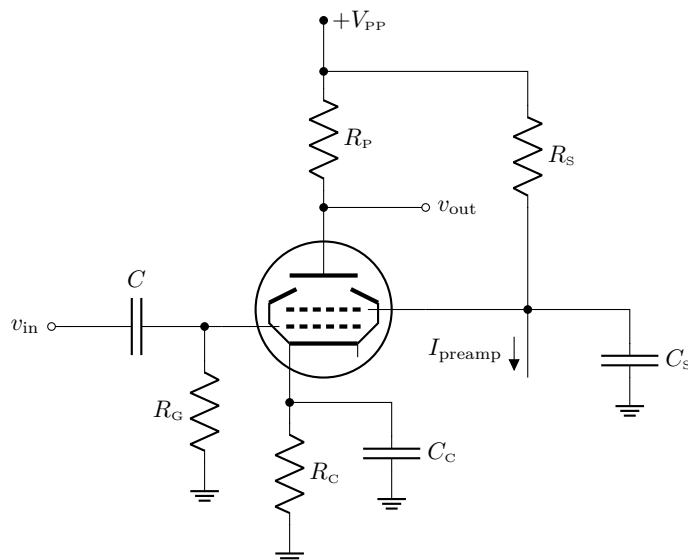
<sup>6</sup>The model used here comes from D. Munro, “PSpice Model: 6V6GT,” <http://www.duncanamps.com/pdf/6v6spicemod.pdf>, a PSpice-language mathematical model of the tube behavior. Another model using the same equations but different parameters was posted by Robert McLean at <http://www.audiobanter.com/showthread.php?t=70465>, and this model gives similar results.



The screen current can be substantial in such a tube, but the screen is a relatively delicate device compared to the plate. For this reason, these tube circuits often employ a resistor to limit the screen current to protect and prolong the life of the screen.

#### 6.4.1 Design Example: 6V6 Power Output Stage

As an example of a pentode-type amplifier circuit, consider the beam-power-tube amplifier power amplifier shown below, based on the 6V6 tube.



The circuit is similar to the 12AX7 preamplifier before, but now there is a screen resistor  $R_S$  that limits the screen current and voltage. A small current  $I_{\text{preamp}}$  is also tapped at the screen to power a preamplifier circuit and further drop the screen voltage.

Unfortunately, the analysis of this power amplifier is more complicated than the triode case because of the extra degree of freedom of the screen voltage. Also the screen current  $I_S$  can be significant. Since there are now three currents (plate, screen, cathode) in the tube that add as

$$I_C = I_P + I_S, \quad (6.44)$$

the load-line graphical analysis is more complicated because the voltage drops across  $R_P$  and  $R_C$  arise due to different currents:

$$V_{PP} - V_{PC} = I_P R_P + I_C R_C. \quad (6.45)$$

This leads to a load “line” which isn’t a line anymore, because  $I_S$  is not simply related to  $I_P$ :

$$V_{PP} - I_P R_P - I_C R_C = V_{PC}. \quad (6.46)$$

Further,  $I_S$  determines the ground-referenced screen voltage via

$$V_S = V_{PP} - (I_S + I_{\text{preamp}})R_S, \quad (6.47)$$

but the screen–cathode voltage is also determined by the cathode bias:

$$V_{SC} = V_S - I_C R_C. \quad (6.48)$$

We still have the plate and control-grid bias voltages, as in the triode case:

$$\begin{aligned} V_{PC} &= V_P - I_C R_C \\ V_{GC} &= -I_C R_C. \end{aligned} \quad (6.49)$$

This leads to a relatively complicated coupled system of equations. The graphical analysis in particular is complicated by the need to consult graphs of both the plate and cathode currents to even set the proper bias. In practice, crude approximations such as assuming  $I_S$  to be a fixed fraction of  $I_C$  are sufficient for approximate design work. But if a mathematical model of the tube is available, it is more convenient nowadays to let a computer grind out the numerical solution of these equations.

As an example, suppose we assume the following parameter set:

- $R_P = 284.2 \Omega$  (i.e., the dc resistance of the primary coil of an output transformer),
- $R_C = 470 \Omega$ ,
- $R_S = 1 \text{ k}\Omega$ ,
- $R_G = 220 \text{ k}\Omega$  (again, the grid resistor only really matters for the input impedance and to tie the grid to ground; the grid current is negligible provided it is biased negative with respect to the cathode),
- $V_{PP} = 360 \text{ V}$ , and
- $I_{\text{preamp}} = 2 \text{ mA}$  to power a small preamplifier.

The numerical solution to the above equations, including the 6V6 model, gives

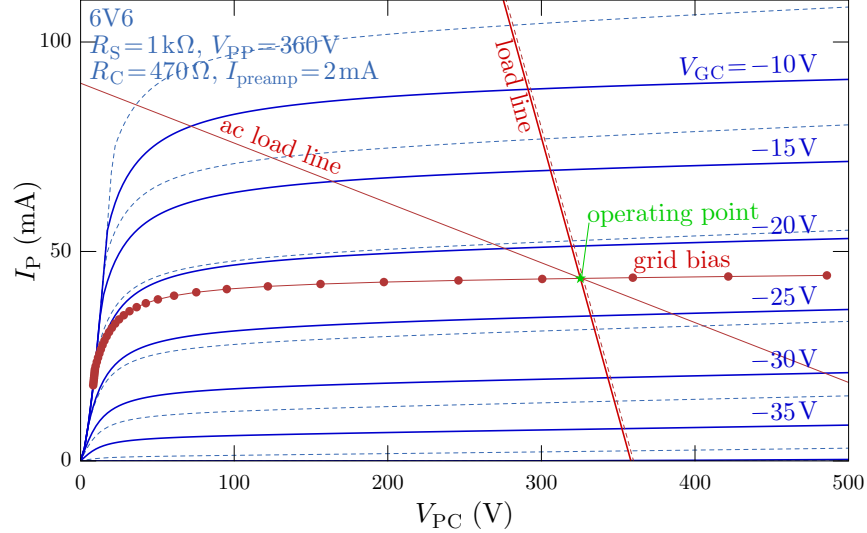
$$\begin{aligned} I_P &= 43.6 \text{ mA} \\ I_S &= 3.75 \text{ mA} \\ I_C &= 47.3 \text{ mA} \\ V_{PC} &= 325 \text{ V} \\ V_{GC} &= -22.2 \text{ V} \\ V_S &= 354 \text{ V} \\ V_P &= 347 \text{ V}. \end{aligned} \quad (6.50)$$

To check the tube operation, these results imply power dissipations of

$$\begin{aligned} P_P &= V_{PC} I_P = 14.2 \text{ W} \\ P_P &= V_{SC} I_S = 1.25 \text{ W} \end{aligned} \quad (6.51)$$

at the plate and screen, respectively. The 6V6 has specified maximum powers of 12 W and 2 W, respectively, so this design is pushing the limits of the plate power. (If the plate power exceeds the specification by too much, the tube will “red plate,” or the anode will glow red with heat, dramatically shortening the tube’s life.)

All this and more is illustrated in the plot below. The characteristic curves for this arrangement are shown as solid curves; the dashed curves show the corresponding characteristics when  $V_{SC}$  is held at a constant 330 V, illustrating the importance of numerically solving for the correct  $V_{SC}$  in the analysis.



The (dc) load curve is plotted from Eq. (6.46); the dashed counterpart is the load line obtained by ignoring the small contribution of  $I_S$ . Clearly the screen current makes a difference, but here it doesn’t change much to ignore it.

The grid-bias curve is more difficult to construct graphically here because it only intersects one characteristic curve (for  $V_{GC} = -20$  V). The plotted points are separated by  $\Delta V_{GC} = 0.1$  V, from  $V_{GC} = -19$  V to  $-22.5$  V, from left to right. Recall that this curve is determined by the cathode current, which is determined by the plate and screen voltages. The intersection with the load line again gives the operation point.

A different, ac load line is also shown. Since  $R_p$  corresponds to the primary winding of an output transformer, it has a different ac response than the dc resistance. The plotted curve assumes a nominal ac impedance of 7 k $\Omega$ . This load line shows how changes to the gate voltage translate into output-current changes. Note that if  $V_{GC}$  rises much above  $-10$  V, the curves bunch together, resulting in clipping of the signal; the same is true if  $V_{GC}$  goes below  $-35$  V or so.

Computing the ac gain directly is somewhat easier than setting the dc bias, because the bypass capacitors allow us to assume that the cathode voltage  $V_C$  and screen voltages  $V_S$  and  $V_{SC}$  are fixed. The transconductance can be evaluated by numerically differentiating the tube model. The result here is

$$g_m = \frac{\partial I_P}{\partial V_{GC}} = 4330 \mu\text{S}, \quad (6.52)$$

which is similar to values specified in the data sheet. The plate resistance then follows from adapting the definition (6.15) to read

$$\frac{1}{r_p} := \frac{\partial I_P}{\partial V_{PC}}. \quad (6.53)$$

This derivative can be evaluated numerically as well, and comes to

$$r_p = 72.5 \text{ k}\Omega. \quad (6.54)$$

Because of the extra dependence on the screen voltage, the relation  $g_m r_p = \mu$  of Eq. (6.20) no longer holds exactly. For example, the 6V6 model assumes a grid  $\mu$  of 390 (and a screen  $\mu$  of 43), but  $g_m r_p = 314$  from the numerical calculations

The ac voltage gain from Eq. (6.26) is then

$$G = -\frac{g_m r_p Z_P}{r_p + Z_P} = -27.7, \quad (6.55)$$

after restoring  $g_m r_p$  in place of  $\mu$  and taking  $R_p$  to be the ac impedance  $Z_P$ . Note that the “naïve” result  $G = -g_m Z_P$  from Eq. (6.14) works reasonably well here, predicting  $G = -30.3$ , because the plate resistance is indeed substantially larger than the plate impedance, by more than a factor of 20. In any case, the gain here is lower than in the triode preamplifier example.

### 6.4.2 Triode Connection

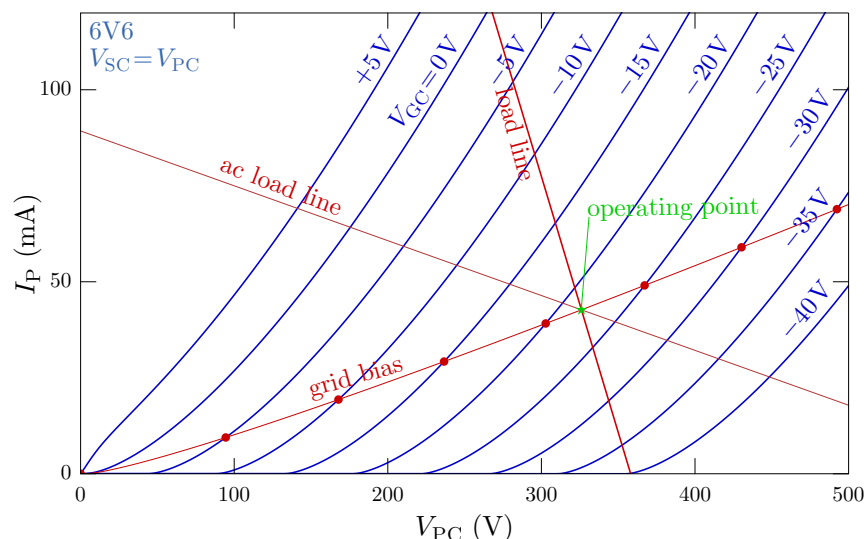
The characteristic curves in the triode-amplifier example of Section 6.2.3 appeared quite different from the beam-power-tetrode example of Section 6.4.1. Roughly speaking the triode is operated in a regime where it behaves like a variable resistor, controlled by the grid voltage. This is evident in the gain expressions (6.25) and (6.26), where the plate resistance (along with the magnification) determine the tube’s contribution to the amplifier gain. The pentode-type circuits, however, when the screen is held at a fixed potential, are different: they typically run in a regime with large plate resistance, making the tube output behave more like a current source (since the plate current is insensitive to  $V_{PC}$ ). In this regime the tube’s contribution to the amplifier gain comes in the form of the transconductance, as in Eq. (6.28).

Pentodes (and beam power tubes) work well enough, but sometimes triode characteristics are more desirable, for example for their distortion or low-noise properties. One common trick for achieving this with a pentode-type tube is the **triode connection**, where the screen is shorted to the plate (typically this is done via a resistor of  $\sim 100 \Omega$ , which limits screen current and suppresses high-frequency instabilities).

As an example, consider the same power-amplifier stage as in the previous section, but with a triode-connected 6V6 (with a direct plate–screen short for simplicity). The analysis is essentially the same, except Eq. (6.47) is replaced by the condition  $V_s = V_p$ . The dc operating point turns out to be nearly the same (all the values below are within a few percent of the previous values):

$$\begin{aligned} I_P &= 42.7 \text{ mA} \\ I_S &= 3.68 \text{ mA} \\ I_C &= 46.4 \text{ mA} \\ V_{PC} &= 326 \text{ V} \\ V_{GC} &= -21.8 \text{ V} \\ V_s &= 348 \text{ V} \\ V_P &= 348 \text{ V}. \end{aligned} \quad (6.56)$$

The characteristic curves are shown below, with the same load curves, and the appropriate grid-bias curve.



As might be expected, the curves here resemble more those of the 12AX7 (in the plot on p. 176) than the pentode-type curves (from the plot on p. 182).

For the ac analysis, the transconductance is similar to the pentode-operation case, with

$$g_m = 4310 \mu\text{S} \quad (6.57)$$

(formerly  $g_m = 4330$ ), but the plate resistance drops dramatically to

$$r_p = 2.1 \text{ k}\Omega \quad (6.58)$$

(from a former  $72.5 \text{ k}\Omega$ ), because it measures the plate-current change caused by variation in the plate voltage, which also varies the screen voltage (the screen voltage has a stronger effect via its magnification). The resulting gain is then

$$G = -\frac{g_m r_p Z_P}{r_p + Z_P} = -6.8, \quad (6.59)$$

(reduced from  $-27.7$ ). This is fairly typical behavior obtained from triode-connecting a pentode-type tube: much-decreased gain, but with (potentially) more desirable characteristics. For example, along the ac load line, the constant- $V_{GG}$  curves bunch together at the extremes, indicating distortion. In the triode case, they appear to remain evenly spread out, which is one consequence of the lower gain. However, another characteristic of the triode configuration is the distortion tends to turn on smoothly, which leads to a “warmer” triode sound, compared to a “clean” and possibly “harsh” pentode sound in audio circuits.

## 6.5 Circuit Practice: Simple Tube Amplifier

As an example of a more complete vacuum-tube circuit, consider below an amplifier based on three tubes for electric guitar or microphone. This amplifier design is similar to various “Champ” amplifiers made by Fender Musical Instruments Corp., and the voltages are nominal values (with 20% stated acceptable tolerance) reported on schematics of those instruments.<sup>7</sup>

<sup>7</sup>Gerald Weber, *A Desktop Reference of Hip Vintage Guitar Amps* (Kendrick Books, 1996) (ISBN: 0964106000).





### 6.5.2 Second Preamplifier Stage

The second-stage preamplifier, after the volume control, is essentially identical to the first, except for the addition of the  $47\text{-}\Omega$  resistor and the feedback path from the output transformer. We will save the feedback for later, but the extra resistance yields modified values of  $V_{\text{PC}} = 226\text{ V}$  (vs.  $225\text{ V}$ ) and  $V_{\text{GC}} = -1.59\text{ V}$  (vs.  $-1.61\text{ V}$ ). Again, these are close to the nominal values. The gain expression (6.25) gives  $G = -54$ , using  $R_{\text{P}} = 68.8\text{ k}\Omega$  (including the  $220\text{-k}\Omega$  input resistor of the power stage) and  $R_{\text{C}} = 47\text{ }\Omega$ . The data sheet for these conditions gives  $r_{\text{p}} = 60\text{ k}\Omega$ ,  $g_{\text{m}} = 1650\text{ }\mu\text{S}$ , and  $\mu = 101$ , giving a gain of  $G = -52$ .

### 6.5.3 Power Amplifier Stage

The power-amplifier stage is intended to match the prototype power-amplifier circuit of Section section:6V6amp, including the  $1\text{ mA}$  of current that supplies each preamplifier stage. The output transformer acts as a load; the transformer is a Hammond 1750C.<sup>9</sup> The primary dc resistance is specified as  $284.2\text{ }\Omega$ , but can vary considerably (one example measured  $267\text{ }\Omega$ ). The transformer specifies an ac impedance ratio of  $7000 : 3.2$ , which means that a  $3.2\text{-}\Omega$  speaker load appears as a  $7\text{-k}\Omega$  load at the 6V6 plate. (The ac impedances are nominal values, varying widely with frequency because of the frequency response of the transformer and especially of the speaker.) The ac gain for this stage is  $G = -27.7$ .

### 6.5.4 Feedback Loop

The main remaining feature of the circuit, besides the power supply, is the  $2.7\text{-k}\Omega$  resistor connecting the transformer secondary to the cathode resistor chain of the second preamplifier stage. This provides overall *negative* feedback in the circuit. That the feedback is negative is not entirely obvious, and so it's worth quickly tracing it through. When the input voltage on the second preamp increases; the input to power tube goes down (owing to the negative gain of the common-cathode amp); the plate voltage on power tube goes up; this reduces current in output-transformer primary (i.e., the change in current points towards the dot); this increases current in secondary (i.e., the change in current again points to dot); and finally this increases cathode voltage at the preamp, which counteracts effect of input (in a way similar to an unbypassed cathode resistor).

To work out the feedback gain, note that the transformer *voltage* ratio is  $\eta_{\text{T}} = \sqrt{3.2/7000}$ . (Recall that the voltage ratio of a transformer is the same as the turns ratio, while the current ratio is the inverse of the turns ratio; hence the impedance  $Z = V/I$  ratio is the *square* of the turns ratio.) Then the net open-loop gain from the input to the second preamplifier to speaker is

$$G_{\text{net}} = \eta_{\text{T}} G_{\text{pre},2} G_{\text{power}} = \sqrt{3.2/7000} \times (-54) \times (-27.7) = 32, \quad (6.60)$$

where  $G_{\text{pre},2}$  is the voltage gain of the second preamplifier, and  $G_{\text{power}}$  is the voltage gain of the power-amplifier stage. To account for the effect of the feedback, we can replace the  $47\text{-}\Omega$  and  $2.7\text{-k}\Omega$  resistors by the  $R_{\text{f}} = 46\text{ }\Omega$  parallel (Thévenin) resistance and a Thévenin voltage of  $v_{\text{f}} = \eta_{\text{div}} G_{\text{net}} v_{\text{in}}$ , where  $v_{\text{in}}$  is the input to the second preamplifier stage, and  $\eta_{\text{div}} = 0.0171$  is the voltage-divider ratio due to the  $47\text{-}\Omega$  and  $2.7\text{-k}\Omega$  resistors.

With this modification, Eq. (6.23) is modified to read (setting  $R_{\text{C}} = 0$  since it is bypassed with a capacitor

$$\begin{aligned} i_{\text{PC}} &= g_{\text{m}}(v_{\text{in}} - v_{\text{C}}) - \frac{R_{\text{P}} + R_{\text{f}}}{r_{\text{p}}} i_{\text{PC}} \\ &= g_{\text{m}}(v_{\text{in}} - i_{\text{PC}} R_{\text{f}} - v_{\text{f}}) - \frac{R_{\text{P}} + R_{\text{f}}}{r_{\text{p}}} i_{\text{PC}} \\ &= g_{\text{m}}(v_{\text{in}} - i_{\text{PC}} R_{\text{f}} - \eta_{\text{div}} G_{\text{net}} v_{\text{in}}) - \frac{R_{\text{P}} + R_{\text{f}}}{r_{\text{p}}} i_{\text{PC}}. \end{aligned} \quad (6.61)$$

Solving for  $i_{\text{PC}}$  gives

$$i_{\text{PC}} = \frac{g_{\text{m}}(1 - \eta_{\text{div}} G_{\text{net}}) v_{\text{in}}}{1 + g_{\text{m}} R_{\text{f}} + (R_{\text{P}} + R_{\text{f}})/r_{\text{p}}}. \quad (6.62)$$

<sup>9</sup>The transformer data sheet: <http://www.hammondmfg.com/pdf/EDB1750C.pdf>.

Then with  $v_{\text{out}} = -i_{\text{PC}}R_{\text{P}}$  the closed-loop ac voltage gain for the preamplifier becomes

$$G_{\text{pre},2,\text{CL}} = -\frac{\mu(1 - \eta_{\text{div}}G_{\text{net}})R_{\text{P}}}{r_{\text{p}} + (\mu + 1)R_{\text{f}} + R_{\text{P}}} = -24, \quad (6.63)$$

again using the parameters  $g_{\text{m}} = 1880 \mu\text{S}$ ,  $\eta_{\text{div}} = 0.0171$ ,  $G_{\text{net}} = 32$ ,  $R_{\text{P}} = 68.8 \text{ k}\Omega$ ,  $R_{\text{C}} = 1.5 \text{ k}\Omega$ ,  $r_{\text{p}} = 5.19 \text{ k}\Omega$ ,  $\mu = 100$ , and  $R_{\text{f}} = 46 \Omega$ . The preamplifier-stage gain here is reduced from the open-loop value of  $G_{\text{pre},2} = -54$ . The reduction in gain is of course accompanied by a slight improvement in gain flatness with frequency and reduction in distortion (see Section 7.7). However, the improvements are not extreme, since the change in gain is only a factor of 2.

With the negative feedback, the gain from the second preamplifier input to the speaker after the volume control is then  $G_{\text{net,CL}} = \eta_{\text{T}}G_{\text{pre},2,\text{CL}}G_{\text{power}} = 14$ . The gain including the first preamplifier stage ( $G_{\text{pre},1} = -63$ ), when the volume control is maximized, is  $G = 900$  (dropping the sign, which doesn't matter anymore). The gain from first input to the transformer primary (i.e., the 6V6 plate voltage) is  $4.2 \times 10^4$ .

To give a sense of scale of this voltage, a 10-mV input signal gives a much larger signal at the input of the power tube of  $(10 \text{ mV}) \times 24 \times 63 = 15 \text{ V}$ . This is comparable to the gate bias voltage of the 6V6. Thus, at the maximum gain setting, the amplifier is already at the threshold of overdriving the input of the power tube, leading to distortion. The distortion can be desirable for electric-guitar amplifiers, and is easily attainable here: magnetic pickups for electric guitars have typical maximum signal levels in the range of 100–500 mV. Dynamic microphones tend to have smaller, mV-level signals, which would lead to much less distortion in this circuit.

### 6.5.5 Power Supply

The power supply is a straightforward variation of the full-wave rectifier on p. 169. However, the voltage is not entirely obvious. The plate winding of the power transformer<sup>10</sup> is specified at 630 V, center-tapped, for 100 mA load current. This results in 315 V to either anode of the 5Y3, but this is rms: the peak anode voltage is 445 V. The rectifier tube drops 50 V at 125 mA dc forward current, but this does not explain the 360 V nominal output voltage in the diagram. The transformer data sheet, however, conveniently specifies a 363 V (dc) output using a 5Y3 tube at 100 mA, implying a larger voltage drop than suggested by the GE data sheet.<sup>11</sup> The answer to this is on page 4 of the 5Y3 data sheet, which shows the loaded, rectified output voltage, which has more “sag” when filtered by a capacitor: the plot on the data sheet shows about 340 V dc output for a 100-mA dc output current and 315-V rms input (with a 20- $\mu\text{F}$  filter capacitor), and about 380 V output for a 50-mA dc output. These numbers are reasonably consistent with the nominal voltage in the schematic, and reflect differences in characteristics of tubes available from different manufacturers.

The rectifier output is then filtered by a 20- $\mu\text{F}$  capacitor, as in the 5Y3 data sheet, and this supplies the plate voltage for the 6V6. A 1-k $\Omega$  resistor feeds the screen from the same voltage, as we analyzed in the 6V6 power-amplifier example. This drops the screen voltage down to the plate-voltage level. The screen voltage is also bypassed by another 20- $\mu\text{F}$  capacitor, decoupling the screen and plate. The screen voltage then supplies the 2 mA to the 12AX7 plates, with an appropriate voltage drop across a 10-k $\Omega$  resistor. The preamplifier plate voltages are bypassed by yet another 20- $\mu\text{F}$  capacitor, decoupling the preamplifier tubes from the power tube and from each other.

<sup>10</sup>Magnetic Components 40-18027, <http://www.classictone.net/40-18027.pdf>.

<sup>11</sup><http://www.r-type.org/pdfs/5y3gt.pdf>

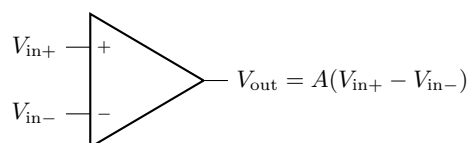


## Chapter 7

# Operational Amplifiers

### 7.1 Op-Amp Basics

We have already talked a bit about differential amplifiers, when we analyzed the transistor differential amplifier. The differential amplifier is an important building block in its own right, and we will spend a fair amount of time looking at them. In particular, the **operational amplifier (op-amp)** is a handy, handy device in analog-circuit design. Recall that the differential amplifier ideally subtracts the input voltages, and multiplies the difference by a gain factor (here,  $A$ ) to produce an output signal.



The operational amplifier is basically a differential amplifier, but with large gain (with larger=better, where op-amps are concerned). The gain  $A$  here is called the **open-loop gain** of the op-amp, for reasons that will become more apparent soon. For real op-amps, gains typically range from around 46 dB on the low end (for low-quality amplifiers, or amplifiers where other engineering considerations compromised the gain), to around 140 dB on the high end. (Remember that 20 dB corresponds to a factor of 10 in voltage, so 140 dB means  $A = 10^7$ .) For more specific examples, let's summarize the gains of a few classic op-amps:

- **741C**: has BJT inputs,  $A = 86$  dB, a famous, old op-amp that is cheap, and not so great anymore (there are better choices, even among cheap op-amps).
- **LF411**: has JFET inputs,  $A = 88$  dB, a cheap op-amp, which is really not bad, and a good “default op-amp” in noncritical applications.
- **OPA111B**: has JFET inputs,  $A = 120$  dB, a venerable, precision, low-noise op-amp, but *expensive*.

#### 7.1.1 Usage: Open-Loop

There are two basic ways to use an op-amp: in the first, we take advantage of the high gain, and in the second, we throw away some of the high gain. The first application, where we use the full gain of the op-amp, is called **open-loop mode**, and the op-amp behaves as a **comparator**. To explain this, note that the op-amp is an active device, and requires a power supply; often op-amps are powered by split  $\pm 15$ -V power supplies, so the output can go either positive or negative. The output, of course, cannot exceed the power-supply voltages (and can typically the output range, or **output swing**, of the op-amp is a volt or two less than the supply range). Now when the output formula,

$$V_{out} = A(V_{in+} - V_{in-}) \tag{7.1}$$

(op-amp output formula)

predicts that  $V_{\text{out}}$  should be outside the supply-voltage range, what *really* happens is that the output **rails**; for example, if the output *should* be +50 V, but the power supplies are  $\pm 15$  V, the output will *rail* at +15 V (or more likely a bit lower, say around 14.3 V).

Then the **comparator** action of an op-amp is as follows:

1. If  $V_{\text{in}+} > V_{\text{in}-}$  by at least a few mV, then the  $V_{\text{out}}$  rails at the *positive* supply voltage.
2. If  $V_{\text{in}+} < V_{\text{in}-}$  by at least a few mV, then the  $V_{\text{out}}$  rails at the *negative* supply voltage.

That is, the op-amp *compares* the input voltages, and swings the output to the appropriate rail to indicate which one is bigger. There are specialized op-amps, called **comparators**, that are optimized to do this, and we will consider comparator circuits in more detail later. For now, note that regular op-amps can be used in this way, though this is not the most common usage.

### 7.1.2 Usage: Closed-Loop

In **closed-loop mode** or **negative-feedback mode**, op-amps have some connection from the output to the inverting (–) input. The net effect is to reduce the gain. Why bother to have an op-amp with extremely high open-loop gain only to reduce it in closed-loop mode? Well, it turns out that when you do this, the resulting circuit behaves *well*, in the sense that its behavior will be (mostly) independent of the device properties. This becomes more true with increasing open-loop gain. This trick of using **negative feedback** is a really nice trick, and opens up a lot of possibilities for cool circuits.

## 7.2 Op-Amp “Golden Rules”

In the simplest method for analyzing op-amp circuits, we will assume that we are dealing with an *ideal* op-amp. This means that we will assume the following two rules:<sup>1</sup>

1. No current flows into or out of the inputs. (Current *can* of course flow into or out of the *output*, as well as the power-supply terminals, which we haven’t bothered to label thus far.)
2. Either  $V_{\text{in}+} = V_{\text{in}-}$ , or the output is “railed” (pinned to one of the supply voltages). Basically we are assuming that the open-loop gain  $A$  is so large, that the only way for the output to *not* be railed is for  $V_{\text{in}+}$  and  $V_{\text{in}-}$  to be almost the same, because

$$V_{\text{in}+} - V_{\text{in}-} = \frac{V_{\text{out}}}{A} \approx 0. \quad (7.2)$$

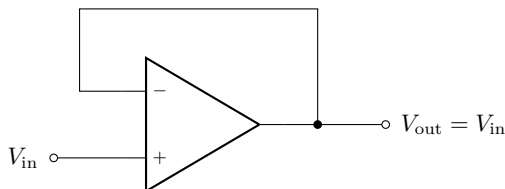
Remember that for the output not to rail, there must generally be negative feedback. In this case *the output does whatever it needs to do* to make sure the input voltages are approximately equal.

Using just these two rules, we can analyze the basic behavior of most op-amp circuits.

## 7.3 Basic Op-Amp Circuits

### 7.3.1 Unity-Gain Buffer/Follower

The first circuit is simple: the **unity-gain buffer**, or **voltage follower**. Here,  $V_{\text{out}} = V_{\text{in}}$ , and this is mostly useful to buffer a high output impedance or a low input impedance. The improvement here on the transistor followers is that this circuit works for dc voltages, not just biased ac voltages.



<sup>1</sup>These are a less succinct version of the “golden rules” of ideal op-amp behavior according to Horowitz and Hill.

To analyze this, note that the output is shorted to the inverting input, providing negative feedback. Then

$$V_{\text{in-}} = V_{\text{out}}. \quad (7.3)$$

But the second op-amp rule says that  $V_{\text{in+}} = V_{\text{in-}}$ , so  $V_{\text{in-}} = V_{\text{in}}$ , or

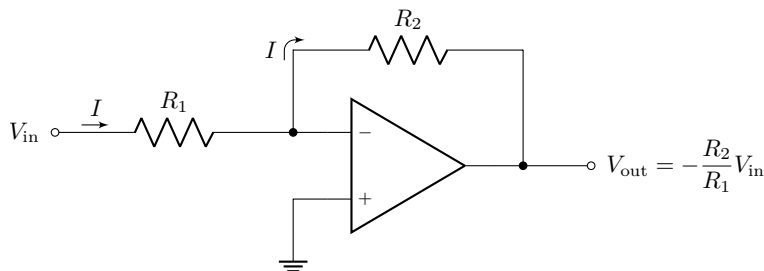
$$V_{\text{out}} = V_{\text{in}}. \quad (7.4)$$

(unity-gain buffer)

Hence, this circuit has unity gain for dc signals. We will return later to the question of the input and output impedances of this circuit, but we expect the input impedance to be high, considering that ideally no current flows into the input. But typically, we expect input impedances of  $\sim 10^8 \Omega$  for op-amps with BJT inputs, and  $\sim 10^{12} \Omega$  for op-amps with JFET inputs. The output impedances can be in the range of  $m\Omega$  for good op-amps. So this circuit is handy for buffering high-impedance sensors, for example. Another common use is to derive reference voltages in circuits: For example, if you need a voltage of +5-V somewhere in a circuit, you can use a voltage divider between  $\pm 15$ -V power supplies, and buffer the output of the divider so you don't have to worry about loading it down.

### 7.3.2 Inverting Amplifier

The **inverting amplifier** is useful as a basic amplifier *with* gain, and not only because it has many useful variations.



To analyze this circuit, assume a current  $I$  flows into the input. Golden rule 2 says that the inverting-input voltage  $V_{\text{in-}}$  must be zero, so

$$I = \frac{V_{\text{in}}}{R_1}. \quad (7.5)$$

Then according to the first rule, no current goes into the input, so it must all go through the feedback resistor  $R_2$ . Then applying Ohm's law across the feedback resistor,

$$0 \text{ V} - V_{\text{out}} = IR_2, \quad (7.6)$$

or putting in the previous expression for  $I$  and solving for the output voltage,

$$V_{\text{out}} = -\frac{R_2}{R_1} V_{\text{in}}. \quad (7.7)$$

(inverting amplifier)

Defining the **closed-loop gain**  $G$  as

$$G := \frac{V_{\text{out}}}{V_{\text{in}}}, \quad (7.8)$$

(inverting amplifier)

we can write the closed-loop gain of the inverting amplifier as

$$G = -\frac{R_2}{R_1}. \quad (7.9)$$

(closed-loop gain, inverting amplifier)

This must be smaller than the open-loop gain, otherwise the assumption of negative feedback breaks down.

### 7.3.2.1 Stability

Now it's useful to try to get an idea of how feedback works in this circuit to maintain the advertised output voltage. An alternative to the above derivation is to think of  $R_1$  and  $R_2$  as forming a voltage divider with output  $V_{in-}$ ; the voltage divider is unloaded because no current flows into the inverting input. Hence

$$V_{in-} = V_{in} \frac{R_2}{R_1 + R_2} + V_{out} \frac{R_1}{R_1 + R_2}. \quad (7.10)$$

Since the op-amp inputs must be equal, we can just set  $V_{in-} = 0$ , obtaining again the result (7.7). But suppose instead that the output has some small error  $\epsilon$ , so that

$$V_{out} = -\frac{R_2}{R_1} V_{in} + \epsilon. \quad (7.11)$$

Putting this back into the voltage-divider result (7.10), the new voltage at the

$$V_{in-} = \frac{R_1}{R_1 + R_2} \epsilon, \quad (7.12)$$

where notice that any contribution from  $V_{in}$  cancelled.

For the op-amp to work properly, it should be *stable*: that is, the op-amp should tend to remove errors like this. To see if this is the case, let's substitute this new input voltage into the (finite-gain) output formula (7.1), which says  $V_{out} = A(V_{in+} - V_{in-})$ :

$$\Delta V_{out} = -\frac{R_1}{R_1 + R_2} A \epsilon, \quad (7.13)$$

We are notating this as a *change* in the output voltage, because this merely represents the *tendency* of the op-amp's output. Since  $A$  is large, this is a strong tendency for the op-amp to move its output in opposition to  $\epsilon$  (so if  $\epsilon$  is a small, positive error, the op-amp will feel a strong “desire” to reduce it's voltage). Thus, any deviation away from  $\epsilon = 0$  makes the op-amp want to drive the output towards  $\epsilon = 0$ , correcting the error. So not only does the inverting amplifier satisfy the golden rules, but it is **stable**, in that the op-amp will correct for deviations away from this value.

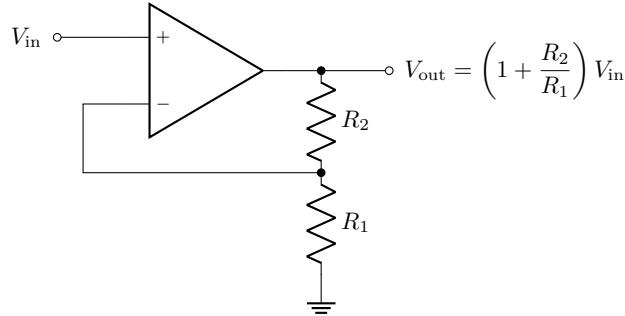
For now, though, this example allows us to be a little more precise about what we mean by “negative feedback.” If we had instead build the inverting amplifier with the op-amp inputs swapped, the golden rules would still predict the same output voltage. However, the stability analysis would have led to the same expression (7.12), but for  $V_{in+}$  instead of  $V_{in-}$ , and the op-amp's desired change in  $V_{out}$  in Eq. (7.12) would have had the opposite sign. This would mean that the op-amp would tend to strongly *reinforce* any small error, running away until the output rails.

One of the main powers of negative feedback is robustness to deviations ( $\epsilon$ ) from the proper output (e.g., due to power-supply noise). If the feedback is positive instead of negative (as in the “backwards” inverting amplifier), then more complicated behaviors like oscillation may result; we will return to this useful case later. In fact, even in the stability analysis above, it may seem that the zealous tendency of the op-amp to correct itself may lead to an overshooting output and thus an oscillation. In practice, if the op-amp's response is quick enough, the circuit will move quickly to equilibrium (as the tendency of the op-amp to self-correct will disappear as  $\epsilon \rightarrow 0$ ). However, a delayed between an output change and the input response can certainly lead to oscillation. Again, we will treat this later in terms of the varying frequency response of the op-amp circuit.

One final note is in order here, because technically there is some sleight of hand in the above derivation. Since we are using the finite gain op-amp equation  $V_{out} = A(V_{in+} - V_{in-})$  to analyze the stability, we should have used the same relation to determine the output  $V_{out}$  in the first place. We will come back to this kind of solution later in Section 7.7, but for now it is sufficient to note that the result of redoing the correct analysis is that with the input error  $\epsilon$ , we would have found that  $V_{out}$  is the expected solution plus the amplified error (see Problem 7.18 for an example). The “cheat” in the above derivation is in assuming that the output signal represented the *change* in output. However, this simple approach is usually sufficient to determine the (dc) amplifier stability, provided we're careful with the interpretation.

### 7.3.3 Noninverting Amplifier

It is also possible to build a **noninverting amplifier** with gain, as shown below.



Using that  $R_1$  and  $R_2$  form a voltage divider,

$$V_{in} = V_{in+} = V_{in-} = \frac{R_1}{R_1 + R_2} V_{out}. \quad (7.14)$$

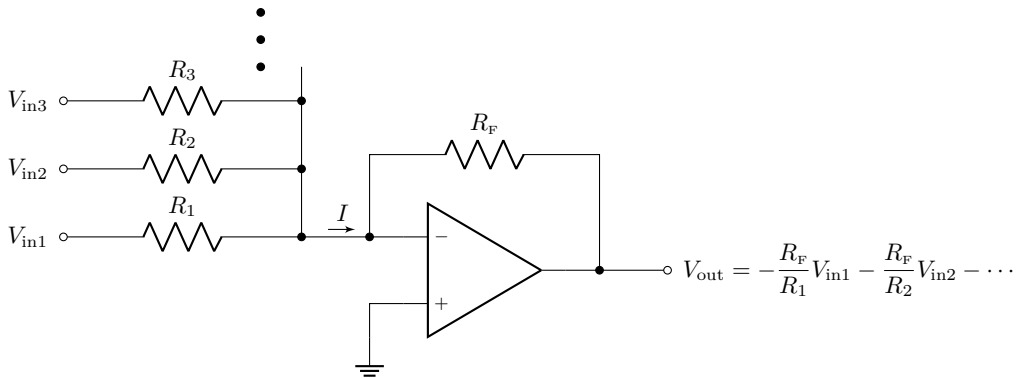
Thus, the closed-loop gain is

$$G = \frac{R_1 + R_2}{R_1} = 1 + \frac{R_2}{R_1}. \quad (\text{closed-loop gain, noninverting amplifier}) \quad (7.15)$$

Note that in the inverting-amplifier case, sub-unity gains are possible if  $R_2 < R_1$ , but here, the smallest possible gain is unity. Like the unity-gain buffer, this circuit enjoys high input impedance; the input impedance in the inverting case is just  $R_1$ .

### 7.3.4 Summing (Inverting) Amplifier

A useful variation on the inverting amplifier is the **summing amplifier**, which combines multiple input voltages with different gains to obtain the output. (Note that you *can't* easily add voltages in passive circuits, which makes this circuit useful.)



The key to this circuit is that each input voltage and input resistor makes a current; at the big junction, all these currents add and go through the feedback resistor  $R_F$ . That is, the total current  $I$  is

$$I = \frac{V_{in1}}{R_1} + \frac{V_{in2}}{R_2} + \frac{V_{in3}}{R_3} + \dots \quad (7.16)$$

Then the output voltage is

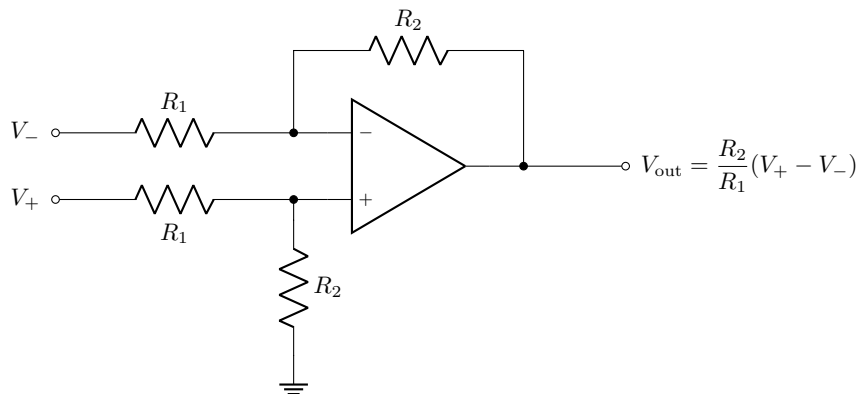
$$V_{out} = -IR_F = -\frac{R_F}{R_1} V_{in1} - \frac{R_F}{R_2} V_{in2} - \frac{R_F}{R_3} V_{in3} - \dots \quad (\text{summing amplifier}) \quad (7.17)$$



The inverting nature of the amplifier can be inconvenient, but easily fixed by following up with another inverting amplifier.

### 7.3.5 Circuit Practice: Differential Amplifier

Another classic operational amplifier is shown below. This takes the difference of two input signals, and implements a closed-loop gain given by the ratio of resistors. Normally, we save the circuit practice for the end of the chapter, but here you should work this circuit out right away to review the concepts thus far before going on.



This circuit relies on the resistor values being well-matched for accuracy. Show that this circuit behaves as advertised.

**Solution.** First, we have a voltage divider at the noninverting input:

$$V_{in+} = \frac{R_2}{R_1 + R_2} V_+. \quad (7.18)$$

Similarly, we have a voltage divider between two voltages at the inverting input:

$$V_{in-} = \frac{R_2}{R_1 + R_2} V_- + \frac{R_1}{R_1 + R_2} V_{out}. \quad (7.19)$$

Setting these two voltages equal, we get

$$R_2 V_+ = R_2 V_- + R_1 V_{out}, \quad (7.20)$$

or

$$V_{out} = \frac{R_2}{R_1} (V_+ - V_-), \quad (7.21)$$

as desired.

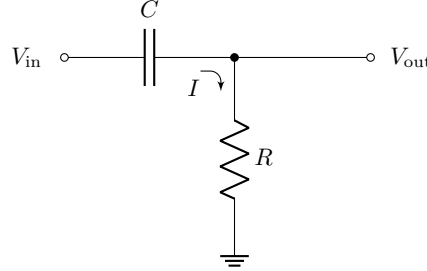
## 7.4 Op-Amp Filters

When we were studying *passive*, linear circuits with resistors, capacitors, and inductors, we saw that the reactive elements (capacitors and inductors) acted like frequency-dependent resistors. The resulting circuits attenuated the input signal in a frequency-dependent way, leading to passive **filters** for signals. We can do the same thing with op-amps: they open up new possibilities as well as straightforward improvements on the passive circuits.

The most important and fundamental op-amp filters are the op-amp versions of the passive integrator and differentiator. The op-amp versions have their own problems, but mainly because they have overall much more ideal behavior.

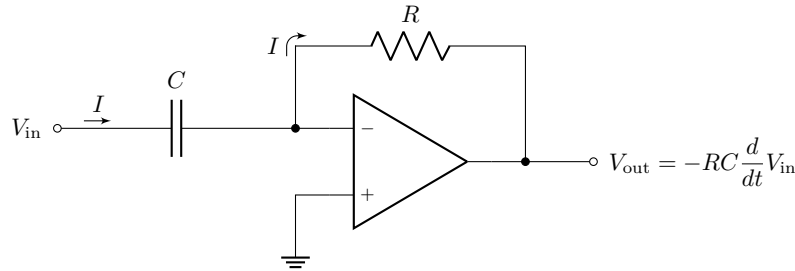
### 7.4.1 Op-Amp Differentiator

Recall the passive differentiator (Section 2.2.2), shown below.



Remember that this works intuitively as follows: The voltage across the capacitor is proportional to the charge ( $Q = CV$ ). Current  $I$  flows from the capacitor through the resistor, where current is the derivative of charge ( $I = dQ/dt$ ). The resistor converts current to voltage ( $V = IR$ ), so the output is the derivative of the input. *But*, for the output to be the derivative of the input, we had to assume here that the output voltage is small, otherwise the voltage across the capacitor is  $V_{in} - V_{out}$ , not merely  $V_{in}$ .

We can improve this, essentially by using the op-amp to decouple the capacitor and resistor voltages as follows.



Intuitively, the current  $I$  flows in through the input, and the voltage across the capacitor is  $V_{in} - V_{in-} = V_{in}$ , enforced by the op-amp. This same current  $I$  gets converted to a voltage, via the feedback resistor.

To show this quantitatively, let's use the fact that at fixed frequency  $\omega$ , this is just an inverting amplifier. The gain is  $G = -R_2/R_1$ , where we should replace  $R_2$  by  $R$ , and  $R_1$  by  $X_C = i/\omega C$ . Then

$$V_{out} = -\frac{R}{X_C} V_{in} = i\omega RC V_{in}. \quad (7.22)$$

Now remember that at fixed frequency  $\omega$ , we can identify  $d/dt \equiv -i\omega$ , so

$$V_{out} = -RC \frac{d}{dt} V_{in}. \quad (7.23)$$

(op-amp differentiator)

Note that the  $RC$  time and the time derivative conspire to make the voltage units come out right. Note also that the frequency-dependent gain

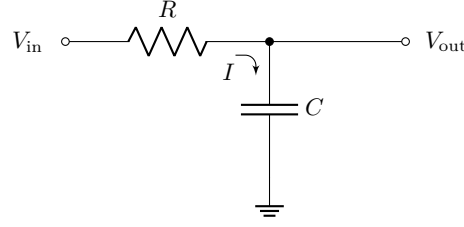
$$G(\omega) = i\omega RC \quad (7.24)$$

(gain, op-amp differentiator)

always increases with frequency, whereas the passive differentiator (high-pass filter) had a gain (transfer function) that leveled off at unity for frequencies above the  $\omega_{3dB}$  (Section 2.3.7).

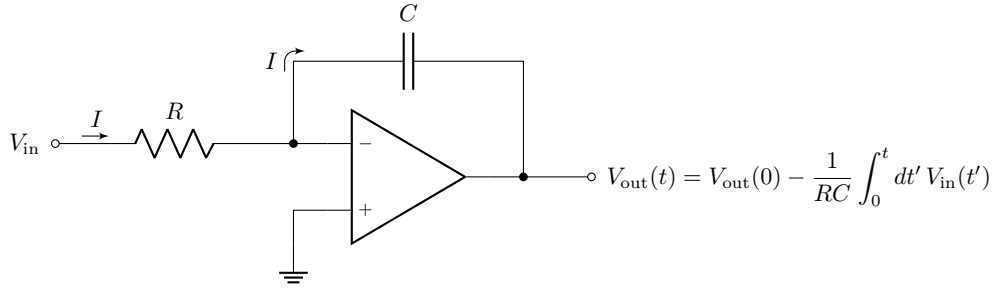
### 7.4.2 Op-Amp Integrator

The passive and active (op-amp) integrators are similar. Recall the passive integrator below (Section 2.2.1).



Here, the resistor converts the input voltage to a current ( $V_{in} = IR$ ), and the capacitor develops a voltage proportional to charge ( $Q = CV$ ), which is the integral of current. Hence the output voltage is the integral of the input voltage. But we had to assume small  $V_{out}$ , so that the voltage across the resistor is  $V_{in} - V_{out} \approx V_{in}$ .

Again, the op-amp helps here by decoupling the capacitor and resistor voltages, while connecting their currents, by maintaining the inverting input at virtual ground. The op-amp integrator is shown below.



This works more like we said: The input resistor converts the voltage  $V_{in}$  to current  $I$  ( $V_{in} = IR$ , with no need for small  $V_{out}$ ), and the capacitor integrates the current to store charge, which produces an output voltage.

For the quantitative analysis, we again use the inverting-amplifier gain  $G = -R_2/R_1$ , with  $R_2$  replaced by  $X_C$ , and  $R_1$  replaced by  $R$ . Then

$$V_{out} = -\frac{X_C}{R} V_{in} = -\frac{i}{\omega RC} V_{in}, \quad (7.25)$$

so that the frequency-dependent gain is

$$G(\omega) = -\frac{i}{\omega RC}. \quad (7.26)$$

(op-amp integrator gain)

Rearranging the factor of  $\omega$ ,

$$-i\omega V_{out} = -\frac{1}{RC} V_{in}. \quad (7.27)$$

Then using  $d/dt \equiv -i\omega$ ,

$$\frac{d}{dt} V_{out} = -\frac{1}{RC} V_{in}. \quad (7.28)$$

Integrating, we have

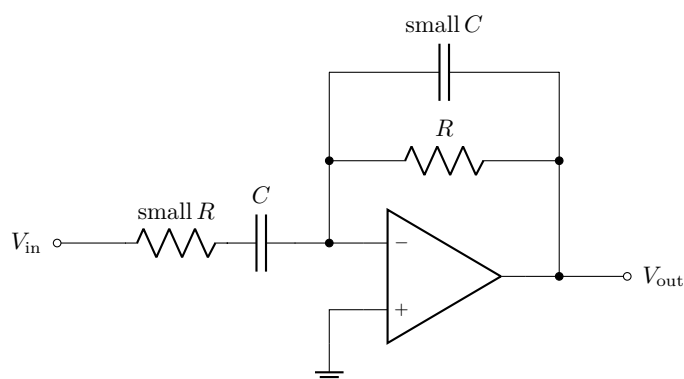
$$V_{out}(t) = V_{out}(0) - \frac{1}{RC} \int_0^t dt' V_{in}(t'). \quad (7.29)$$

(op-amp integrator)

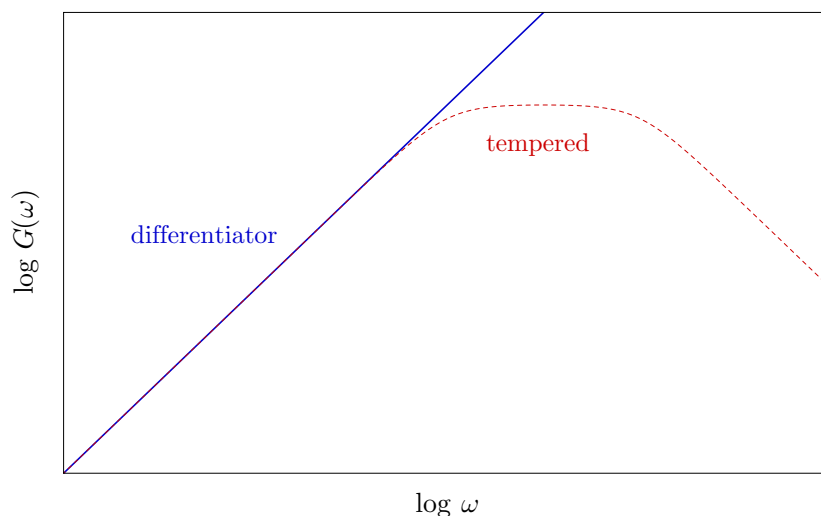
Like the differentiator, this amplifier gives the integral with an overall minus sign. It also depends on the initial output state.

### 7.4.3 Differentiator Issues

From Eq. (7.24), the main problem with the differentiator is that it has a gain that increases as  $\omega$ , so the gain becomes arbitrarily large for large frequencies. This causes potential problems in two ways. First, differentiators suffer from bad high-frequency noise, and second, the high gain can possibly de-stabilize the amplifier. The solution is to add an extra, small, parallel capacitance in the feedback loop, and an extra, small, series resistance in the input.



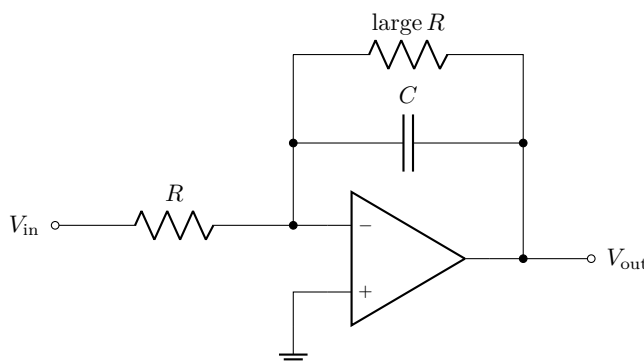
How do we think about this? First of all, the input network crosses over from capacitive to resistive at an input RC frequency. Above this frequency, the resistance dominates, and in combination with the feedback resistor, the amplifier acts like an inverting amplifier, where the gain is flat with frequency. The feedback network also defines a second RC frequency; above this frequency, the capacitor bypasses the resistor. Then the capacitor, in combination with the input resistor, makes the op-amp behave as an integrator, where the gain decreases with frequency. The net effect is shown schematically in the gain plot below.



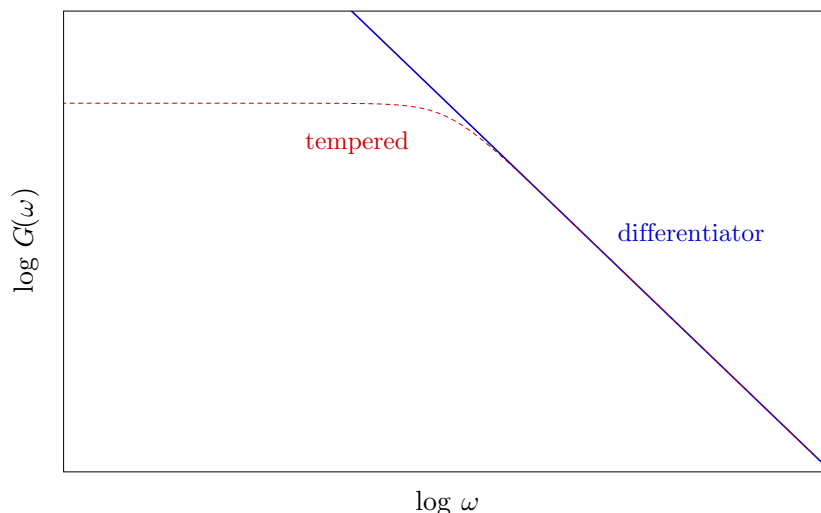
The solid line is the differentiator response, and the dashed line shows the “tempered” differentiator, including the effects of the input resistor (when the gain flattens) and the feedback capacitor (when the gain “rolls off,” or decreases).

#### 7.4.4 Integrator Issues

The integrator has the *opposite* problem: from Eq. (7.26), the gain scales as  $\omega^{-1}$  over all frequencies, so the gain at dc diverges. This essentially means that the integrator has no “natural” dc level, and any dc input (even a spurious dc input) will eventually rail the op-amp. The fix for this is to put a large resistance in parallel with the feedback capacitor, as shown below.



The feedback network then defines an RC frequency (a small frequency, since the resistance is large), below which the resistor dominates the feedback impedance, and the op-amp acts as an inverting amplifier. This levels the low-frequency gain to some finite value, as illustrated schematically below.



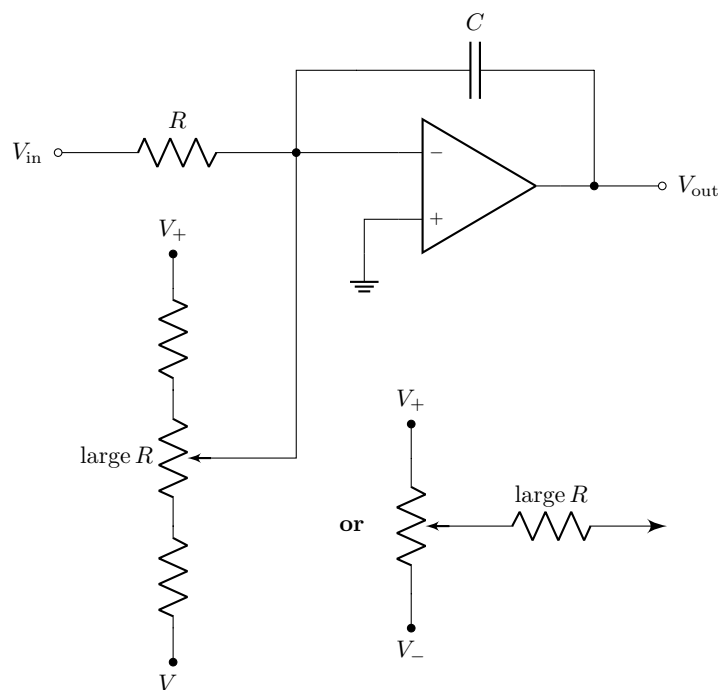
This shows the normal  $-6$ -dB scaling of the integrator as the solid line, with the dashed line showing the “tempered” behavior with the feedback resistor rolling off the low-frequency gain.

### 7.4.5 Sources of Integrator Error

We have mentioned that the basic integrator is very sensitive to spurious dc inputs, due to the divergent dc gain. What are these spurious dc offsets? There are two main sources for op-amps: input bias current, and input offset voltage.

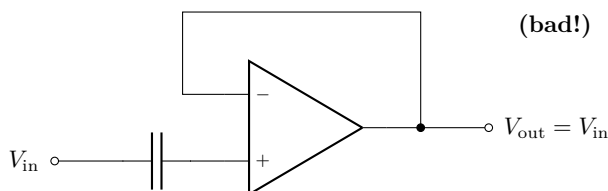
#### 7.4.5.1 Input Bias Current

In the ideal op-amp we stated the golden rule that no current flows into the inputs. However, this isn't *quite* true. The idea behind the *input bias current* is that in fact a small current flows into (or out of) the inputs, which makes sense, as the op-amp inputs drive internal transistors, which either require current to work or allow a bit of leakage current to flow. For BJT-input op-amps, the input bias current is  $\sim 10$  nA, while for JFET-input op-amps, the input bias current is  $\sim 10$  pA. For example, the precision OPA602C with JFET inputs has a 1-pA input bias current (compared to the 741C with BJT inputs at 500 nA). Thus, JFET-input op-amps are the clear winner in this regard. To some extent, it is possible to compensate for the input bias current by injecting a small, adjustable current at the inverting input.

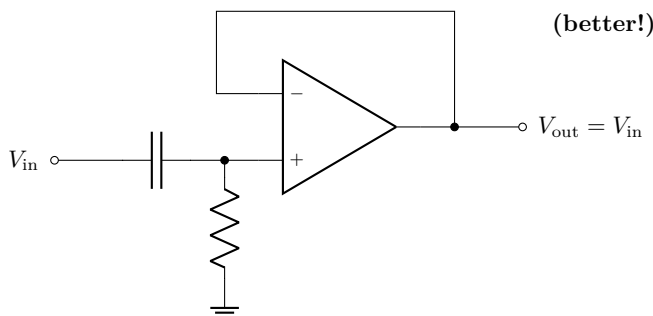


Here,  $V_{\pm}$  are the power supplies. The idea here is essentially the same as for the summing amplifier: the adjustable voltage and resistance (or Thévenin resistance of the potentiometer) causes a current to flow—a small one, if the resistance is large. This can be adjusted to cancel the input bias current, for example by adjusting it until the output is stationary when  $V_{in} = 0$ . However, this isn't perfect: the input bias current depends on temperature, for example, so compensating at one temperature doesn't guarantee compensation everywhere. If this is a concern, it is far better to start with a good op-amp, rather than try to “fix” a crappier op-amp.

Another side effect of the input bias current is that the inputs need some dc path to ground. So, for example, it would be a bad idea to build an ac-coupled follower like this:



The inverting input is fine, as the path to “ground” is via the output. However, the noninverting input has no path to ground, and the input bias current will charge the capacitors until the inputs go out of range with respect to the power supplies, causing real problems. The fix is to use an input high-pass network, as shown below.



Now, the resistor supplies the path to ground. There is some offset voltage error given by the product of the bias current and the resistance, but this can be made small (in the previous circuit, this error was huge because the impedance to ground was effectively very large).

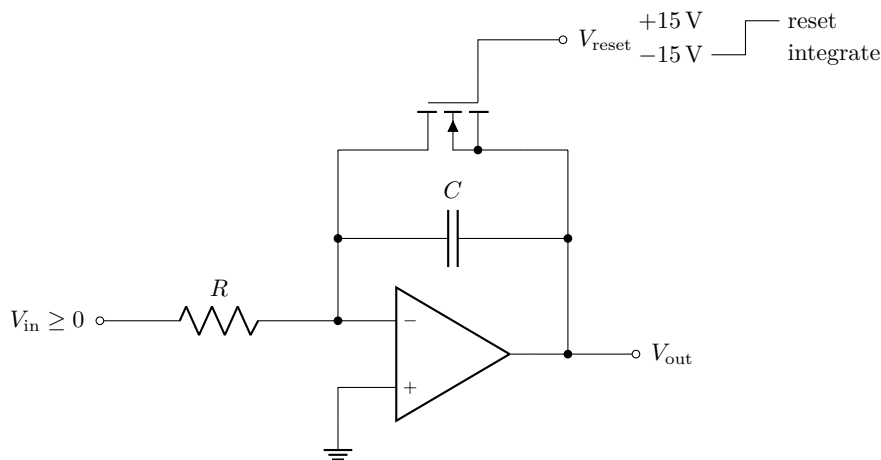
### 7.4.5.2 Input Offset Voltage

The other main effect is the **input offset voltage**. For an ideal op-amp, the output is zero when the inputs are exactly equal. But for a *real* op-amp, the output is zero when the inputs are *almost* the same, or said another way, when there is a small “error” voltage  $\delta V$  between the inputs. This is the input offset voltage, and is due to manufacturing variation when producing op-amp devices. Typically, the input offset voltage ranges from  $\sim 10\mu\text{V}$  to a few mV, with BJT inputs faring better than FET inputs. For example, the FET-input, precision OPA602C has 0.1 mV typical input offset voltage, and 0.25 mV max (compare to the 741C, which is 2 mV typical; this is not as much worse, when compared to the bias current). In the integrator, the net effect is that a zero  $V_{\text{in}}$  causes integration (the input should be set to the input offset voltage for no integration to occur). The compensation circuit above can also compensate for this effect, because it is equivalent to summing another input voltage with  $V_{\text{in}}$ . Most op-amps also have pin connections for a potentiometer to allow nulling of the input offset voltage (typically, for a single op-amp in a dual-inline package, a trim-pot is connected across pins 1 and 5). Again, while this can be trimmed, the drift will be of the same order as the uncompensated error, so in critical applications, it’s better to choose an op-amp with a low offset voltage, rather than try to correct for the offset voltage of a “bad” op-amp.

## 7.4.6 Integrator Applications

The integrator is a widely useful circuit. One example is in feedback-control circuits (circuit to generate a stable voltage, current, temperature, etc.)—it turns out integration is useful in obtaining stable operation, a point to which we will return later.

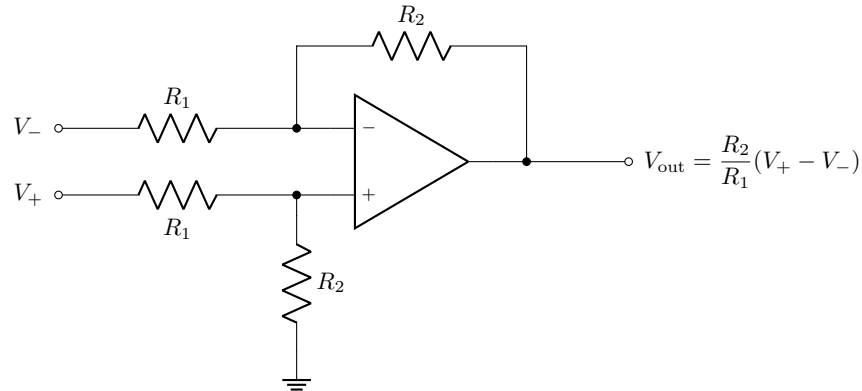
Literal integration of signals is also a useful task. For example, suppose we have an optical-pulse signal from a laser pulse on a photodetector. If we integrate the signal, we can get the pulse energy (or pulse **fluence**). In this case, it is useful to be able to reset the integrator just before we expect to receive each pulse, and the circuit below takes care of this.



The MOSFET here acts as an analog switch, that dumps the capacitor charge when it is necessary to reset the integrator. Note that we are assuming a positive input voltage, and thus a negative output voltage. If the output voltage may have either sign, a second, reversed MOSFET may be necessary to prevent dumping the capacitor charge during integration.

## 7.5 Instrumentation Amplifiers

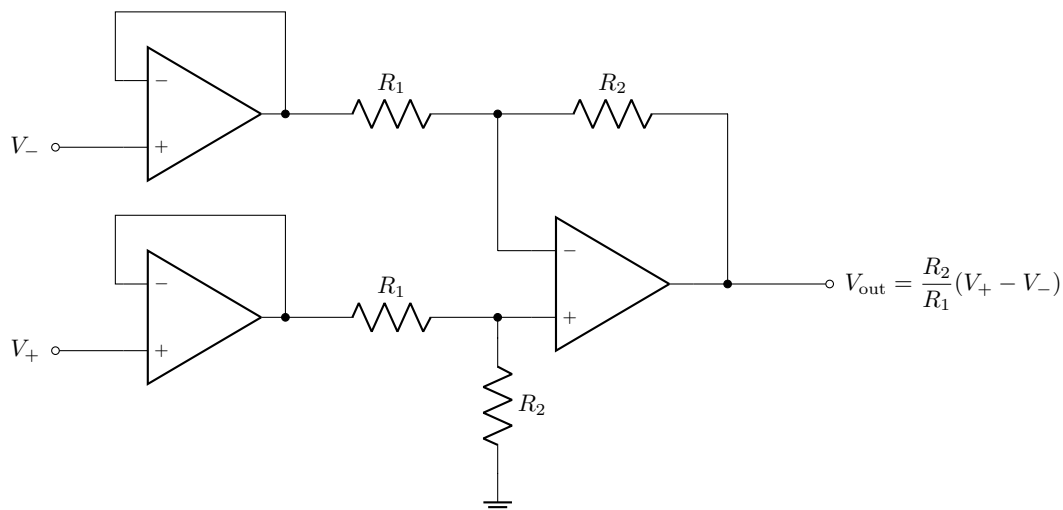
Recall the op-amp differential amplifier, shown below, from Section 7.3.5.



This circuit works fine, but has some disadvantages:

- This circuit does not have high input impedances, especially for large gains. Specifically, you should convince yourself that the noninverting input has an input impedance  $R_1 + R_2$ , while the inverting input has an input impedance  $R_1$ . But, for example, if  $R_1 = 1\text{ k}\Omega$  and  $R_2 = 10\text{ k}\Omega$ , so that  $G = 10$ , the input impedance is at worst  $1\text{ k}\Omega$ .
- The circuit requires accurately matched resistor pairs to achieve a high CMRR. To obtain a CMRR of 80 dB, the resistors must have a tolerance of around 0.01% at  $G = 1$ ; resistors this accurate are typically wirewound, but these don't work well at high frequencies.
- Any source impedances add to the input  $R_1$ 's. That is, the sources must act as ideal voltage sources, otherwise the gain and CMRR may be affected. For example, if  $R_1 = 1\text{ k}\Omega$ , the source impedance must be  $10\text{ }\Omega$  or less for 1% gain accuracy. Even worse, the source impedances must be *matched* to ensure good CMRR.

The main solution to these problems, especially that of source impedances interacting with the resistances in the op-amp circuit, is to simply buffer the inputs. A differential amplifier with buffered (high-impedance) inputs is called an **instrumentation amplifier**, and the basic circuit is shown below.

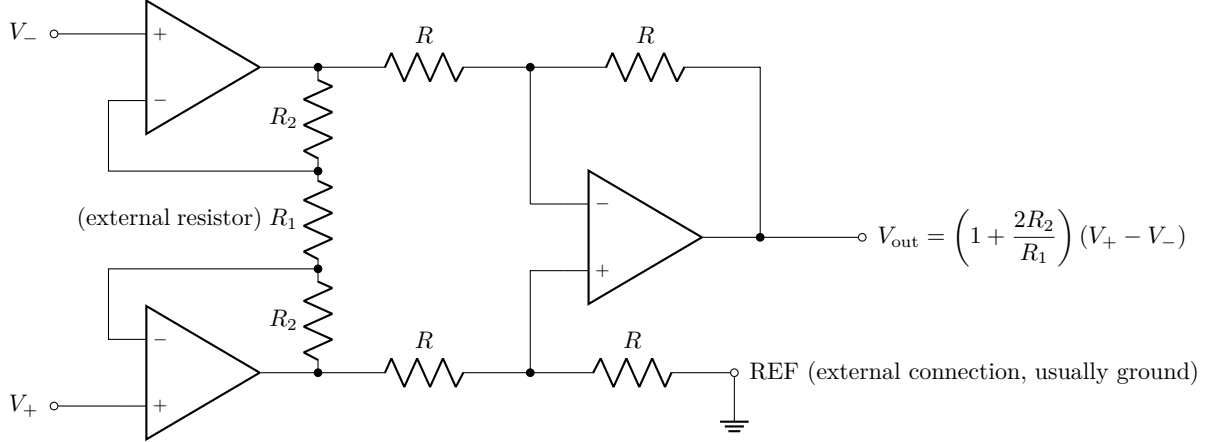


The point is that this entire circuit should come in a single package, with laser-trimmed (matched) resistor networks, for good performance and to make life easy. This also guarantees that any errors in the matched resistor pairs due to temperature drifts is kept to a minimum.



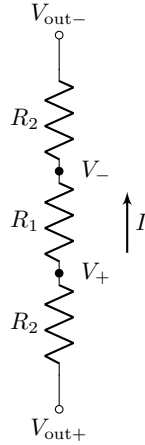
### 7.5.1 “Classic” Instrumentation Amplifier

What you *usually* find packaged as an instrumentation amplifier is actually a bit different than the above circuit. The inputs are still buffered for high impedance, but the input buffers are hooked up in a resistor chain as shown below, and the resulting difference is computed by a unity-gain differential amplifier.



The main advantage here is that the gain can be set by changing only *one* resistor, here  $R_1$ . This resistor is usually not included in the package, but rather the package has pins for an externally connected resistor for a user-settable gain. (You should be able to see that when  $R_1$  is omitted, the input amplifiers reduce to buffers, and the output is just the difference of the inputs.) Another feature to note is that the ground connection of the differential amplifier is usually given as a “reference” (REF) pin on the package. This allows the subtraction to be referenced to another voltage besides ground, which is sometimes convenient. We will give an example below in Section 7.5.2.3.

To analyze the circuit above, let’s focus on just the first resistor chain.



We have drawn in the voltages, assuming the op amps enforce the equality of their inputs, and we are labeling the outputs of the two input buffers as  $V_{\text{out}\pm}$ . No current flows into or out of the buffer-op-amp inputs, so we will assume a current  $I$  flows from  $V_{\text{out}+}$  to  $V_{\text{out}-}$ . Using Ohm’s law across  $R_1$ ,

$$I = \frac{V_+ - V_-}{R_1}. \quad (7.30)$$

Then the voltage drop across the top resistor in the chain gives

$$V_{\text{out}-} = V_- IR_2 = V_- - \frac{R_2}{R_1}(V_+ - V_-), \quad (7.31)$$

and across the bottom resistor in the chain,

$$V_{\text{out}+} = V + IR_2 = V_+ + \frac{R_2}{R_1}(V_+ - V_-). \quad (7.32)$$

Now the differential amplifier takes the difference between  $V_{\text{out}+}$  and  $V_{\text{out}-}$ , so

$$V_{\text{out}} = V_{\text{out}+} - V_{\text{out}-} = (V_+ - V_-) + \frac{2R_2}{R_1}(V_+ - V_-), \quad (7.33)$$

or

$$V_{\text{out}} = \left(1 + \frac{2R_2}{R_1}\right)(V_+ - V_-). \quad (7.34)$$

(instrumentation-amplifier output)

As a gain, this reads

$$G = \left(1 + \frac{2R_2}{R_1}\right). \quad (7.35)$$

(instrumentation-amplifier gain)

Thus we see again if  $R_1$  is omitted ( $R_1 = \infty$ ), the gain is unity, while other values can serve to increase the gain above unity.

A good example of an instrumentation amplifier is the INA128 from Burr-Brown, one of a family of “INXXX” instrumentation amplifiers. In the INA128, the internal  $R_2$  resistors are  $50\text{ k}\Omega$ , so the external resistor  $R_G$  (i.e.,  $R_1$ ), sets the gain via

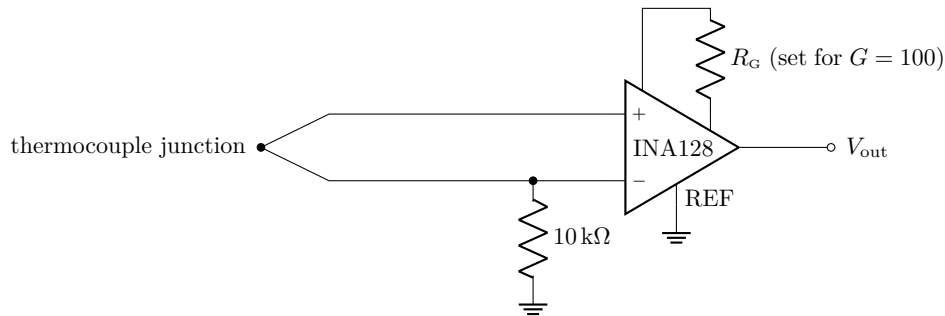
$$G = 1 + \frac{50\text{ k}\Omega}{R_G}. \quad (7.36)$$

The CMRR is 86 dB at unity gain ( $R_G = \infty$ ), and 125 dB at  $G = 100$ . The input impedance is  $\sim 10^{10}\ \Omega$ , and the input bias current is 2 nA. These cost  $\sim \$8$  each, depending on the grade (quality) and package.

## 7.5.2 Instrumentation-Amplifier Applications

### 7.5.2.1 Thermocouple Amplifier

The instrumentation amplifier turns out to be really useful in a number of applications, particularly in amplifying high-source-impedance sensors that require high gain. On example, shown below, is a simplistic thermocouple amplifier.



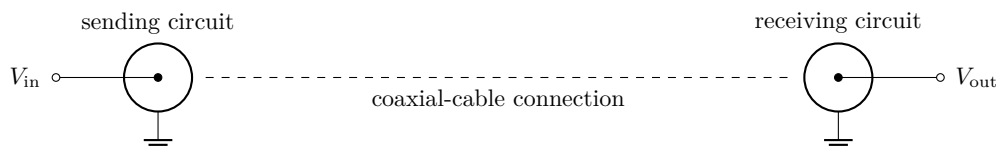
A **thermocouple** is a junction of two dissimilar metals that develops a voltage that is related (in a nonlinear way) to the temperature. Typical signal levels from thermocouples range from  $\sim 10\ \mu\text{V}$  to  $100\ \mu\text{V}$ , so significant gain is useful here. Note that this circuit is simplistic in the sense that it is only useful for measuring *relative* temperature changes. Absolute calibration requires more information, because the thermocouple wires must make a transition to copper on the way to the amplifier, and these copper junctions (and their temperature) contribute to the output voltage. The underlying physical mechanism is the **Seebeck effect**: a temperature gradient over the length of a wire (of a *single* material) causes charges to migrate creating a potential difference from one end to the other. The voltage is material dependent, so using two different wires in the thermocouple junction gives information about the temperature *difference* between the

thermocouple junction at one end and the junctions to copper at the other end. If the temperature of the copper-junction block is fixed (e.g., by controlling its temperature with an ice bath), then the thermocouple voltage is sufficient to determine the absolute temperature of the thermocouple junction. Otherwise, the mere knowledge of the copper-junction temperature is sufficient to deduce the thermocouple-junction temperature, although the relationship between voltage and thermocouple temperature is more complicated. The removal of the dependence on the copper-junction temperature is called **cold-junction compensation**, and there are special amplifiers that can emulate the cold junction electronically—a good example is the AD594.

One other thing to note is the 10-k $\Omega$  resistor, which is necessary for the circuit to function. Why?

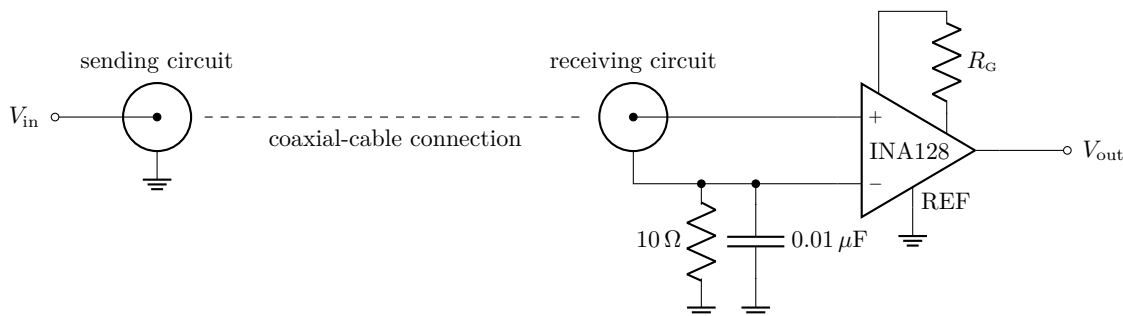
### 7.5.2.2 Differential Transmission for Noise Rejection

Another good application of the instrumentation amplifier is as a receiving input for a signal that is transferred between two instruments. For low noise, the signal is usually sent via **coaxial cable**—a center conductor shielded by a cylindrical outer conductor, which is usually used as a ground connection, so the grounded jacket protects the center (signal) conductor from external interference. The simplest way to send a signal between instruments is just to connect it as shown below—use a jack on either instrument, with the coaxial cable in between, the outer conductor grounded on either end.



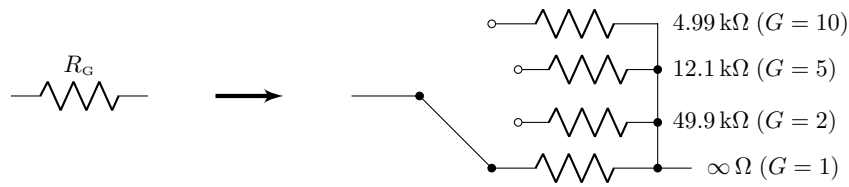
This is a bad idea, though, because it introduces a **ground loop**: both instruments are generally grounded via their power supplies, and the grounds are *also* connected via the coaxial-cable jacket. This means that the ground has a big loop for a path, and changing electromagnetic fields can induce an EMF in this loop. The symptoms are noise pickup, such as 60-Hz buzz in audio systems, or radio-frequency interference in wide-bandwidth circuits.

One solution is to power one instrument from a battery, so that it “floats” (i.e., there is no ground connection). This is not always convenient, but it is also possible to use an isolation transformer to break the ground connection via the power supply. (Isolation transformers for this purpose are commercially available, but usually require a safety ground connection to be defeated before it truly provides ground isolation.) However, having instruments grounded is otherwise desirable for safety and noise immunity, so another solution is to use an instrumentation amplifier as a differential receiver on the receiving instrument, as shown below.



This also breaks the ground loop, and any induced interference, which is common to the ground and signal conductors in the cable, will be cancelled in the subtraction. (This also works for twisted-pair cable, in place of the coaxial cable, as is usually used in Ethernet networks.) Note that the “ground” conductor on the received is still tied to ground via a resistor and capacitor. This prevents large induced input swings, in order to protect the amplifier (and remember we need the resistor to provide a dc path to ground).

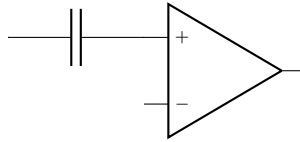
Another useful trick is to use the instrumentation amplifier to provide variable gain on the input, which is easy to accomplish by replacing  $R_G$  by a switch-selectable array of gain resistors (e.g., using a rotary switch to select gain), as shown below.



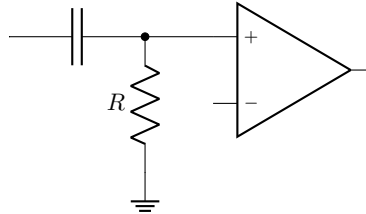
A logarithmic spacing of gains (1, 2, 5, 10, 20, 50, 100, ...) covers a wide gain range in a useful way.

### 7.5.2.3 AC-Coupled Inputs with High Impedance

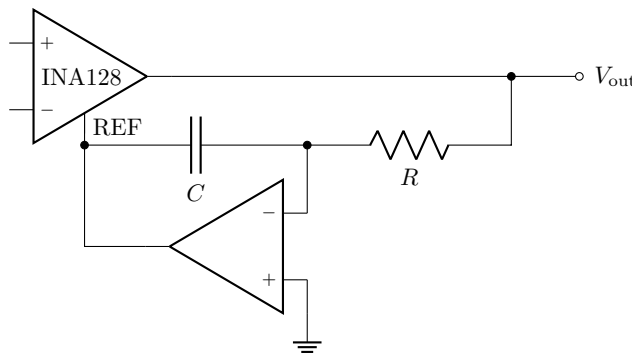
Sometimes it is useful to have an ac-coupled input, but with high input impedance. For example, we may have a sensor with a weak signal and high source impedance that requires large amplification, but if it has a large dc bias, a large gain would take the dc bias out of range. If we just use a capacitor to block dc on the input, this is bad, remember, because we need a dc path to ground on each input.



An improvement is to introduce a resistor, which makes a high-pass filter. The resistor prevents problems from input bias currents.



However, in the pass band, the input impedance is limited to  $R$ , which will not be nearly as good as a “bare” input on a decent op-amp. A clever solution is to use an integrator to feed back the output of an instrumentation amplifier to the reference input.



Assume for simplicity that the instrumentation amplifier is set for unity gain, so

$$V_{\text{out}} = V_+ - V_- + V_{\text{REF}}. \quad (7.37)$$

That is, the reference voltage is the reference for ground, so the amplifier sets  $V_{\text{out}} - V_{\text{REF}}$  to the difference  $V_+ - V_-$ . Note that the only steady state occurs when the integrator is at steady state, which is when the integrator input ( $V_{\text{out}}$ ) is zero. If  $V_{\text{out}} \neq 0$ , the integrator will build up a voltage until it cancels the steady output voltage on a time scale  $RC$  (with negative feedback due to the inverting nature of the integrator).

High-frequency signals are not affected by the integrator, because the integrator has no time to “catch up” to the rapidly changing output (which is saying that the integrator gain is suppressed at high frequencies as  $1/\omega$ ). To make this input switchable between ac and dc, a switch can short across the capacitor to change the instrumentation amplifier to dc mode.

## 7.6 Practical Considerations

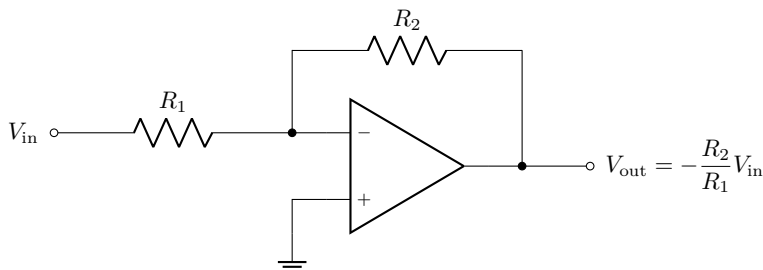
In this section, we will deal some more with some deviations of op-amps from their ideal behavior, and a few common tricks to mitigate these effects.

### 7.6.1 Input-Bias Currents and Precision Amplifiers

Before, when we were dealing with integrating amplifiers (Section 7.4.4), we talked about input-bias currents and how they can cause problems by charging up the integrating capacitor, even with a zero input voltage. We also talked about how it is critical to have a dc path to ground for each input to prevent similar charging problems. However, input bias currents can *still* cause problems (albeit usually less serious) in “regular” op-amp circuits like inverting and noninverting amplifiers.

#### 7.6.1.1 Inverting Amplifier

For example, let’s return to the inverting amplifier (Section 7.3.2).



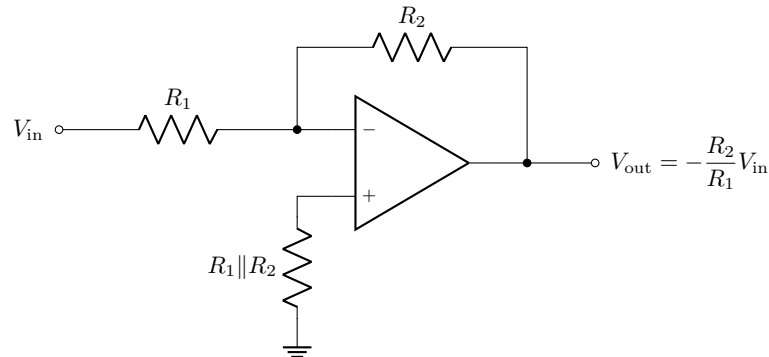
As a numerical example, consider a 741C (a lower-grade op-amp that will make such problems obvious), with a worst-case input bias current of  $0.5\ \mu\text{A}$ . If we take  $R_1 = 10\ \text{k}\Omega$  and  $R_2 = 1\ \text{M}\Omega$  (for  $G = -100$ ), then the inverting input sees a Thévenin-equivalent input resistance of  $R_{\text{Th}} = R_1 \parallel R_2 \approx 10\ \text{k}\Omega$ . The input bias current then leads to a bias voltage of  $(0.5\ \mu\text{A})(10\ \text{k}\Omega) = 5\ \text{mV}$  at the input. With a gain of  $(-100)$ , this leads to an output bias voltage of  $5\ \text{V}$  worst-case, which is pretty bad! (That is, for  $V_{\text{in}} = 0$ , it is permissible according to the 741C’s specs to have  $|V_{\text{out}}|$  as much as  $5\ \text{V}$ .)

How do we get around this? There are a few approaches:

- Since this is a dc-bias issue, this is no problem for *ac* circuits—just make sure to design for zero dc gain in the circuit (using a blocking capacitor, for instance).
- The main problem in the example was an underwhelming performer of an op-amp (by modern standards, anyway). Of course, we can do better, and recall that FET-input op-amps are superior to BJT-input op-amps in terms of bias currents. Precision BJT-input op-amps are a good option, too. For example, the FET-input LF411 has  $I_{\text{bias}} \approx 0.2\ \text{nA}$ , so the output error is only  $0.2\ \text{mV}$  in the above example.
- For BJT-input op-amps, or in high-precision circuits, we can also reduce errors by balancing input impedances, which is best illustrated by a few examples.

### 7.6.1.2 Balanced Input-Impedances: Inverting Amplifier

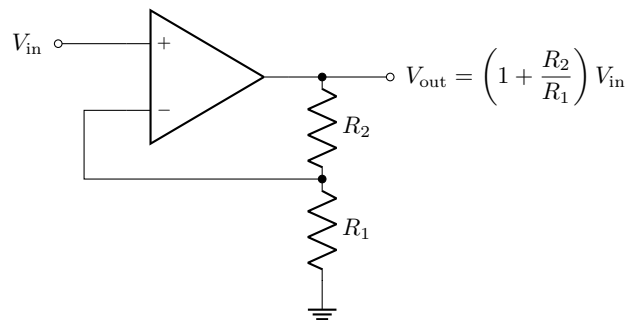
To balance the input impedances of the inverting amplifier, recall that the problem came from the inverting input “seeing” an effective source impedance of  $R_1 \parallel R_2$ . We can simply insert an equivalent source impedance on the *noninverting* input as shown below.



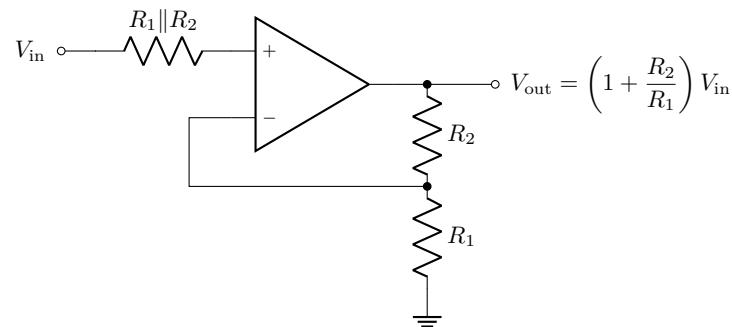
This resistor does nothing for an ideal op-amp, because no current flows through the resistor (and thus the resistor drops no voltage), but for a *real* op-amp, this trick reduces bias errors. Of course, all this assumes  $V_{in}$  acts like an ideal voltage source (i.e., it has a source impedance much smaller than  $R_1$ ), otherwise the source impedance of  $V_{in}$  must also enter into the compensation scheme.

### 7.6.1.3 Balanced Input-Impedances: Noninverting Amplifier

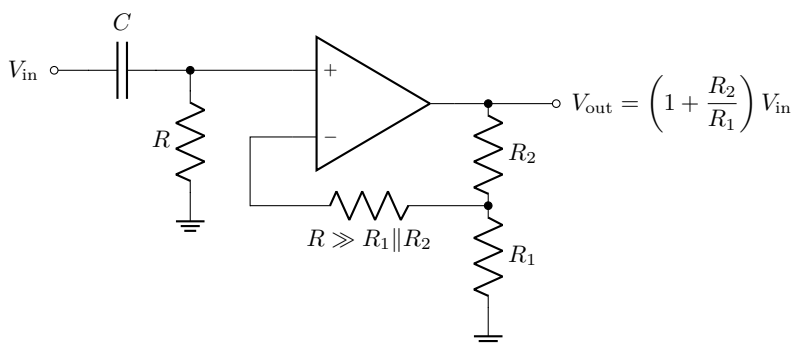
Similar balancing tricks are possible for noninverting amplifiers. Recall the basic noninverting amplifier (Section 7.3.3), shown below.



Here, the inverting input sees a source impedance of  $R_1 \parallel R_2$ , due to the voltage divider. Thus, it is again a simple matter to insert an equivalent source impedance for the other input.



Balancing input impedances is a bit trickier for ac-coupled amplifiers. A good example is shown below.



A high-pass filter is used on the input to block the dc component (remember the resistor  $R$  is necessary to provide a dc path to ground for the input bias current). We need to balance impedance at dc, so the noninverting input sees a source impedance of  $R$ . An easy way to ensure the same impedance for the inverting input is to insert another resistor  $R$  in the feedback loop, and simply ensure  $R \gg R_1 \parallel R_2$ , so we don't have to worry about  $R_1$  and  $R_2$  when setting the resistances.

Impedance-balancing tricks like this, by the way, only work for **voltage-feedback op-amps**—this applies to most available op-amps—because the input-bias currents are similar in such devices. In **current-feedback op-amps** (see the footnote on p. 212), however, the input-bias currents won't be similar, so this trick doesn't work.

#### 7.6.1.4 Input Offset Currents

In addition to input bias current, there is **input offset current**, which is basically the difference between the input bias currents for the two inputs of an op-amp. While input bias current is more or less intrinsic to an op-amp's design, input offset current is due to manufacturing asymmetry. Thus, even with balanced inputs, there will be some bias signal, but much smaller than in the unbalanced case. For example, the 741C has an input offset current of 200 nA (compared to  $0.5 \mu\text{A}$  for input bias current), and the LF411 has an input offset current of 0.1 nA (compared to 0.2 nA for input bias current).

#### 7.6.1.5 Common-Mode Rejection Ratio

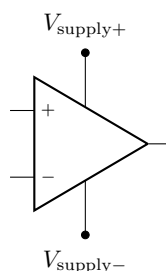
At this point, we can return to another source of error, the **common-mode rejection ratio (CMRR)**, which we introduced in Section 4.10.3 for the transistor differential amplifier, and we discussed it again for instrumentation amplifiers in Section 7.5. The typical CMRR range of op-amps is around 50–125 dB. But now look back at the op-amp inverting and noninverting amplifiers above. The *inverting* amplifier is *insensitive* to bad CMRR, compared to the *noninverting* amplifier, which is relatively *sensitive* to common-mode signals. Why is this? (What are the op-amp input voltages in each case?)

So while the noninverting amplifier looks good because of its very-high-impedance input, it is not quite as precise as the inverting amplifier. The latter is better in high-precision applications.

### 7.6.2 Power Supplies

So far, we haven't talked so much about the power-supply connections of op-amps. We have talked about the rule that no current flows in or out of the inputs (or at least there is only a very small current). Clearly, a current must flow in or out of the output in order for interesting things to happen. This output current must come from somewhere, and that is where the power-supply connections come in.

Op-amps have two supply connections, and they are often powered from **split supplies** of  $\pm 15\text{V}$ . More generally we will call these supply voltages  $V_{\text{supply}\pm}$ , and the explicit connections are shown below.



There are also **single-supply** op-amps, where  $V_{\text{supply}-}$  can be ground. Usually these can *also* be powered by split supplies, but the difference is as follows. Op-amps involve transistors, so the output can usually only swing to within a volt or two of either power-supply rail. Single-supply op-amps are designed with outputs that can swing all the way to the negative rail (so the output can swing to ground if in a single-supply circuit).

### 7.6.2.1 Power-Supply Rejection

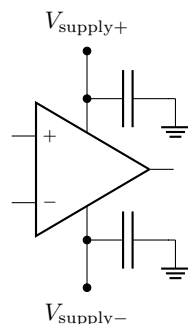
Ideally, the behavior of an op-amp is completely independent of the power supplies. That is, suppose the ideal op-amp output is 5 V, and the op-amp is powered from  $\pm 15$  V. Then suppose the power supply changes to  $\pm 16$  V. The output should be determined only by the circuit inputs and feedback network, and so shouldn't change at all, but in reality, it will change slightly, say to 5.1 V, to fabricate an example.

The (in)sensitivity of an op-amp to power-supply fluctuations is characterized by the **power-supply rejection ratio (PSRR)**. Ideally, the PSRR is very large, meaning the op-amp effectively “rejects” fluctuations in the power supply. The PSRR is defined with respect to the op-amp *inputs*, and is the ratio of the change in the power supply to the corresponding effect on the op-amp, referenced to a change at the input. (This accounts for the fact that power-supply fluctuations will have larger effects on the output for circuits with high gain.) As an example, suppose we have a PSRR of 120 dB, which is a ratio of  $10^6$  in voltage. Then a 1-V change in the power supply corresponds to a 1- $\mu$ V change *at an input* to the op-amp. We must then get into the specific connections of the op-amp circuit (i.e., the gain) to determine the change in the output voltage. (So a unity-gain follower would also see a 1- $\mu$ V change *at the output*.)

As a real example, the LF411 has a PSRR at dc of 100 dB typical, 80 dB minimum. The PSRR is better for the + supply than for the - supply, and the PSRR is worse at higher frequencies, dropping to  $\sim 90$  dB at 100 Hz, and  $\sim 30$  dB at 100 kHz.

### 7.6.2.2 Power-Supply Bypass Capacitors

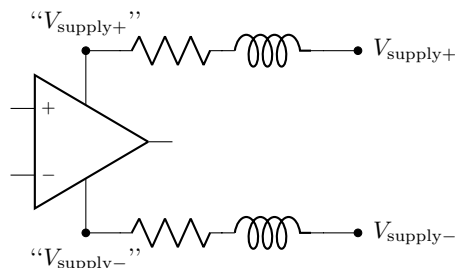
One common trick for improving the behavior of op-amp circuits is to use “bypass capacitors” on the power-supply leads of op-amps. The basic connection is shown schematically below.



The values of these bypass capacitors are not critical, but typically they would be 0.01- $\mu$ F (monolithic) ceramic, 0.1- $\mu$ F (monolithic) ceramic, or 1- $\mu$ F tantalum (polarized) capacitors. These capacitors should also be placed physically as close as possible to the op-amp power-supply pins.



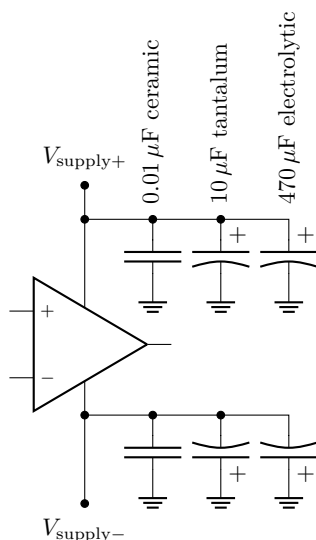
Why use bypass capacitors? In a circuit, for example, on a **printed circuit board (PCB)**, there may be long wires (or PCB traces) connecting the op-amp supply pins to the power supplies. Schematically this is shown below, where the wires have some intrinsic resistance and inductance.



The inductance is particularly problematic, as it means the path to the power supply has high impedance. The capacitor acts to short-circuit, or bypass, the inductance of the power-supply lead by providing a low-impedance path to ground at high frequencies. Otherwise, what can happen is as follows. A sudden change in the output of an op-amp (in response to an input change) implies a quick change in the power-supply currents. Inductance in the power-supply leads means that the voltage will drop if, for example, the op-amp is suddenly demanding more power-supply current. As we have seen, this can lead to output inaccuracies, because the op-amp PSRR is worse at high frequencies. In the worst case, the op-amp may even self-oscillate.

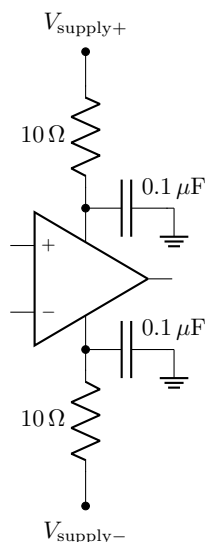
Another way to think of this is at the capacitor acts as a “charge reservoir” that tries to stabilize the power-supply voltages, temporarily supplying extra current as necessary when the op-amp demands it. Again, this is most effective at high frequencies, where the PSRR is bad anyway. Thus we also see the importance of using small, ceramic or tantalum capacitors (which have low inductance and respond well at high frequencies), and placing the capacitors very close to the op-amp (to bypass as much of the lead inductance as possible). On a PCB, each op-amp has its own local bypass capacitors, and the connection to ground should have very low impedance (e.g., to a **ground plane**, or a grounded copper layer that covers most of one side or layer of a PCB).

For critical applications, for example for a high-current amplifier, an even better approach is to use multiple, parallel bypass capacitors on each power-supply pin, as shown below.



Here the large 470-μF capacitor acts as a large charge reservoir, as appropriate for the high-current circuit, but only responds well at relatively low frequencies due to a high intrinsic inductance. Progressively smaller capacitors with lower inductance will help stabilize the power-supply voltage at higher frequencies. In this case, the smallest capacitors should be located closest to the op-amp, with the location of the large capacitors not so critical.

Another good approach for a low-current, or low-level, precision op-amp circuit is as shown below.



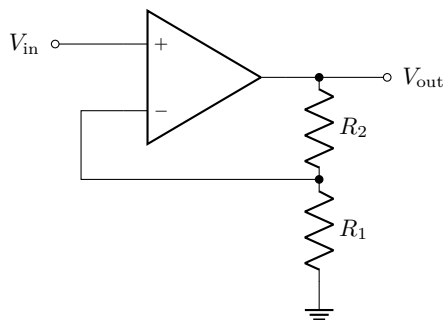
The idea is to add extra resistance to the power-supply lines, before the bypass capacitor. Obviously, this trick is limited to low-current circuits, so the voltage drop due to the  $10\text{-}\Omega$  resistor is small. This resistor enhances the effect of the bypass capacitor and gives improved isolation from power-supply fluctuations. This is an especially useful technique if high- and low-current amplifiers coexist in the same circuit, where feedback from the high-current to low-current amplifiers may cause the circuit to self-oscillate. Even better is to make sure the high-current and low-current ICs do not share the same power-supply lines (they should use separate lines, connected only at a point where the supply voltage is regulated, for example, by a 3-terminal regulator).

## 7.7 Finite-Gain Analysis

So far, we have assumed that the *open-loop* gain  $A$  of the op-amp is arbitrarily large. We have mainly made this assumption via the rule that the input voltages are the same in negative-feedback mode. In a *real* op-amp, this open-loop gain is high, but finite, ranging from about  $\sim 50\text{--}146\text{ dB}$ . So now let's relax the assumption of infinite open-loop gain to (1) see what the effects are, and (2) to see why having a large gain is a good thing in an op-amp.

### 7.7.1 Noninverting Amplifier

As a first example, let's return to the noninverting-amplifier circuit of Section 7.3.3, shown below.



Recall that the result with infinite open-loop gain was

$$V_{\text{out}} = \left(1 + \frac{R_2}{R_1}\right) V_{\text{in}}. \quad (7.38)$$

But now, we will use the op-amp rule<sup>2</sup>

$$V_{\text{out}} = A(V_{\text{in}+} - V_{\text{in}-}), \quad (7.39)$$

where  $A$  is the open-loop gain, instead of just assuming  $V_{\text{in}+} = V_{\text{in}-}$ . To simplify notation a bit, let's define the voltage-divider fraction

$$\eta := \frac{R_1}{R_1 + R_2}, \quad (7.40)$$

so that the divided output voltage  $\eta V_{\text{out}}$  is fed back to the inverting input,

$$V_{\text{in}-} = \eta V_{\text{out}}. \quad (7.41)$$

Also, we have the input

$$V_{\text{in}} = V_{\text{in}+}, \quad (7.42)$$

so putting these equations into the op-amp rule (7.39), we have

$$V_{\text{out}} = A(V_{\text{in}} - \eta V_{\text{out}}). \quad (7.43)$$

Solving for the output voltage, we have

$$V_{\text{out}} = \frac{A}{1 + \eta A} V_{\text{in}},$$

(noninverting amplifier, finite open-loop gain) (7.44)

which defines the *closed-loop* gain

$$G = \frac{A}{1 + \eta A}.$$

(closed-loop gain, noninverting amplifier with finite open-loop gain) (7.45)

Note that in the limit  $\eta A \gg 1$ , this expression reduces to the original formula

$$G_{\infty} = \eta^{-1} = 1 + R_2/R_1, \quad (7.46)$$

so this analysis reproduces the ideal-op-amp limit.

### 7.7.1.1 Gain Limits and Error

Notice that  $G \leq A$ , so that the *open-loop* gain  $A$  limits the *closed-loop* gain  $G$ —negative feedback can only *reduce* the gain. Additionally,  $G < G_{\infty}$ , so the *ideal* closed-loop gain always limits the *real* closed-loop gain. More specifically, assuming  $A$  to be large ( $A \gg G_{\text{infty}}$ ),

$$G = \frac{A}{1 + \eta A} = \frac{\eta^{-1}}{1 + (\eta A)^{-1}} = \frac{G_{\infty}}{1 + G_{\infty}/A} \approx G_{\infty} \left(1 - \frac{G_{\infty}}{A}\right). \quad (7.47)$$

Then the fractional “error” in the gain is

$$\frac{\delta G}{G_{\infty}} \approx -\frac{G_{\infty}}{A}. \quad (7.48)$$

---

<sup>2</sup>This voltage rule applies to **voltage-feedback op-amps**, and thus applies to *most* op-amps. There is another class of op-amps called **current-feedback op-amps**, which are sensitive to a small input current between the inputs. That is, Eq. (7.39) is replaced by  $V_{\text{out}} = Z_A I_{\text{in}}$ , where  $I_{\text{in}}$  is the input bias current flowing from the + input to the – input of the op-amp, and  $Z_A$  is the (large) open-loop transimpedance. For more information, see James Karki, “Voltage Feedback Vs Current Feedback Op Amps,” Texas Instruments Application Report (1998) <http://www.ti.com/lit/an/slva051/slva051.pdf>.

As an example, a decent op-amp has  $A$  of 100 dB, which corresponds to  $A = 10^5$ . For a  $G_\infty = 10$  setup, the fractional error is

$$\frac{\delta G}{G_\infty} \approx -\frac{10}{10^5} = -10^{-4} = -0.01\%, \quad (7.49)$$

which is pretty small. For reasonably high values of  $A$ , this error is usually negligible compared to the error due to the feedback-resistor tolerances.

### 7.7.1.2 Insensitivity to Gain Variation

Another handy result that we obtain for large  $A$  is that if  $A$  is sufficiently large, then  $G$  is insensitive to variations in  $A$ . Starting with the closed-loop-gain expression above in terms of  $G_\infty$  and  $A$ ,

$$G = \frac{G_\infty}{1 + G_\infty/A}, \quad (7.50)$$

then

$$\frac{\partial G}{\partial A} = -\frac{G_\infty}{(1 + G_\infty/A)^2} (-G_\infty/A^2) = \frac{G}{A^2} \frac{G_\infty}{(1 + G_\infty/A)} = \frac{G}{A} \frac{1}{(1 + A/G_\infty)}. \quad (7.51)$$

Then the variation  $\delta G$  in the closed-loop gain  $G$  is

$$\frac{\delta G}{G} = \frac{\partial G}{\partial A} \delta A, \quad (7.52)$$

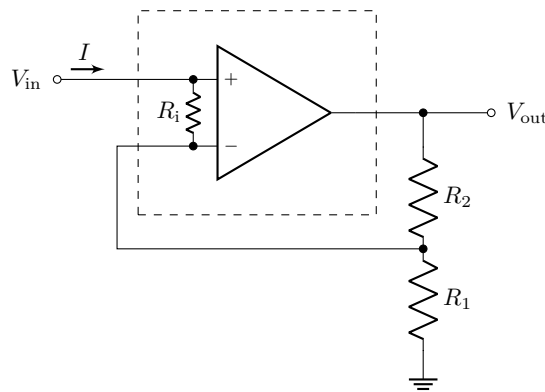
so

$$\frac{\delta G}{G} = \frac{1}{(1 + A/G_\infty)} \frac{\delta A}{A}, \quad (7.53)$$

so the fractional error  $\delta G/G$  is much smaller than the fractional open-loop variation  $\delta A/A$  by a factor of  $1/(1 + A/G_\infty) \approx G_\infty/A$  for large  $A$ . This is a nice property because, for example, feedback with large  $A$  reduces variations of gain with frequency, for a flatter response in a closed-loop amplifier. It also reduces nonlinearity and distortion, which you can roughly think of as variations in gain with signal amplitude.

## 7.7.2 Feedback and Input Impedance

Negative feedback with large open-loop gain also helps quite a bit with input and output impedance. Going back to the noninverting amplifier, we can construct an explicit model for input impedance  $R_i$  as shown below.



The dashed line encompasses the “real” amplifier, which consists of an ideal op-amp and a resistor modeling the input impedance. As before, the closed-loop gain (7.45) is

$$G = \frac{A}{1 + \eta A}. \quad (7.54)$$

Since,  $V_{in+} = V_{in}$ , and  $V_{in-} = \eta V_{out}$ , we can take the input current  $I$  to be

$$I = \frac{V_{in+} - V_{in-}}{R_i} = \frac{V_{in} - \eta V_{out}}{R_i} = \frac{1 - \eta A / (1 + \eta A)}{R_i} V_{in} = \frac{V_{in}}{(1 + \eta A) R_i}. \quad (7.55)$$

Then the effective input impedance  $Z_{in}$  is

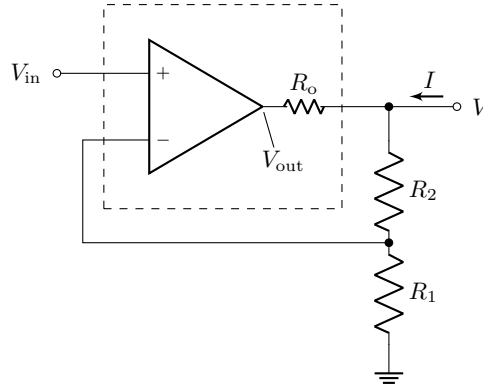
$$Z_{in} = (1 + \eta A) R_i = \left(1 + \frac{A}{G_{\infty}}\right) R_i, \quad (\text{input impedance, noninverting amplifier}) \quad (7.56)$$

which is much larger than the intrinsic input impedance  $R_i$  if there is large open-loop gain  $A \gg G_{\infty}$ .

For example, the modest 741C has  $R_i = 2 \text{ M}\Omega$  typical,  $300 \text{ k}\Omega$  minimum, which is not great. The open-loop gain  $A$  is typically  $2 \times 10^5$ , minimum  $1.5 \times 10^4$ , which is not bad. If  $G_{\infty} = 10$ , then  $Z_{in} = 4 \times 10^{10} \Omega$  typical, or  $5 \times 10^8 \Omega$  minimum. These are pretty high input impedances, and they can be much higher with a precision op-amp.

### 7.7.3 Feedback and Output Impedance

To model the effects of feedback on the output impedance, we can again introduce an explicit model, including an output resistance  $R_o$ .



Again, the dashed box represents the “real” amplifier, with an ideal op-amp and the output resistor. We will call the output of the ideal amplifier  $V_{out}$ , while the “real” output is  $V$ . There is also a current  $I$ , which we define as flowing *into* the output. Then setting a null input  $V_{in} = 0$ , we have  $V_{in+} = 0$ , and now  $V_{in-} = \eta V$ , so the op-amp rule (7.39) gives

$$V_{out} = -\eta A V. \quad (7.57)$$

Then the current is

$$I = \frac{V - V_{out}}{R_o} = \frac{V(1 + \eta A)}{R_o}. \quad (7.58)$$

Thus, the output impedance  $Z_{out} = V/I$  is

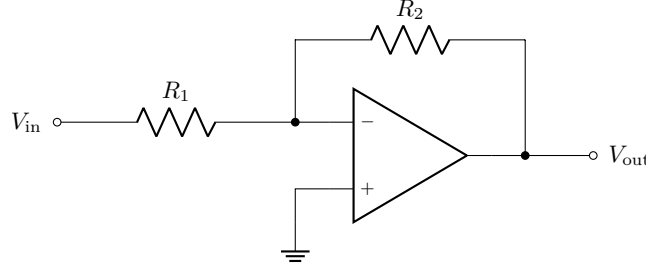
$$Z_{out} = \frac{R_o}{1 + \eta A} = \frac{R_o}{1 + A/G_{\infty}}. \quad (\text{output impedance, noninverting amplifier}) \quad (7.59)$$

This should be much smaller than the intrinsic output impedance  $R_o$ , provided we have large open-loop gain,  $A \gg G_{\infty}$ .

For example, the modest 741C has  $R_o = 75 \Omega$ . With a typical open-loop gain  $A$  of  $2 \times 10^5$ , and  $G_{\infty} = 10$ , then  $Z_{out} = (75 \Omega) / (2 \times 10^4) \approx 4 \text{ m}\Omega$ , which is quite small.

### 7.7.4 Circuit Practice: Finite Gain in the Inverting Amplifier

For practice in dealing with finite op-amp gain, consider the noninverting amplifier, with finite open-loop gain  $A$ . (Again, it's best to do this before continuing, so we won't defer this until the end of the chapter.)



Show the following:

(a) The finite- $A$  gain is

$$G = -\frac{(1 - \eta)A}{1 + \eta A}, \quad \eta := \frac{R_1}{R_1 + R_2}.$$

(closed-loop gain, inverting amplifier with finite open-loop gain) (7.60)

(b) Take the  $A \rightarrow \infty$  limit of Eq. (7.60) and show that

$$G_\infty = -\frac{R_2}{R_1}.$$

(closed-loop gain, inverting amplifier with infinite open-loop gain) (7.61)

(c) The input impedance is

$$Z_{\text{in}} = R_1 + \frac{R_2}{1 + A}.$$

(input impedance, inverting amplifier with finite open-loop gain) (7.62)

(Ignore any intrinsic input impedance of the op-amp, which we will assume is much larger than  $R_1$ .)

(d) The output impedance is

$$Z_{\text{out}} = \frac{R_o}{1 + \eta A},$$

(input impedance, inverting amplifier with finite open-loop gain) (7.63)

where  $R_o$  is the intrinsic output impedance of the op-amp, as in the noninverting case.

#### Solution.

(a) First, the noninverting input has  $V_{\text{in}+} = 0$ . The inverting input has a voltage determined by a voltage divider between  $V_{\text{in}}$  and  $V_{\text{out}}$ :

$$V_{\text{in}-} = \eta V_{\text{out}} + (1 - \eta)V_{\text{in}}. \quad (7.64)$$

Remember  $\eta = R_1/(R_1 + R_2)$ , so as a sanity check,  $\eta \rightarrow 1$ ,  $V_{\text{in}-}$  becomes connected to  $V_{\text{out}}$ , and as  $\eta \rightarrow 0$ ,  $V_{\text{in}-}$  becomes connected to  $V_{\text{in}}$ , which makes sense. Then using Eq. (7.39),

$$V_{\text{out}} = -\eta A V_{\text{out}} - (1 - \eta)A V_{\text{in}}. \quad (7.65)$$

Solving for  $V_{\text{out}}$ ,

$$V_{\text{out}} = -\frac{(1 - \eta)A}{(1 + \eta A)} V_{\text{in}}. \quad (7.66)$$

This is the result we wanted, with  $G$  the coefficient of  $V_{\text{in}}$ .

(b) As  $A \rightarrow \infty$ ,  $G = -(1 - \eta)A/(1 + \eta A) \rightarrow -(1 - \eta)/\eta = -R_2/R_1$ .

(c) Suppose a current  $I$  flows into the  $V_{\text{in}}$  terminal. Then

$$I = \frac{V_{\text{in}} - V_{\text{in}-}}{R_1} = \frac{V_{\text{in}} - \eta V_{\text{out}} - (1 - \eta)V_{\text{in}}}{R_1} = \frac{\eta(V_{\text{in}} - V_{\text{out}})}{R_1} = \frac{\eta[(1 + \eta A) + (1 - \eta)A]}{(1 + \eta A)R_1} V_{\text{in}} = \frac{\eta(1 + A)}{(1 + \eta A)R_1} V_{\text{in}}, \quad (7.67)$$

where we used the solution (7.66). Then  $Z_{\text{in}} = V_{\text{in}}/I$ , so

$$Z_{\text{in}} = \frac{(1 + \eta A)R_1}{\eta(1 + A)} = \frac{R_1}{\eta(1 + A)} + \frac{AR_1}{(1 + A)} = \frac{R_1 + R_2}{(1 + A)} + \frac{AR_1}{(1 + A)} = \frac{R_2}{(1 + A)} + R_1. \quad (7.68)$$

Note that this reduces to  $R_1$  as  $A \rightarrow \infty$ .

(d) Here, we set  $V_{\text{in}} = 0$  and call the output  $V$ , with a current  $I$  going into the output terminal.  $V_{\text{out}}$  is the output voltage of the ideal op-amp, before the intrinsic resistor  $R_o$ , as in the noninverting case. Then  $V_{\text{in}-} = \eta V$ , and

$$V_{\text{out}} = A(V_{\text{in}+} - V_{\text{in}-}) = -AV_{\text{in}-} = -\eta AV. \quad (7.69)$$

So the current is

$$I = \frac{V - V_{\text{out}}}{R_o} = \frac{1 + \eta A}{R_o} V, \quad (7.70)$$

and so the output impedance  $Z_{\text{out}} = V/I$  is

$$Z_{\text{out}} = \frac{R_o}{1 + \eta A}. \quad (7.71)$$

Note that this decreases to zero as  $A \rightarrow \infty$ .

## 7.8 Bandwidth

The **bandwidth** of an amplifier refers to the frequency range over which the response (gain) is reasonably flat. For electronic amplifiers, one characteristic is that the gain must fall off above some frequency—no amplifier can work at arbitrarily high frequencies.

Recall that the closed-loop gain  $G$  and the open-loop gain  $A$  are related, in that the latter bounds the former:

$$G \leq A. \quad (7.72)$$

Now let's consider the frequency dependence of the gain. In particular, the open-loop gain  $A(\omega)$  typically has a “one-pole response,” like that of a low-pass filter:

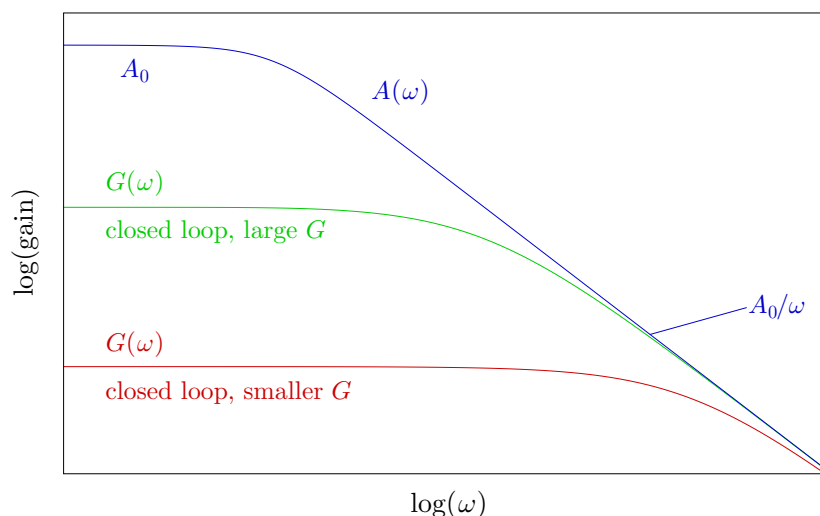
$$A(\omega) = \frac{A_0}{\sqrt{1 + (\omega/\omega_{3\text{dB}})^2}}. \quad (7.73)$$

Here, the cutoff frequency  $\omega_{3\text{dB}} = 1/RC$  for a low-pass filter, where  $R$  and  $C$  are typically set by intrinsic transistor resistance and stray or added (internal) capacitance in the op-amp. Note that asymptotically,

$$A(\omega) \sim \frac{1}{\omega} \quad (7.74)$$

for large  $\omega$ , for a scaling of  $-6\text{dB/octave}$ .

Then the closed-loop gain is bounded by the open-loop gain, so that as the open-loop gain falls off, so does the closed-loop gain. This is illustrated schematically below for two different dc gains  $G$  and an open-loop gain  $A(\omega)$ .



Note that as the dc gain becomes smaller, the bandwidth (frequency range over which the gain is roughly constant) becomes wider. Since  $A$  cuts off as  $\omega^{-1}$ , the closed-loop gain  $G(\omega)$  meets  $A(\omega)$  at a frequency that scales in the same way as  $\omega^{-1}$ . That is, the bandwidth scales as  $1/G_0$ , where  $G_0$  is the dc gain. Said differently, the product of the dc gain  $G_0$  and the bandwidth is a constant, and this is often quoted as the **gain–bandwidth product (GBWP)**, or the **unity-gain bandwidth**. For example, the 741C has a GBWP of 1.5 MHz. Generally speaking, op-amps tend to be slow, especially at high gains, compared to discrete transistors.

### 7.8.1 Slew Rate

A concept closely related to bandwidth is the **slew rate**, which is the maximum rate of change of the output. Intuitively, this should be proportional to the GBWP, but this is somewhat more complicated because the same signal, but with different amplitudes, would involve different slew rates, even if they have the same frequency spectrum. So for rapidly changing signals, an op-amp with a particular slew rate may be able to follow the signal at low amplitudes, but it may be harder for the op-amp to follow the same signal at larger amplitudes.

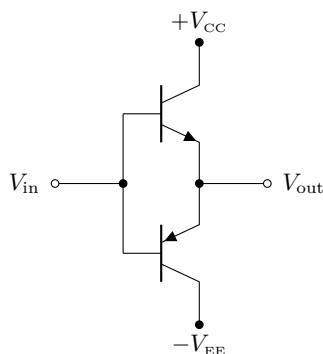
As a concrete example, the 741C has a modest slew rate of  $0.5 \text{ V}/\mu\text{s}$ . Slew rates can be much high; for example, the BUF634 unity-gain buffer has a slew rate of  $2000 \text{ V}/\mu\text{s}$ .

Generally speaking, the speed of an op-amp (either in terms of slew rate or GBWP) is controlled by the internal capacitance, which is usually fixed by an internal compensation capacitor, but also by intrinsic emitter resistance. Recalling that  $r_e \propto 1/I_C$ , generally speaking, a larger quiescent current (idling current) for an op-amp gives a higher slew rate or a wider GBWP. There is thus a trade-off between power and speed—some op-amps, like the OPA602, have a programmable quiescent current so the user can choose exactly where to make this trade-off.

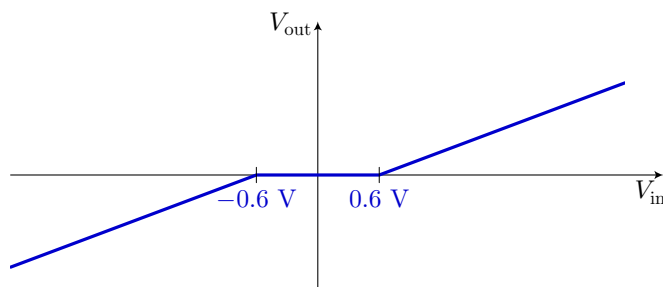
#### 7.8.1.1 Slew Rate and Power-Boosted Op-Amps

There are certain circuits where the slew rate of an op-amp is critical to its performance. One example is a “power-boosted” op-amp, where transistors are used to boost the output current capacity of an op-amp. The motivation for this circuit comes from the following “push-pull” current amplifier that we studied in Section 4.13.

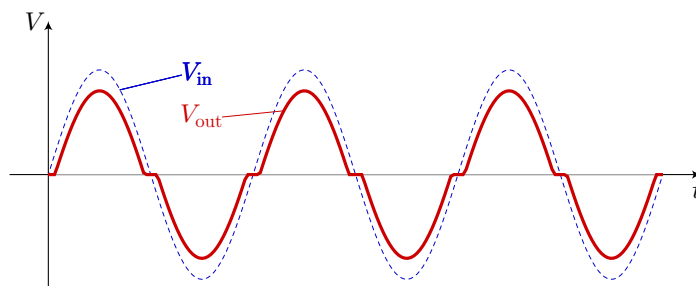




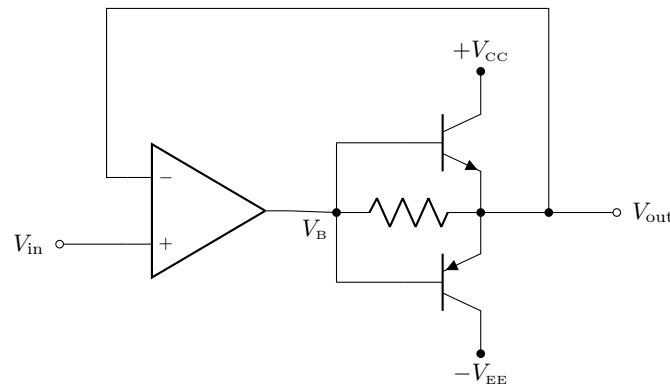
To review, this is basically a stacked pair of emitter followers. The potential advantage is operation with input signals of either polarity. The problem, though, is that one of the transistors will conduct, and the emitter (output) voltage must be a diode drop closer to zero than the base (input) voltage. That is, a graph of the output voltage responding to input voltage is schematically as in the graph below, if we assume the simple model that the base-emitter voltage drop is a constant 0.6 V (or less).



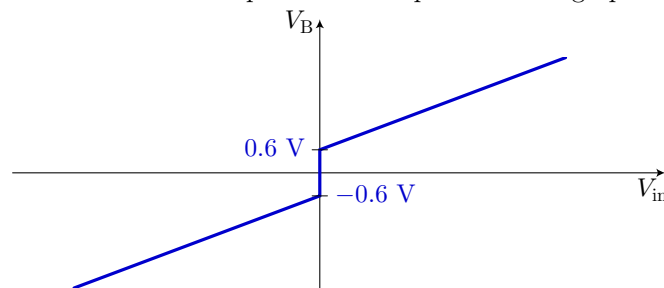
The problem in using this as an amplifier is that it leads to **crossover distortion**, because the base-emitter drop changes as the signal crosses through zero. An example is shown in the graph below of crossover distortion of an input sine wave.



One nice solution, at least in principle, is to use an op-amp, and enclose the push-pull transistor pair in the feedback loop of the op-amp, as in the circuit below. Note the addition of a base-emitter resistor (of the order 1 k $\Omega$  or so), which provides a feedback path when neither transistor is conducting—remember that omitting the negative-feedback path is a good way to annoy an op-amp.



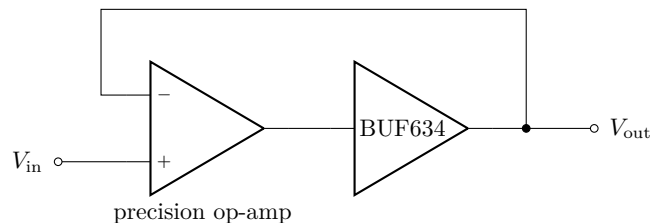
This circuit acts as a unity-gain buffer with high current-driving capacity, because the op-amp does whatever it needs to do to ensure that  $V_{out}$  is the same as  $V_{in}$ . And to do this, it must “undo” the crossover distortion, so the base voltage  $V_B$  in this circuit must respond to the input as in the graph below.



This response combined with the crossover distortion results in, in principle, a distortion-free output.

However, the problem with this conclusion is that it assumes that the op-amp has a long time to settle to the “correct” value. But with a rapidly changing input signal, the op-amp must jump discontinuously by 1.2 V when the input signal crosses through zero, which in practice can be problematic. This can lead to larger distortion and “glitching” with faster input signals.

The solution to this problem that we discussed before in Section 4.13 is to bias the transistors into conduction. However, the design of the bias-corrected circuit can be involved, especially to avoid thermal problems and proper selection and matching of bias voltages. A handier way (but with more limited options) to solve this problem is to use a high-current buffer amplifier, where engineers have already taken care of the effort of properly biasing the push-pull pair. One example is the circuit below, which can handle 250-mA output signals via a BUF634 unity-gain buffer.



The idea is to use a precision op-amp, and take the feedback from the output of the BUF634. This way, we get the high-current capacity of the “slave” BUF634, combined with the precision of the “master” op-amp. One caveat, which we will explore in more depth, is that the buffer amplifier must have a much wider bandwidth than the master op-amp.

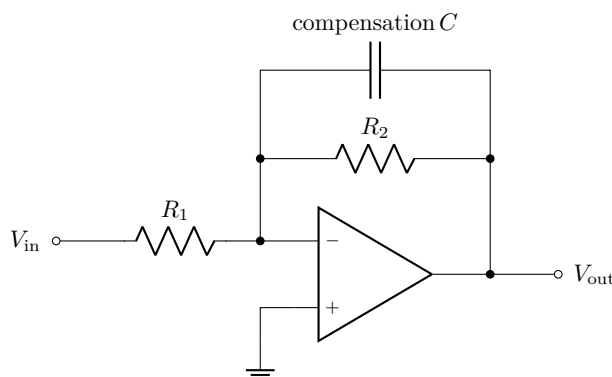
## 7.8.2 Stability and Compensation

So far, we have talked about the frequency-dependent gain, but the frequency-dependent phase is also critical. As we have noted, for most purposes we can regard an op-amp as having a gain behavior similar to a low-pass

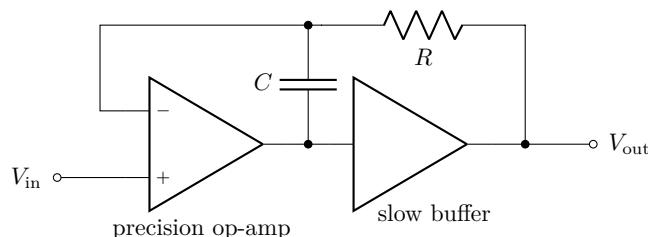
filter.

Recall that for an RC filter, the relative phase of the output is  $0^\circ$  in the limit  $\omega \rightarrow 0$ , and changes to a  $90^\circ$  phase lag as  $\omega \rightarrow \infty$ . If we have multiple, cascaded filters, at high frequencies, intuitively we can think of having a  $90^\circ$  phase **per “pole”**, or roughly speaking, per RC pair. This can pose a problem for op-amp circuits that require negative feedback. Due to phase shifts and time delays in the feedback loop at high frequencies can add up to a  $180^\circ$  phase shift. However, negative feedback in combination with a phase shift of  $180^\circ$  (or in fact anything between  $90^\circ$  and  $270^\circ$ ) is in fact *positive* feedback. This can turn into unstable behavior (oscillation) if the gain of the circuit exceeds unity in the frequency range where the feedback becomes positive feedback.

Thus comes the idea of **compensation**. Most op-amps include an internal capacitor to “roll off” the gain, and in particular to ensure that the gain is less than unity at frequencies where large phase shifts may cause problems. There are also **uncompensated op-amps**, which need an external capacitor or an appropriate reactive load to achieve stability. An example is the inverting amplifier below, with an explicit compensation capacitor to cut off the gain at high frequencies.



Another example is a variation on the BUF634 buffered op-amp circuit from the previous section. If the slave buffer amplifier is slower than the master op-amp, this may cause a problem because the phase shift due to the slower buffer may cause the master to become unstable or oscillate. A solution to use a slow buffer is the circuit below.



For small frequencies ( $\omega \ll \omega_{3\text{dB}} = 1/RC$ ), the feedback comes from the buffer output, while at high frequencies ( $\omega \gg \omega_{3\text{dB}} = 1/RC$ ), the feedback comes from the output of the master op-amp. This arrangement avoids problems with the phase shift and maintains stability of the amplifier.

### 7.8.2.1 Op-Amp Output and Capacitive Loads

To examine a problematic situation in a bit more detail, let's return to the bandwidth argument of Section 7.8, but now keep the complex phase in the open-loop gain:

$$\tilde{A}(\omega) = \frac{A_0}{1 - i\omega/\omega_{3\text{dB}}}. \quad (7.75)$$

Again, this response has the form of a low-pass filter [Eq. (2.44)], where the op-amp response “rolls off” due to a single capacitor. Then, taking the example of the noninverting amplifier of Section 7.7.3, the op-amp

output impedance (7.59) becomes

$$Z_{\text{out}}(\omega) = \frac{R_o}{1 + A(\omega)/G_\infty} = \frac{R_o(1 - i\omega/\omega_{3\text{dB}})}{1 + A_0/G_\infty - i\omega/\omega_{3\text{dB}}}. \quad (\text{output impedance, noninverting amplifier}) \quad (7.76)$$

Since typically  $A_0 \gg G_\infty$ , we can consider the intermediate range of frequencies between  $\omega_{3\text{dB}}$  and  $(A_0/G_\infty)\omega_{3\text{dB}}$ , where the impedance reduces to

$$Z_{\text{out}}(\omega) = \frac{R_o(-i\omega/\omega_{3\text{dB}})}{A_0/G_\infty} = -i\omega \frac{R_o G_\infty}{\omega_{3\text{dB}} A_0} \quad \left(1 \ll \frac{\omega}{\omega_{3\text{dB}}} \ll \frac{A_0}{G_\infty}\right). \quad (7.77)$$

This has the form of an inductive reactance [see Eq. (2.40)], with effective inductance

$$L_{\text{eff}} = \frac{R_o G_\infty}{\omega_{3\text{dB}} A_0}. \quad (7.78)$$

This inductive regime can span a wide range. With  $\omega_{3\text{dB}}/2\pi \sim 100$  Hz,  $R_o \sim 100 \Omega$ ,  $A_0 \sim 10^6$ , and  $G_\infty = 10$ , this works out to a frequency range of  $\sim 100$  Hz to  $\sim 10$  MHz, with an effective inductance of  $\sim 1.6 \mu\text{H}$ .

Outside this frequency band, the output impedance is simply resistive. For small frequencies, we have

$$Z_{\text{out}}(\omega) = \frac{R_o}{1 + A_0/G_\infty} \quad \left(\frac{\omega}{\omega_{3\text{dB}}} \ll 1 < \frac{A_0}{G_\infty}\right), \quad (7.79)$$

while for very large frequencies,

$$Z_{\text{out}}(\omega) = R_o \quad \left(1 < \frac{A_0}{G_\infty} \ll \frac{\omega}{\omega_{3\text{dB}}}\right). \quad (7.80)$$

Note that in practice, other stray capacitances will become important in the high-frequency range.

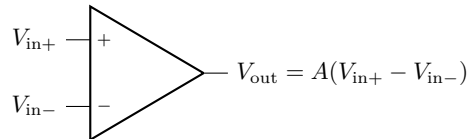
The inductive output of the op-amp in the regime of Eq. (7.77) can give rise to problems if the output of the op-amp drives a capacitive load. Then  $L_{\text{eff}}$  and the capacitance form a resonant circuit. This can lead to a resonant peak in the gain profile at frequency

$$\omega_0 = \frac{1}{\sqrt{L_{\text{eff}} C}} = \sqrt{\frac{\omega_{3\text{dB}} A_0}{R_o C G_\infty}}. \quad (7.81)$$

For the same parameters, this would lead to a resonance at  $\omega/2\pi \sim 100$  kHz for a load capacitance  $C = 1 \mu\text{F}$ . If the effect of the resonance is sufficiently strong (i.e., sufficiently large  $Q$ ), the circuit can become unstable and oscillate near this resonance frequency.

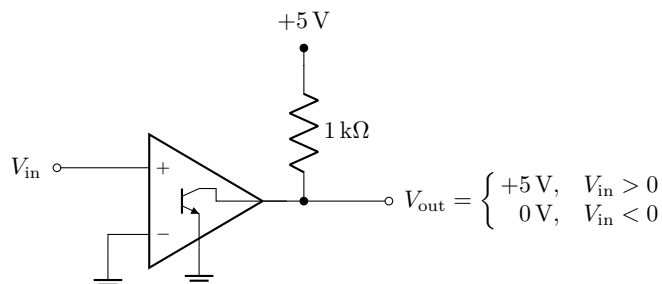
## 7.9 Comparators

Recall that op-amps are basically high-gain differential amplifiers.



We have mostly concentrated on closed-loop operation (feedback from output to the inverting input), which forces the inputs to have basically the same voltage. In open-loop operation (no feedback), the inputs are not the same, and if they are different by even a small amount ( $\sim \text{mV}$ ), the output rails one way or the other to reflect the difference. This open-loop operation is useful in some contexts, and op-amps that are specifically designed for this purpose are called **comparators**.

Specialized comparators (vs. using regular op-amps in the same role) have some advantages. For example, stability is not a concern, because comparators are not generally used with negative feedback. Thus, they need no compensation, and are instead optimized for very high slew rates. In fact, a common configuration for a comparator is the **open-collector output**. The common LM311 comparator, with open-collector output, is shown below, connected as in typical usage.



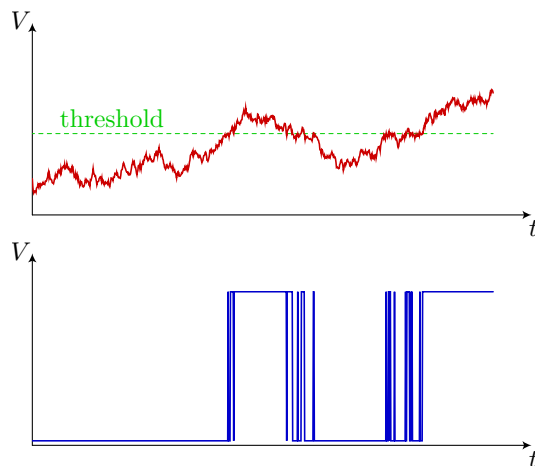
The usual output of the op-amp in the LM311 drives the base of an output transistor, whose collector is connected to the output. If  $V_{in} > 0$ , then the transistor acts as an open circuit, causing the output to go high to  $+5\text{ V}$  via the  $1\text{ k}\Omega$  pull-up resistor. If  $V_{in} < 0$ , the transistor acts as a short, causing the output to fall to zero.

Comparators are useful in interfacing analog signals to digital circuits, which only recognize two states (HIGH voltage and LOW voltage). The comparator simply compares the analog signal to some reference voltage, and “tells” the digital circuit whether the analog signal is above or below the reference, but using the correct digital voltages. The states of 0 and  $+5\text{ V}$  as in the LM311 example above are appropriate for TTL logic, for example. More complex interfaces are certainly possible, and we will return to this later when we discuss analog-to-digital conversion.

Beyond digital interfacing and analog-to-digital conversion, other applications of comparators include oscillators and drivers for alarms or indicators (LEDs, buzzers, beepers) based on an input sensor (e.g., for temperature or water level).

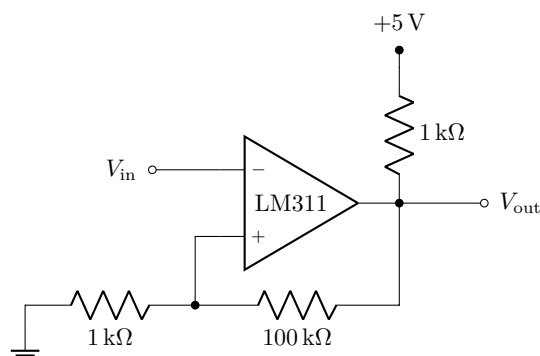
### 7.9.1 Schmitt Trigger

One problem with comparators arises with noisy input signals. Consider the noisy input voltage below going into a comparator with the reference voltage shown. What we want from the comparator is a signal that reflects when the input signal goes above or below the reference. The corresponding output is shown in the lower graph.



But what we see is that due to the noise, the output signal makes many (spurious) transitions whenever the signal crosses a reference, whereas we would expect a smooth input signal to make only one transition at each crossing.

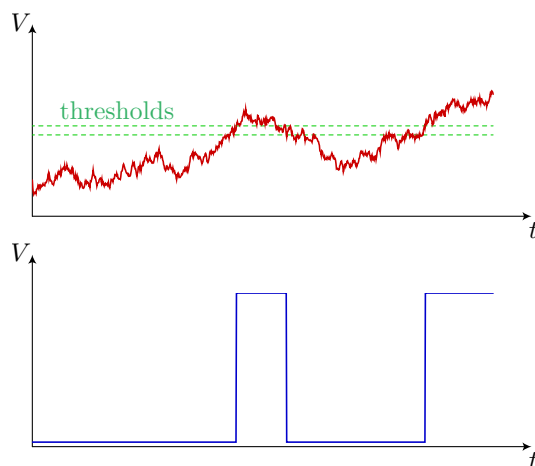
A solution to this is *positive feedback*, which introduces hysteresis. The circuit below, based on the LM311, uses feedback to the noninverting input.



Again, the output swings between 0 V and 5 V, depending on the inputs. Now look at the two cases.

1. If  $V_{in}$  is low, then  $V_{out}$  is high (+5 V), and the trigger point is about 50 mV.
2. If  $V_{in}$  is high, then  $V_{out}$  is low (0 V), and the trigger point is 0 mV.

The trigger point depends on the output, and thus to the input; in other words,  $V_{in}$  “repels” the trigger point, and this gives the circuit immunity to noise at the level of about 50 mV or less. The schematic operation of the Schmitt trigger, from introducing the two effective trigger points, is shown below on the same signal. The hysteresis suppresses the spurious transitions. (Note that the output is inverted compared to the discussion of the LM311 circuit, so it compares more closely to the comparator output in the previous graphs.)



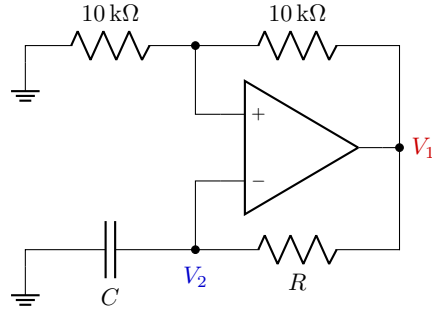
Of course, other nominal trigger levels besides 0 V are possible, by replacing the 1-k $\Omega$  resistor with a voltage divider. The Thévenin resistance of the divider acts in place of the 1-k $\Omega$  resistor.

## 7.10 Positive Feedback and Oscillator Circuits

Besides the Schmitt trigger, positive feedback is useful in op-amp oscillators. We will study two examples of positive-feedback oscillators here: a relaxation oscillator and a phase-shift oscillator.

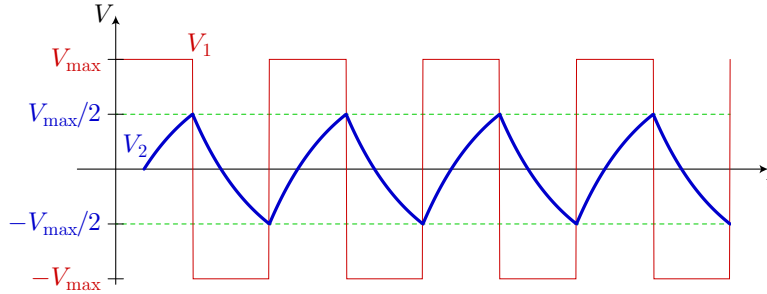
### 7.10.1 Relaxation Oscillator

One good example of a positive-feedback oscillator is the **relaxation oscillator**, shown below.



Here, the amplifier is standard op-amp, acting as a comparator in open-loop mode. We will assume the output rails are  $+V_{\max}$  and  $-V_{\max}$ . There is a 50% voltage divider feeding the noninverting input, similar to the Schmitt trigger above. This sets the trigger points of the comparator to  $+V_{\max}/2$  and  $-V_{\max}/2$ .

Now consider the output of the oscillator at the two points  $V_1$  and  $V_2$ , shown below.



If the output  $V_1$  is positive, the RC circuit charges  $V_2$  until the inverting-input voltage exceeds the  $V_{\max+}/2$  trigger point, at which point  $V_1$  goes negative, and the charging proceeds in the opposite direction until  $V_2$  reaches  $V_{\max-}/2$ , and the cycle repeats.

To treat this more quantitatively, the interval between the switching times is the time from RC decay of  $V_2$  from  $+V_{\max}/2$  to  $-V_{\max}/2$ . The process is (RC) exponential decay starting from  $+V_{\max}/2$  to  $-V_{\max}$ , so we are waiting for the decay to  $1/3$  of the initial voltage, thinking of  $-V_{\max}$  as “ground.” That is, if  $\Delta t$  is the time interval, then

$$e^{-\Delta t/RC} = \frac{1}{3}, \quad (7.82)$$

so

$$\Delta t = RC \log 3 \approx 1.1 RC. \quad (7.83)$$

The period  $T$  is  $2\Delta t$ , so we have

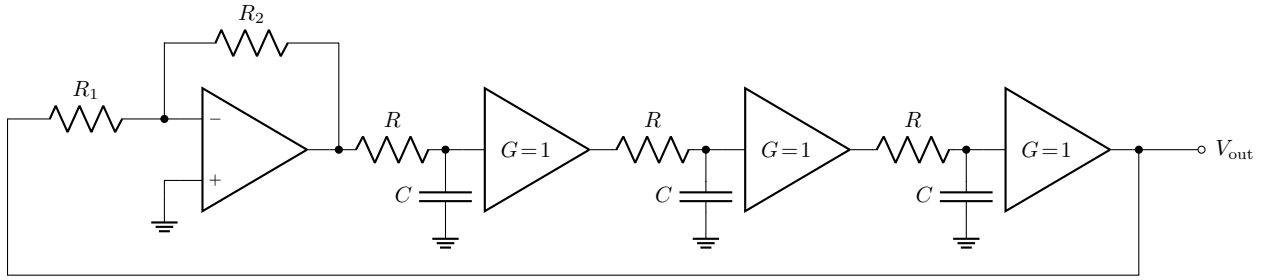
$$T = RC \log 9 \approx 2.2 RC \quad (7.84)$$

(relaxation-oscillator period)

for the period of the relaxation oscillator. The output can be either a quasi-triangle wave or a square wave, depending on which point serves as the output.

### 7.10.2 Buffered Phase-Shift Oscillator

Another example of an oscillator is shown below. It produces a sine wave at a frequency determined by the RC low-pass filters. The buffers are op-amps connected as unity-gain followers.



Note that the first op-amp is connected as an inverting amplifier, and the output  $V_{\text{out}}$  feeds back into the inverting amplifier. There are 3 RC filters in the feedback loop. The oscillation condition is that the phase shift of each RC filter is  $60^\circ$ , so the total RC phase shift is  $180^\circ$ . In combination with the action of the inverting amplifier, this is a total phase of  $0^\circ$ , which means that we have positive feedback.

The correct phase shift only happens at one frequency, which we can find by setting the low-pass-filter phase [Eq. (2.61)]

$$\phi = -\tan^{-1}(\omega RC) \quad (7.85)$$

to  $\phi = 60^\circ$ . The solution is the angular frequency

$$\omega = \frac{\tan 60^\circ}{RC} = \frac{\sqrt{3}}{RC} \approx \frac{1.732}{RC}. \quad (7.86)$$

This corresponds to a frequency  $f = \omega/2\pi$ , or

$$f = \frac{\sqrt{3}}{2\pi RC} \approx \frac{0.276}{RC}. \quad (7.87)$$

(oscillation frequency, phase-shift oscillator)

For example, if  $R = 10 \text{ k}\Omega$  and  $C = 0.01 \text{ }\mu\text{F}$ , then  $f = 2.76 \text{ kHz}$ .

Recall that the low-pass amplitude transfer function is [Eq. (2.45)]

$$T(\omega) = \frac{1}{\sqrt{1 + (\omega RC)^2}}. \quad (7.88)$$

At the oscillation frequency,  $\omega RC = \sqrt{3}$ ,  $T(\omega) = 1/2$ , so the transfer of 3 RC sections is  $1/8$ . Thus, to ensure oscillation, we should set  $R_2/R_1 = 8$  or a bit higher.

## 7.11 Amplifier Noise Analysis

Returning to op-amp circuits with *negative* feedback, let's consider amplifier circuits again. In certain situations, such as where a low-level signal must be amplified or some system (such as a laser) must have its state be precisely controlled, it is critical to have low levels of electronic noise. In the design of low-noise analog circuits, it is necessary to understand and to be able to characterize sources of noise. Previously, we discussed some noise-mitigation techniques related to reducing pickup from external sources (Sections 7.5.2.2 and 7.6.2.2); these techniques amount to better isolation of the critical circuit from coupling to other circuits. However, here we will be interested in *intrinsic* sources of electronic noise—those that can't be avoided completely no matter how well they are isolated.

To begin, one critical aspect of intrinsic noise is the frequency content of the noise. This is characterized through a frequency spectrum. Broadly speaking, and without getting too far into the details of where this comes from, we can associate a **voltage-noise spectral density**  $S_V(f)$ , which encodes the contribution from a particular frequency to a total noise spectrum. (Note that we are using the “regular” frequency  $f$ , not the angular frequency  $\omega$ , because this avoids an extra factor of  $2\pi$  that would have to be tracked.) Then we can write the noise voltage  $v_n$  associated with the spectrum as

$$v_n^2 = \int_0^\infty |S_V(f)|^2 df. \quad (7.89)$$



That is, the voltage noise should be thought of as a (hopefully small), fluctuating voltage signal;  $v_n$  characterizes its size in the root-mean-square sense as

$$v_n = \sqrt{\overline{v^2(t)}}, \quad (7.90)$$

where  $v(t)$  is the explicit, time-dependent noise signal, and the bar  $(\overline{\phantom{x}})$  denotes a time average. From Eq. (7.89), the dimensions of  $S_V(\omega)$  must be  $V/\sqrt{\text{Hz}}$ , an odd-looking unit that commonly comes up in noise analysis (and op-amp data sheets). Correspondingly, one can write down an analogous relation for the rms current noise  $i_n$  in terms of a noise spectrum  $S_I(f)$  as

$$i_n^2 = \int_0^\infty |S_I(f)|^2 df, \quad (7.91)$$

where  $S_I(f)$  is measured in  $A/\sqrt{\text{Hz}}$ . If the voltage refers to a difference across a pure resistance  $R$ , then the two noise quantities are related by  $v_n = i_n R$ . More generally, however, if there is a frequency-dependent impedance  $Z(f)$ , we should instead write

$$S_V(f) = S_I(f)Z(f), \quad (7.92)$$

because it is the noise spectra that are related, with the impedance affecting each frequency differently.

It's worth elaborating on a couple of details before continuing. First, from Eq. (7.92), evidently the spectra may be complex functions. There are some subtleties in handling this properly, but for our purposes this is irrelevant; any complex phase will get discarded in computing the magnitude of the noise signals in Eqs. (7.89) or (7.91). The second, related detail is in the squaring of the spectrum *before* we sum over the contributions of different frequencies. There are a couple of useful ways to think about this. First, since these are noise fluctuations associated with different frequencies, the signals at different frequencies should not have any special phase relation, and so it is the *power* or *energy* of the different components that should add, not the voltage or current. [More mundanely, adding different frequency components and *then* squaring to find the energy only gives extra “cross” terms of the form  $\sin(2\pi f_1 t) \sin(2\pi f_2 t)$ , which average to zero after enough time if  $f_1 \neq f_2$ .] However, it is also worth noting that this adding of squared signals mirrors what happens when you combine independent, random quantities: they add **in quadrature**. To see what we mean by this, consider two random variables  $X_1$  and  $X_2$ . For simplicity, suppose they have zero mean (but a similar argument carries through for nonzero mean):  $\overline{X_1} = \overline{X_2} = 0$ . Then the variance of the sum  $X_1 + X_2$  is

$$\overline{(X_1 + X_2)^2} = \overline{X_1^2 + X_2^2 + 2X_1X_2} = \overline{X_1^2} + \overline{X_2^2} + 2\overline{X_1X_2} = \overline{X_1^2} + \overline{X_2^2}, \quad (7.93)$$

where we used  $\overline{X_1X_2} = \overline{X_1}\overline{X_2}$  for *independent* random variables. This means that the *variances* of random variables add. (This is the same idea behind error-propagation, adding independent contributions to the total statistical error.) In the case of summing over a spectrum in Eq. (7.89), you can think of each frequency component as something like a random variable, because it is a quantity varying sinusoidally, but with some random phase. The total fluctuation at any time is the sum over all of these “random” contributions, but these contribution should be added in quadrature to characterize their magnitude.

### 7.11.1 Sources of Noise

There are many intrinsic sources of noise, but we will discuss only the most important ones to worry about (at least for *intrinsic* noise).

#### 7.11.1.1 Johnson–Nyquist Noise

One of the most important sources of noise in circuits is associated with thermal motion of electrons in resistive materials. The associated noise goes by the name of **Johnson noise** or **Johnson–Nyquist noise**.<sup>3</sup>

<sup>3</sup>Johnson noise is so named after the first person to measure it: J. B. Johnson, “Thermal Agitation of Electricity in Conductors,” *Physical Review Letters* **32**, 97 (1928) (doi: 10.1103/PhysRev.32.97). It is also called Johnson–Nyquist noise, named additionally after the first to describe it theoretically: H. Nyquist, “Thermal Agitation of Electric Charge in Conductors,” *Physical Review Letters* **32**, 110 (1928) (doi: 10.1103/PhysRev.32.110).

This noise corresponds to a voltage noise density<sup>4</sup>

$$S_V(f) = \sqrt{4R(f)k_B T}. \quad (7.94)$$

(Johnson–Nyquist voltage noise density)

The only frequency dependence here is in the resistance, and Johnson noise is commonly treated within the idealization of a frequency-*independent* resistance, in which case the spectral density is a constant in frequency. This is clearly a problem if you naïvely go plugging this density into Eq. (7.89), because the noise voltage will diverge (and an infinity of noise isn’t good for *anybody*). A flat noise density corresponds to **white noise** (which is simply defined as noise having constant power with frequency), and is an idealization that requires some careful handling. The idea is that the noise fluctuations should only be physically important over some bandwidth—measurement electronics always have some bandwidth limit, not to mention that at some frequency the assumption of a constant resistance must break down (the resistance must in fact tend to zero for consistency with causality). As a simple way to handle this, assume that we only detect the noise up to a certain “brick-wall bandwidth”  $B$ . This means that the detector has unit-gain response from frequencies from dc to  $B$ , and zero response above this frequency. Then we should modify Eq. (7.89) to read

$$v_n^2 = \int_0^B |S_V(f)|^2 df, \quad (7.95)$$

leading to an rms noise voltage

$$v_n = \sqrt{4Rk_B TB} \approx 0.128 \frac{\text{nV}}{\sqrt{\Omega \cdot \text{Hz}}} \cdot \sqrt{RB} \quad (7.96)$$

(rms Johnson–Nyquist noise)

at 25°C, or equivalently, a noise current

$$i_n = \sqrt{\frac{4k_B TB}{R}} \approx 0.128 \frac{\text{nA}}{\sqrt{\text{Hz}/\Omega}} \cdot \sqrt{\frac{B}{R}}, \quad (7.97)$$

(rms Johnson–Nyquist noise)

That is, the amount of noise we detect depends on the detection bandwidth for white noise, and (according to this nonphysical idealization) the noise increases without bound as the bandwidth increases. Note also that, as a side benefit, when we add different noise voltages or noise currents, and we add them in quadrature as prescribed by Eq. (7.93), the expressions are naturally compatible with the rules for adding resistances in series (when  $v_n$ ’s add) and for adding parallel resistances (when  $i_n$ ’s add).

To put in some more specific sample numbers, at  $T = 293$  K, a 10 kΩ resistor measured over a  $B = 10$  kHz bandwidth has an intrinsic, thermal rms voltage noise of  $1.2 \mu\text{V}$ . At the same temperature, a 1-MΩ resistor measured over a  $B = 1$ -MHz bandwidth has an rms voltage noise of 0.12 mV, which is starting to become significant on the scale of laboratory voltages.

The brick-wall-filter response is not terribly physical, but more physical filter responses are straightforward to incorporate by including the detector response in the noise integral (7.89). For example, for a simple low-pass filter response,  $S_V(f)$  should be multiplied by a factor  $1/(1 - if/f_{3\text{dB}})$ , where  $f_{3\text{dB}}$  is the 3-dB frequency of the filter. Thus, Eq. (7.89) should be modified to become

$$v_n^2 = \int_0^\infty \frac{|S_V(f)|^2}{|1 - if/f_{3\text{dB}}|^2} df = 4Rk_B T \int_0^\infty \frac{df}{1 + (f/f_{3\text{dB}})^2} = 2\pi Rk_B T f_{3\text{dB}}. \quad (7.98)$$

Notice that this is equivalent to Eq. (7.96), provided we identify

$$B = \frac{\pi f_{3\text{dB}}}{2}, \quad (7.99)$$

under the assumption of white noise.

---

<sup>4</sup>The derivation is quite mathematical, but for one version see Daniel Adam Steck, *Quantum and Atom Optics*, available online at <http://steck.us/teaching>, Section 14.3.8.2.

### 7.11.1.2 “ $kT/C$ ” Noise

Another noise that is basically just Johnson noise is associated with capacitance. Consider an  $R$ – $C$  low-pass filter, and recall that the filter’s corner frequency is, from Eq. (2.47),

$$f_{3\text{ dB}} = \frac{1}{2\pi RC}. \quad (7.100)$$

Putting this into Eq. (7.98) gives

$$v_n = \sqrt{\frac{k_B T}{C}} \approx \frac{64.2 \text{ nV} \cdot \sqrt{\mu\text{F}}}{\sqrt{C}}. \quad (7.101)$$

( $k_B T/C$  noise)

at 25°C. Surprisingly, the dependence on the resistance cancels completely out, so the result is effectively associated with only the capacitance. But, of course, some kind of resistance must be present to generate the noise—reactances don’t generate Johnson noise, only resistive impedances. Intuitively, a larger resistance creates more noise, but the noise is filtered to a narrower bandwidth by the smaller cutoff frequency. Importantly, larger capacitances are associated with lower noise,

### 7.11.1.3 Shot Noise

Another type of noise is associated with flowing current: you might think of electrical current as a smooth flow (like the flow of water), but the current itself is the motion of discrete charge carriers. So if you sit at a certain point on a wire and watch current flowing past (to the extent this is possible), with the “smooth current,” the charge that passes by you rises linearly with time:  $Q(t) = It$ . However, counting the discrete passage of electrons instead,  $Q(t)$  rises in little jumps, and it’s only the *average* slope that works out to  $I$ . This intrinsic noise in the current is called **shot noise**, and to a good approximation shot noise behaves as white noise (provided many carriers are passing by in any given observation time). The current shot noise is<sup>5</sup>

$$i_n = \sqrt{2eI|B|}, \quad (7.102)$$

(current shot-noise)

where  $e$  is the electron charge,  $I$  is the (mean) current, and  $B$  is again the brick-wall bandwidth. This expression assumes that the electrons in the current arrive *independently*; if electron motion is correlated, this can affect the magnitude of the shot noise. Shot noise is an important effect in, for example, vacuum tubes and in semiconductor-junction devices, or in any situation that is sensitive to current noise.

Shot noise applies in other situations where there is a mean flow of discrete particles. For example, laser light comes in discrete photons of energy  $\hbar\omega$ , where  $\omega$  is the optical frequency. Adapting the formula (7.102) gives

$$\delta P = \sqrt{2\hbar\omega PB} \quad (7.103)$$

(optical shot-noise)

as the rms power fluctuation, given a mean optical power  $P$ .

### 7.11.1.4 $1/f$ Noise

Another important class of noise is  **$1/f$  noise**, which has a spectral density of the form

$$S_V(f) = \sqrt{\frac{k}{f}}, \quad (7.104)$$

( $1/f$  noise density)

for some constant  $k$ . Note that the noise is so named because the noise *power* density (i.e.,  $S_V^2$ ). This noise, also called **pink noise** or **flicker noise**, was first noticed as a component of the shot noise in a

<sup>5</sup>For more details, see Daniel Adam Steck, *Quantum and Atom Optics*, available online at <http://steck.us/teaching>, Section 17.5.3.1.

vacuum tube.<sup>6</sup> However, it is quite ubiquitous, appearing in condensed-matter, biological financial, and other systems. It is self-similar in the sense of having the same power within different frequency octaves (or decades). In electronics, it is generally due to fluctuations in the properties of materials (semiconductors as well as metals). In the original observation, the  $1/f$  component of the noise was attributed to a fluctuation in the emissivity of the oxide cathode material.

The  $1/f$  spectrum (7.104) is even more problematic in some sense than white noise, because in the noise expression (7.89), the frequency integral is now problematic at both ends of the integral. Thus,  $1/f$  noise must be observed only within a finite frequency band from  $f_{\min}$  to  $f_{\max}$ . Thus, the noise voltage becomes

$$v_n^2 = k \int_{f_{\min}}^{f_{\max}} \frac{df}{f} = k \log(f_{\max}/f_{\min}). \quad (7.105)$$

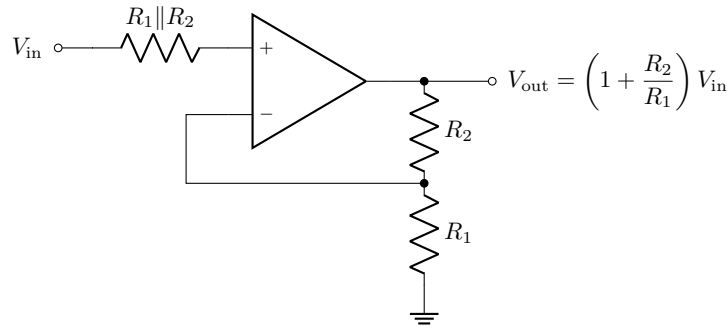
The divergence as  $f_{\max} \rightarrow \infty$  is more “gentle” than for white noise (which diverges linearly in the analogous expression); in fact  $v_n$  is relatively insensitive to the precise (finite) value of  $f_{\max}$ . However, since white noise and  $1/f$  noise typically coexist, the white noise will be the more important consideration above some frequency. (Typically, for intrinsic op-amp noise, which contains both components, the noise spectrum crosses over from  $1/f$  to white somewhere between 10 and 100 Hz.)

For low frequencies  $1/f$  noise will dominate the white noise, however. The lower cutoff  $f_{\min}$  is related to the finite **observation time** of the noise signal (if the signal is observed over some long time  $\tau$ , it is only possible to resolve frequencies down to the order of  $1/\tau$ ). This means it is possible to “cure” the  $1/f$  divergence at low frequencies, but that the noise level depends on the observation time. As in the dc divergence of the response of the op-amp integrator (Section 7.4.2), the divergence as  $f \rightarrow 0$  means intuitively that the response of a circuit or device drifts with time, not staying near any particular state. For example, a current source, given a particular forward voltage, will have a current that is steady on short time scales, but drifts on long time scales.

Thus,  $1/f$  noise is important, but tends to be most important for longer-term observations and drifts. Here, we will mostly be concerned with intrinsic noise superposed on a dynamic signal in either wideband or precision circuits—that is, we will stick mainly to analyzing the effects of white noise.

### 7.11.2 Estimation of Amplifier Noise

As an example of how to estimate the noise in an op-amp circuit, let’s consider the bias-compensated noninverting amplifier from Section 7.6.1.3.



Note that the noise will be superposed with whatever the normal amplified signal will be; so this analysis will apply equally well to the bias-compensated *inverting* amplifier from Section 7.6.1.3—the only difference is the ground and input voltages are swapped in the circuit.

From Eq. (7.94), the resistance at the noninverting input presents a voltage-noise density at the input of

$$S_{V_+} = \sqrt{4(R_1 \parallel R_2)k_B T}. \quad (7.106)$$

<sup>6</sup>J. B. Johnson, “The Schottky Effect in Low Frequency Circuits,” *Physical Review* **26**, 71 (1925) (doi: 10.1103/PhysRev.26.71).

The Thévenin equivalent circuit for the voltage divider presents the same resistance to the inverting input, so we get the same noise density there:

$$S_{V-} = \sqrt{4(R_1 \parallel R_2)k_B T}. \quad (7.107)$$

The op-amp itself will have intrinsic noise; this will generally be specified as a voltage and current noise, referenced to the inputs (that is, noise voltage and current as if they are fed to the input). The op-amp voltage noise,  $S_{V_{\text{opamp}}}$ , will add to the Johnson noises. In this circuit, the op-amp current noise  $S_{I_{\text{in}}}$  will contribute to the noise by being converted to voltage noise by the input impedances: each input will see an extra  $S_{I_{\text{opamp}}}(R_1 \parallel R_2)$  of voltage noise.

All of these components add together to give a total *input*-referenced voltage noise density of

$$\begin{aligned} S_{V_{\text{in}}} &= \sqrt{S_{V+}^2 + S_{V-}^2 + S_{V_{\text{opamp}}}^2 + 2S_{I_{\text{opamp}}}^2(R_1 \parallel R_2)^2} \\ &= \sqrt{8(R_1 \parallel R_2)k_B T + S_{V_{\text{opamp}}}^2 + 2S_{I_{\text{opamp}}}^2(R_1 \parallel R_2)^2}. \end{aligned} \quad (7.108)$$

Note the factor of 2 in the last term, corresponding to adding uncorrelated noise sources. There is some subtlety here, because the input current noise of an op-amp will generally have both incoherent and coherent parts. For example, for a bipolar-input op-amp, if we refer back to the BJT differential amplifier in Section 4.10.4, the current source in the “long tail” contributes current noise that is split between the two input transistors; thus this noise is correlated. In more sophisticated op-amps like the LT1028, there is an input-bias cancelling circuit that injects current into the bases of the input transistors, trying to minimize the bias current from the op-amp inputs; this system also results in correlated current noise into the inputs. On the other hand, both input transistors will have their own shot noise, which is an example of *uncorrelated* noise. Recall that the rule above for adding noise sources in quadrature applies to *uncorrelated* noise sources; correlated noises simply add. The correlated current noises would lead to a factor of 4 instead of 2 in the last term under the radical of Eq. (7.108). In this case, we assume the current noise to be uncorrelated because the input impedances are balanced, so any correlated input current noise gets translated into common-mode voltage noise, which is rejected by the op-amp (the data sheet confirms that the  $1 \text{ pA}/\sqrt{\text{Hz}}$  figure corresponds to balanced input impedance). Unfortunately it is sometimes hard to tell whether or not the current noise is correlated, but note that the LT1028 data sheet gives input current noise plots for both matched and unmatched source resistances (when the input resistances are far unmatched, it is better to treat the noises as correlated). But in any case, this can be another advantage to having an input-impedance-matching resistor as in Sections 7.6.1.2 and 7.6.1.3, at least if the reduction in the effect of input current noise is more significant than the added Johnson noise.

At the *output* of the op-amp, the noise density will simply be the input-referenced noise density multiplied by the circuit gain:

$$S_{V_{\text{out}}} = GS_{V_{\text{in}}} = \left(1 + \frac{R_2}{R_1}\right) \sqrt{S_{V+}^2 + S_{V-}^2 + S_{V_{\text{opamp}}}^2 + 2S_{I_{\text{opamp}}}^2(R_1 \parallel R_2)^2}. \quad (7.109)$$

To continue at this point, we will need to be more specific with example numbers. Suppose we take  $R_1 = 100 \Omega$ ,  $R_2 = 100 \text{ k}\Omega$ , so that  $R_1 \parallel R_2 \approx 100 \Omega$ , and  $G = 1001$ . Suppose also that we choose the LT1028 low-noise, wide-band op-amp.<sup>7</sup> For this amplifier, at  $25^\circ\text{C}$  and  $1 \text{ kHz}$  (and for the better LT1028AM/AC variant),  $S_{V_{\text{opamp}}} = 0.85 \text{ nV}/\sqrt{\text{Hz}}$ , and  $S_{I_{\text{opamp}}} = 1.0 \text{ pA}/\sqrt{\text{Hz}}$ . Then the magnitudes of the various noise components is

$$S_{V+} = S_{V-} = 1.28 \frac{\text{nV}}{\sqrt{\text{Hz}}}, \quad S_{V_{\text{opamp}}} = 0.85 \frac{\text{nV}}{\sqrt{\text{Hz}}}, \quad S_{I_{\text{opamp}}}(R_1 \parallel R_2) = 0.1 \frac{\text{nV}}{\sqrt{\text{Hz}}}. \quad (7.110)$$

In this example, the resistor thermal noise is the most important, and the op-amp current noise is almost negligible. Computing the final noise at the output gives

$$S_{V_{\text{out}}} = 2.0 \frac{\mu\text{V}}{\sqrt{\text{Hz}}}. \quad (7.111)$$

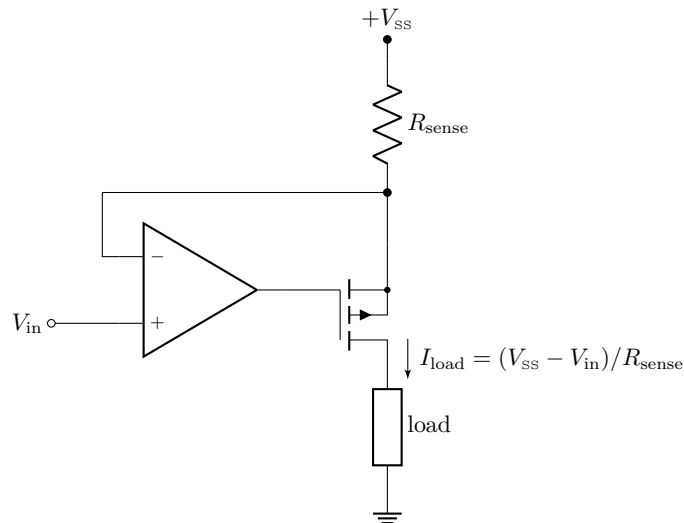
<sup>7</sup><http://cds.linear.com/docs/en/datasheet/1028fd.pdf>. Note that the numbers in the example here match those of the noise calculations in the noise example there.

To compute the size of the voltage fluctuations, it is still necessary to incorporate the operating bandwidth  $B$  of the circuit. Thus, assuming a  $B = 1\text{-MHz}$  (brick-wall) bandwidth, the rms voltage fluctuation would be  $v_n = 2.0\text{ mV}$ .

One important observation here is that because the source resistances were fairly small, the input noise current didn't contribute much to the final noise. More generally, if the source impedances are small, it is most important to employ an op-amp with small input *voltage* noise, while the current noise is less important. Conversely, it's more important to find an op-amp with small input *current* noise if the impedances presented to the inputs are large.

### 7.11.2.1 Example: Precision Current Source

As another example, the circuit below is a precision current source intended for diode lasers—in frequency stabilized diode lasers the noise in the supply current can be one of the most important factors in determining the frequency stability. We will take it for granted here that the current source works as advertised (see Problem 7.30). Also note that this is only the critical current-regulation stage in a larger, low-noise circuit.<sup>8</sup>



Here  $R_{\text{sense}} = 50\ \Omega$ , and the op-amp in the original circuit is an LT1028; to lowest order we expect the noise in the regulated current to be set mainly by these components. Again, the LT1028 specifies a voltage noise density of  $0.85\text{ nV}/\sqrt{\text{Hz}}$  for the better LT1028AM/AC variant at  $1\text{ kHz}$ . This is input-referenced, and translates to a current noise density via the  $50\ \Omega$  sense resistor to  $17\text{ pA}/\sqrt{\text{Hz}}$ . The sense resistor itself contributes Johnson noise with rms voltage noise density of  $S_{V_{\text{sense}}} = (0.91\text{ nV}/\sqrt{\text{Hz}})$ , at  $25^\circ\text{C}$ , corresponding to a current noise density of  $18\text{ pA}/\sqrt{\text{Hz}}$ . Adding the two current noises in quadrature gives  $25\text{ pA}/\sqrt{\text{Hz}}$ . Note that the  $1\text{ pA}/\sqrt{\text{Hz}}$  input noise current density should also be added (in quadrature) to this result, because current flowing through the sense resistor also can flow into the inverting input of the amplifier. However, this is a negligible correction.

The setup here is roughly optimal for the selected amplifier, since the amplifier and Johnson noises are about the same. A larger sense resistor decreases the effect of the amplifier noise, but increases the Johnson noise. A smaller sense resistor has the opposite effect. For example, a  $30\ \Omega$  sense resistor increases the effective amplifier noise to  $28\text{ pA}/\sqrt{\text{Hz}}$ , while decreasing the Johnson noise to  $0.70\text{ nV}/\sqrt{\text{Hz}}$ , which is equivalent to  $23\text{ pA}/\sqrt{\text{Hz}}$ . The total is then  $36\text{ pA}/\sqrt{\text{Hz}}$ , which is worse off, at least using the  $1\text{-kHz}$  noise figure. A smaller resistor may better optimize for larger noise figures at lower frequencies (for example,  $1.0\text{ nV}/\sqrt{\text{Hz}}$  for the LT1028 at  $10\text{ Hz}$ ). However, for laser stability the critical factor is the noise density integrated over a wide band (of hundreds of kHz or several MHz).

<sup>8</sup>This is from the classic current-controller circuit in K. G. Libbrecht and J. L. Hall, "A low-noise high-speed diode laser current controller," *Reviews of Scientific Instruments* **64**, 2133 (1993).

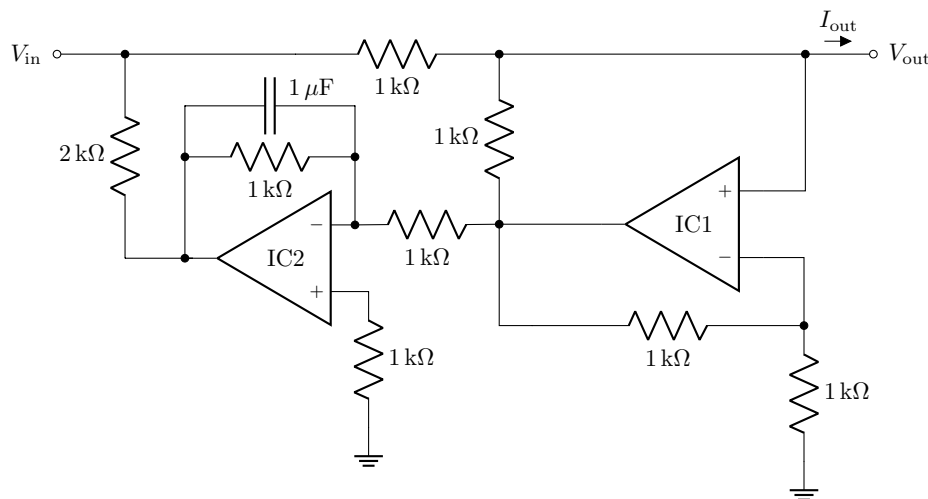
One curiosity of the original circuit is a pair of diodes in the power-supply lines of the op-amp. The diodes on the power-supply inputs don't directly help with noise suppression, but rather are there for protection of the op-amp inputs. For ultralow-noise amplifiers like the LT1028 or AD797, the inputs don't have current-limiting input resistors, but have back-to-back protection diodes. The diodes help to protect the inputs in situations where the amplifier is unpowered (or unusually powered, as in turn-on or turn-off transients).<sup>9</sup>

A final thing to check here is that the control voltage  $V_{in}$  is subject to noise, as it is derived from a voltage divider referenced to  $+V_{SS}$ . The output of the voltage divider is bypassed by a  $22\text{-}\mu\text{F}$  capacitor, so a good estimate of the noise comes from the  $k_B T/C$  noise in Eq. (7.98), which amounts to an rms noise of  $v_n = 14\text{ nV}$ . To compare this to the op-amp input voltage noise, we must invoke a bandwidth. For a rough estimate, suppose we take  $B = 1\text{ MHz}$ , in which case the input voltage noise of the op-amp is  $0.85\text{ }\mu\text{V}$  (i.e., after multiplying by  $\sqrt{B}$ ). On this scale, the control-voltage noise is negligible (and of course can always be made smaller by choosing a larger bypass capacitor).

If you've been paying attention, you might remember that the MOSFET (here specified as a VP0106<sup>10</sup>) acts something like a variable resistor—but to good approximation we are completely ignoring it in this analysis. (*Why?*) (Hint: this is only true assuming the op-amp is doing its job.) The original circuit also has a couple of resistors at the output of the capacitor (a  $10\text{-k}\Omega$  resistor connecting the op-amp output to the MOSFET gate, and a  $100\text{-}\Omega$  resistor in series with a capacitor to “load” the gate). But we can ignore those too. (*Why?*)

The noise estimate here is good as a first cut, but life is more complicated because of the gate capacitance interacting with the finite gain of the op-amp, as modeled in Problem 7.30; these factors combine to produce resonant behavior, which can boost the current noise in the related frequency band.<sup>11</sup>

**Input-stage noise.** The same current source has a modulation-input stage that allows relatively fast modulation of the current via an external control voltage (to sweep the laser current or to modulate it for use in lock-in detection). This input stage is shown below, and the output current  $I_{out}$  is simply added to the current of the main regulating circuit above. This circuit regulates the current to  $I_{out} = V_{in}/(1\text{ k}\Omega)$ , independent of  $V_{out}$  (which depends on the load); the circuit also functions properly when no input is connected (setting  $I_{out} = 0$ ). For the analysis of this circuit, see Problem 7.25.



To analyze the noise here, first consider the noise contributed by the IC1 portion of the circuit. IC1 sees an input resistance of  $500\text{ }\Omega$  at the inverting input, corresponding to  $2.9\text{ nV}/\sqrt{\text{Hz}}$  of Johnson noise. The

<sup>9</sup>This is explained in more detail in John Ardizzone, “Protecting Off-Amps” [http://www.analog.com/library/analogdialogue/archives/42-10/off\\_amps.html](http://www.analog.com/library/analogdialogue/archives/42-10/off_amps.html)

<sup>10</sup><http://ww1.microchip.com/downloads/en/DeviceDoc/VP0106%20C082313.pdf>

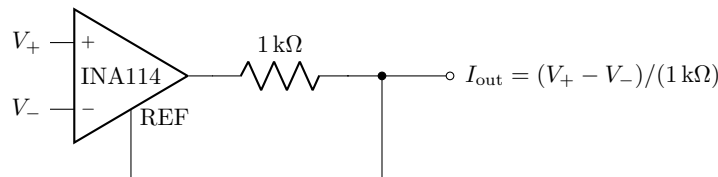
<sup>11</sup>For further analysis with noise measurements, see Chris Seck, Paul J. Martin, Eryn C. Cook, Brian Odom, and Daniel A. Steck, “Noise reduction of a Libbrecht–Hall-style current driver,” *Review of Scientific Instruments* **87**, 064703 (2016) (doi: 10.1063/1.4953330) (arXiv: 1604.00374).

op-amps are OP27's, with  $3.0 \text{ nV}/\sqrt{\text{Hz}}$  at 1 kHz for the better (OP27A/E) grade, and an input noise current density of  $0.4 \text{ pA}/\sqrt{\text{Hz}}$ . Then the total input-referenced noise is  $4.2 \text{ nV}/\sqrt{\text{Hz}}$ , and thus  $8.4 \text{ nV}/\sqrt{\text{Hz}}$  at IC1's output due to the noninverting gain of 2. Across the  $1\text{-k}\Omega$  resistor to  $V_{\text{out}}$ , this converts to  $8.4 \text{ pA}/\sqrt{\text{Hz}}$  of current noise. Adding this to the  $4.0 \text{ pA}/\sqrt{\text{Hz}}$  Johnson noise of the resistor itself gives  $9.3 \text{ pA}/\sqrt{\text{Hz}}$  current noise for the IC1 section. Note that IC1's input noise current should also contribute directly to the  $I_{\text{out}}$  noise, but this contribution is negligible. Notice that adding this current density to the  $25 \text{ pA}/\sqrt{\text{Hz}}$  noise from the main regulating amplifier only increases the noise by about 7%, to  $27 \text{ pA}/\sqrt{\text{Hz}}$ .

That is the noise estimate assuming that the input is held fixed by a noiseless source. With a disconnected input voltage, IC2 also contributes to the current noise. IC2 sees Johnson noises of  $2.9 \text{ nV}/\sqrt{\text{Hz}}$  and  $4.0 \text{ nV}/\sqrt{\text{Hz}}$  at the inverting and noninverting inputs, respectively. Adding the  $3.0 \text{ nV}/\sqrt{\text{Hz}}$  of input voltage noise, and the  $0.4 \text{ pA}/\sqrt{\text{Hz}}$  input current noise (giving  $0.2 \text{ nV}/\sqrt{\text{Hz}}$  and  $0.4 \text{ nV}/\sqrt{\text{Hz}}$  at the inverting and noninverting inputs, respectively, which we can assume worst-case to be correlated), the total input-reference noise is  $5.8 \text{ nV}/\sqrt{\text{Hz}}$ . With a gain of 2 and over the  $3 \text{ k}\Omega$  resistance to  $V_{\text{out}}$ , this is  $4.5 \text{ pA}/\sqrt{\text{Hz}}$ ; added to the  $2.3 \text{ pA}/\sqrt{\text{Hz}}$  of Johnson noise over the same resistance, this is  $5.1 \text{ pA}/\sqrt{\text{Hz}}$  noise from IC2. Together with IC1, this is  $11 \text{ pA}/\sqrt{\text{Hz}}$  of current noise—not too much different from the IC1 part of the noise. So when running this circuit, shorting the input to ground when not in use should reduce the noise in the circuit slightly, but realistically the noise reduction would barely be noticeable.

The point of this analysis is that, even though the resistances here are larger and the op-amps are noisier than in the primary circuit above, the contribution of this input stage is fairly minimal compared to the main circuit because of its small contribution to the total regulated current.

**A cautionary tale.** Instrumentation amplifiers have a reputation as high-quality, low-noise circuits. They readily form a current source, as in the circuit below, and in view of the discussion of instrumentation amplifiers as a differential receiver in Section 7.5.2.2, this would seem to be an ideal basis for a modulation input for current controller.



This circuit uses an INA114 instrumentation amplifier<sup>12</sup> as a current source. Without a gain resistor connected, the gain is unity, so the amplifier output is  $(V_+ - V_-)$  with respect to the reference voltage. Thus, the amplifier maintains  $(V_+ - V_-)$  across the sense resistor, thus giving the advertised current-regulating behavior.

Things look promising, with the data sheet specifying a noise voltage of  $11 \text{ nV}/\sqrt{\text{Hz}}$  at 1 kHz and a gain of 1000. Across the sense resistor this would be  $11 \text{ pA}/\sqrt{\text{Hz}}$  of current noise, which is not bad on the scale of the modulation circuit above. However, a simple sanity check shows this can't be so simple. Recall that an instrumentation amplifier is packed with resistors (see the diagram on p. 202). In the INA114 the output amplifier is a differential op-amp with four  $25 \text{ k}\Omega$  resistors. The Johnson noise of these resistors alone would contribute  $20 \text{ nV}/\sqrt{\text{Hz}}$  to the output voltage, which already exceeds the number quoted by the data sheet.

The complication is that it is better to think of instrumentation amplifiers as having noise sources at both the inputs *and* the output. But data sheets commonly specify these two sources lumped together as a single source referenced at the input. The upshot is that the input noise is added to the noise at the output, but the output noise source is divided by the gain. So for  $G = 1000$ , the noise looks artificially small. Digging deeper into the data sheet, a plot gives the  $G = 1$  input-referenced noise, which is a little over  $60 \text{ nV}/\sqrt{\text{Hz}}$  in the white-noise region of the noise spectrum. Again, that corresponds to  $60 \text{ pA}/\sqrt{\text{Hz}}$  of current noise, which is far worse than the above modulation circuit, and is far worse than the main MOSFET/op-amp current source that this is trying to modulate. Bad instrumentation amp!

The moral of the story here is that instrumentation amps are great, but even better in situations that require amplifying high-impedance signals (which will have their own significant Johnson noise) or in

<sup>12</sup><https://www.ti.com/lit/ds/symlink/ina114.pdf>



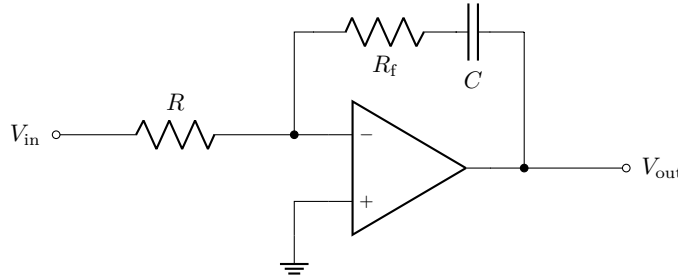
amplifying low level signals with substantial gain (which minimizes the effects of the internal  $25\text{-k}\Omega$  resistors).

## 7.12 Circuit Practice

### 7.12.1 Analog Computers

#### 7.12.1.1 Proportional–Integral Amplifier

As an introduction to the next problem, compute the output voltage in the op-amp circuit below. (It should be proportional to the sum of the integral of the input signal and the input signal itself.) For simplicity, assume  $V_{\text{out}}(0) = V_{\text{in}}(0) = 0$ .



**Solution.** Using the inverting-amplifier result,

$$V_{\text{out}} = -\frac{R_f + X_C}{R} V_{\text{in}} = -\frac{R_f}{R} V_{\text{in}} - \frac{i}{\omega RC} V_{\text{in}}. \quad (7.112)$$

Multiplying through by  $-i\omega$ ,

$$-i\omega V_{\text{out}} = i\omega \frac{R_f}{R} V_{\text{in}} - \frac{1}{RC} V_{\text{in}}, \quad (7.113)$$

and then changing to derivatives,

$$\frac{dV_{\text{out}}}{dt} = -\frac{R_f}{R} \frac{dV_{\text{in}}}{dt} - \frac{1}{RC} V_{\text{in}}. \quad (7.114)$$

Integrating,

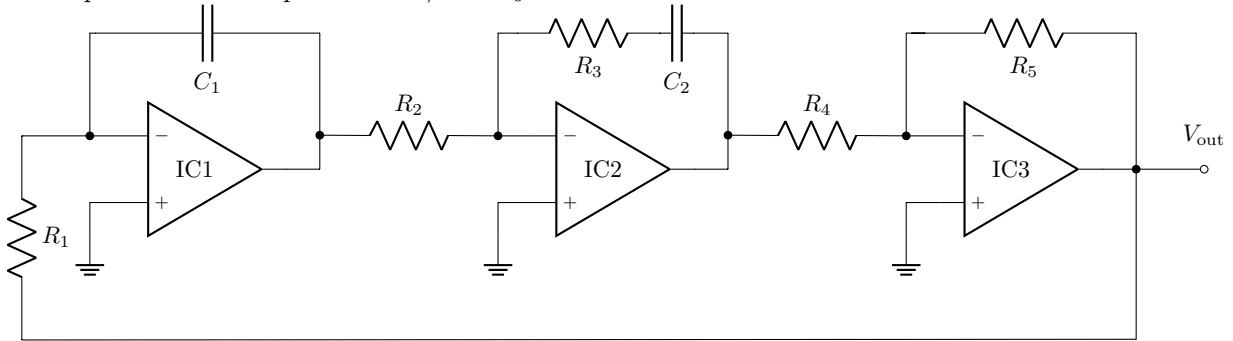
$$V_{\text{out}}(t) = -\frac{R_f}{R} V_{\text{in}}(t) - \frac{1}{RC} \int_0^t V_{\text{in}}(t') dt'. \quad (7.115)$$

#### 7.12.1.2 Damped Harmonic Oscillator

The circuit below is an example of an **analog computer**, in this case a computer that solves a differential equation. In particular, show that this circuit solves the damped-harmonic-oscillator equation,

$$\ddot{x} = -\gamma \dot{x} - \omega_0^2 x. \quad (7.116)$$

Give expressions for the parameters  $\gamma$  and  $\omega_0$  in terms of the  $R$  and  $C$  values.



*Hint:* think of the input to IC1 as  $\ddot{x}(t)$ , and start integrating from there.

Actually, this circuit only solves for  $\ddot{x}(t)$ , while  $x(t)$  is buried in an inaccessible way in IC2. Can you think of a way to modify the circuit, by replacing IC2 with two other op amps, to make  $x(t)$  available?

Note also in particular how  $R_3$  controls the *damping* ( $\gamma$ ) for the circuit.

**Solution.** Suppose the input to the IC1 integrator is  $\ddot{x}$ . Then the output of IC1 (and the input of IC2) is

$$V_{IC1} = -\frac{\dot{x}}{R_1 C_1}. \quad (7.117)$$

Now applying the results of the first problem, the output of IC2 is

$$V_{IC2} = \frac{R_3}{R_1 R_2 C_1} \dot{x} + \frac{x}{R_1 R_2 C_1 C_2}. \quad (7.118)$$

The last op amp just inverts with some gain:

$$V_{IC3} = V_{out} = -\frac{R_3 R_5}{R_1 R_2 R_4 C_1} \dot{x} - \frac{R_5}{R_1 R_2 R_4 C_1 C_2} x. \quad (7.119)$$

Then  $V_{out} = \ddot{x}$ , so we have the equation

$$\ddot{x} = -\frac{R_3 R_5}{R_1 R_2 R_4 C_1} \dot{x} - \frac{R_5}{R_1 R_2 R_4 C_1 C_2} x. \quad (7.120)$$

This is the harmonic-oscillator equation with

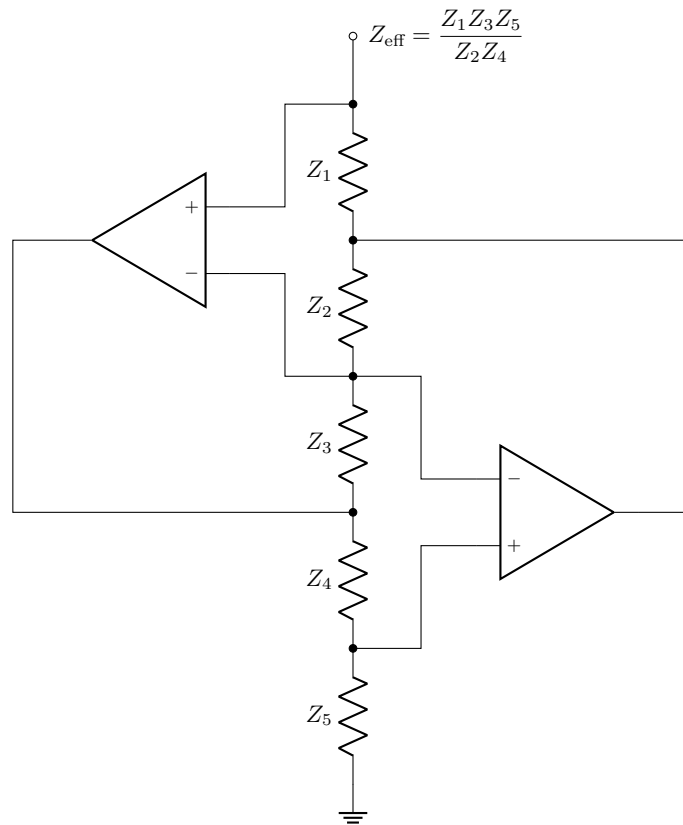
$$\gamma = \frac{R_3 R_5}{R_1 R_2 R_4 C_1}, \quad \omega_0 = \sqrt{\frac{R_5}{R_1 R_2 R_4 C_1 C_2}}. \quad (7.121)$$

Note that we set  $V_{out}$  to  $\ddot{x}$ , so this solves for the second derivative. To make  $x(t)$  available, replace IC2 by two op amps, one an integrator and one an inverting amplifier. Then use IC3 as an inverting summer. The  $x(t)$  signal is then available on the output of the inverting amplifier.

### 7.12.2 Gyrator

The circuit below is an example of a **gyrator**,<sup>13</sup> which presents an effective impedance based on the constituent impedances  $Z_1$ – $Z_5$  (which are not necessarily resistors).

<sup>13</sup>Paul Horowitz and Winfield Hill, *The Art of Electronics*, 2nd ed. (Cambridge, 1989), p. 291 (ISBN: 0521370957).



(a) Show that the effective input impedance is as advertised. Begin by considering a voltage  $V$  and a current  $I$  at the input terminal, and divide them to obtain the impedance.

*Hints:* to get you started, here are a few things to simplify and keep things organized. First, consider currents  $I_1$ – $I_5$ , flowing through “resistors”  $Z_1$ – $Z_5$ , all in the direction of the ground connection.

Second, *assume* that both op amps are operating “normally” with negative feedback (this is *plausible* in this circuit, since the outputs are “closer” in the impedance chain to the negative inputs than to the negative outputs, but strictly speaking this would need to be proven; for simplicity, just assume this is the case).

Third, note that  $V = I_5 Z_5$ . (Why?)

Finally, use what you know about op amps to relate all the different currents together; you don’t need to consider any other currents besides  $I$  and  $I_1$ – $I_5$ , and obviously you want to eliminate all of them but  $I$ .

(b) One of the utilities of this circuit is to realize an effective inductor using only op amps, resistors, and capacitors. This is useful since these components are often better behaved (i.e., closer to ideal) than inductors. Suppose  $Z_4$  is a capacitor in this circuit, with the rest resistors. Show that the result is an effective inductor, and give the effective inductance. What set of (reasonable) components would give you a 1-H inductor? (A pretty big inductor!)

### Solution.

(a) We start with

$$V = I_5 Z_5. \quad (7.122)$$

This is because the voltage drop across the other four resistors must be zero, because they are wrapped between the inputs of the op amps.

Now the voltage between the inputs of the right-hand op amp is zero, so

$$I_3 Z_3 = -I_4 Z_4. \quad (7.123)$$

Then using  $I_4 = I_5$  (no current into the op-amp input), we find

$$V = I_4 Z_5 = -\frac{Z_3 Z_5}{Z_4} I_3. \quad (7.124)$$

Repeating this argument, we have

$$I_1 Z_1 = -I_2 Z_2 \quad (7.125)$$

and  $I_2 = I_3$ , so

$$V = -\frac{Z_3 Z_5}{Z_4} I_2 = \frac{Z_1 Z_3 Z_5}{Z_2 Z_4} I_1. \quad (7.126)$$

Finally,  $I_1 = I$ , so

$$V = \frac{Z_1 Z_3 Z_5}{Z_2 Z_4} I =: I Z_{\text{eff}}, \quad (7.127)$$

which establishes the effective impedance.

(b) With

$$Z_{\text{eff}} = \frac{Z_1 Z_3 Z_5}{Z_2 Z_4} \quad (7.128)$$

and setting  $Z_4 = i/\omega C$  and the other impedances to resistances,

$$Z_{\text{eff}} = -i\omega \frac{R_1 R_3 R_5 C}{R_2}. \quad (7.129)$$

Comparing this to

$$X_L = -i\omega L, \quad (7.130)$$

we have an inductance

$$L_{\text{eff}} = \frac{R_1 R_3 R_5 C}{R_2}. \quad (7.131)$$

With  $C = 0.01 \mu\text{F}$ , we could pick all resistors to be  $10 \text{ k}\Omega$ , which would give  $1 \text{ H}$  of inductance.

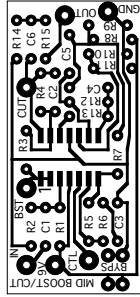
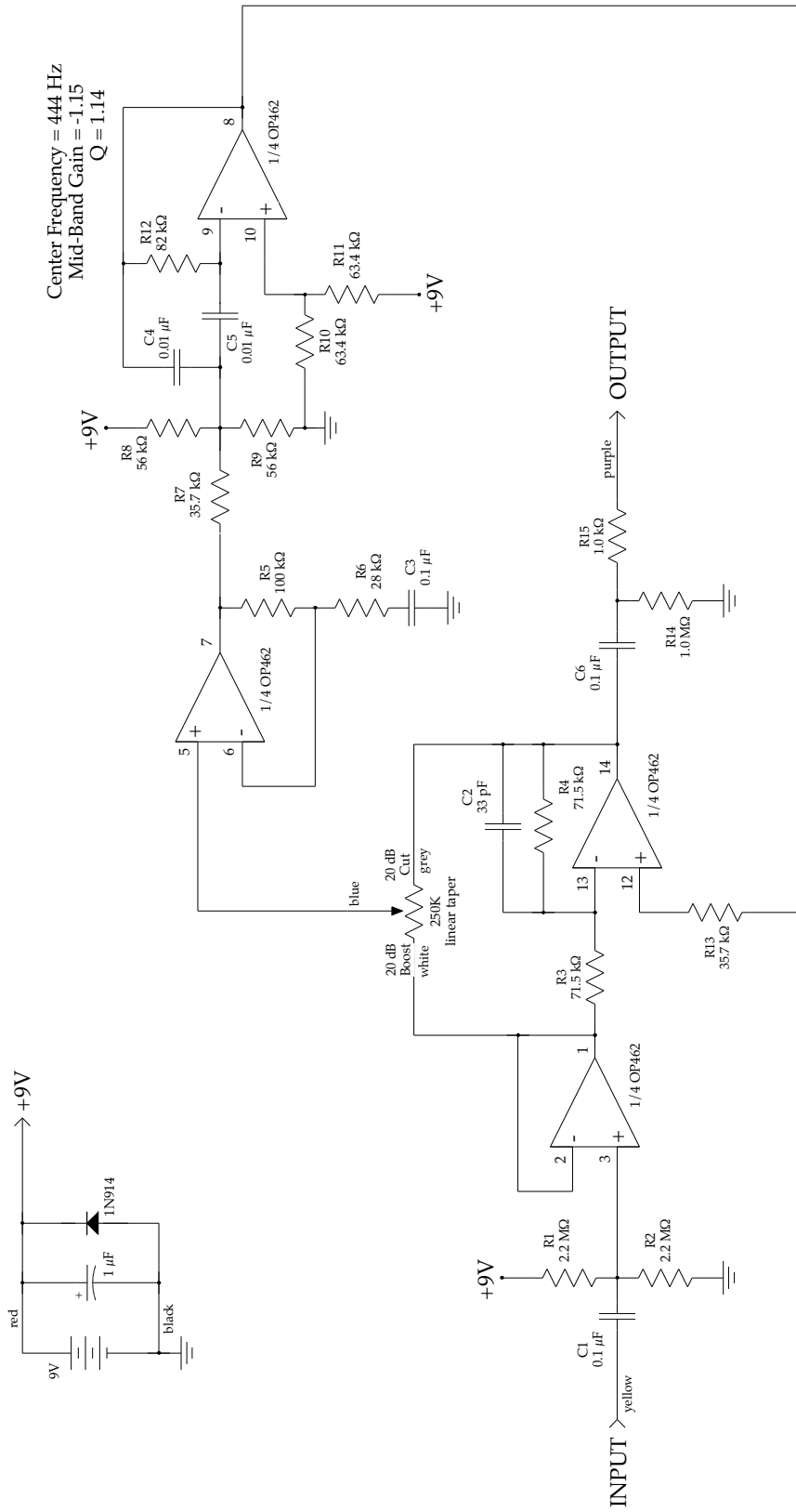
### 7.12.3 Guitar Preamp with Midrange Boost/Cut

The circuit on the next page is a preamplifier for an electric guitar, powered from a single 9-V battery. It is designed for a Fender Stratocaster, and uses one of the “tone knobs” to control a midrange boost or cut (a midrange boost gives a “fat” sound more like a Les Paul guitar, with “humbucking pickups;” a midrange cut gives a clear, “thin” sound, more like the neck pickup on a Fender Telecaster guitar). The PCB design is also shown, printed at actual size. (Compare this circuit to the preamplifier for the Eric Clapton Signature Stratocaster from Section 4.16.4.)

Try to work through the circuit and understand each of the elements, noting the following:

- Since this circuit is powered by a single 9-V battery, but the signal is bipolar, all circuits must be referenced to the “effective ground” of 4.5 V. For example, the input is ac-coupled and biased at this effective ground.
- There are four op-amps, but all packaged in one chip. Hence the “1/4 OP462,” and the IC pin numbers on each op-amp.
- First, convince yourself that the input op-amp (pins 1-3) is an ac-coupled, unity-gain buffer.
- Now, the upper-left op-amp (pins 5-7) functions as a noninverting amplifier. But the capacitor gives a frequency-dependent gain. You should convince yourself that the dc gain is unity, but the ac gain is higher. What is the ac gain? Why do we want unity gain at dc?

- The upper-right-hand op-amp is less obvious, but this is an inverting bandpass filter, as suggested by the center frequency and  $Q$  factor. Note that the noninverting input is biased at virtual ground (normally this input would just be grounded), and that there are two R–C pairs at the input (R7 in parallel with R8 and R9, with C4; and R12 with C5) that together give the bandpass action, because they act something like cascaded low-pass and high-pass filters. (See Problem 7.24.)
- The final op-amp (pins 12-14) combines the filtered signal with the original to give the boost or cut. Capacitor C2 rolls off the gain at high frequencies, where the combination may not be accurate due to different delays at the inputs.
- To control the amount of boost or cut, a potentiometer (with color-coded wires according to the standard Fender convention for the tone knob) interpolates between the buffered input signal and the inverted version produced by the last op-amp.
- The output is ac-coupled and biased to ground.
- Note the diode in the battery/power-supply connection, which protects the op-amp in case the battery is accidentally connected in reverse. A single bypass capacitor stabilizes all op-amps, since they are in a single package. The op-amp operates at low power, so the bypass capacitor value and location are not critical; a larger capacitance is fine since the capacitor does not need to work well at very high frequencies (audio = low frequency in circuits).



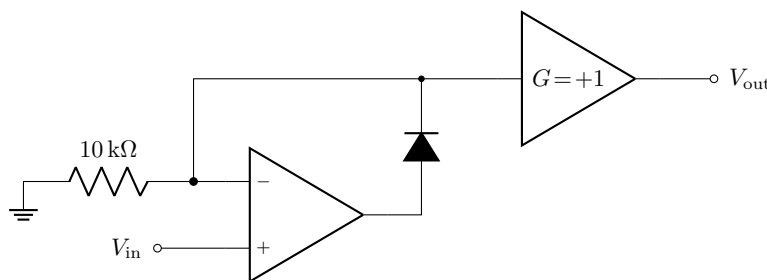
GUITAR PREAMP / MID-BOOST CIRCUIT
D. A. STECK 4 / 8 / 00

### 7.12.4 Active Rectifiers

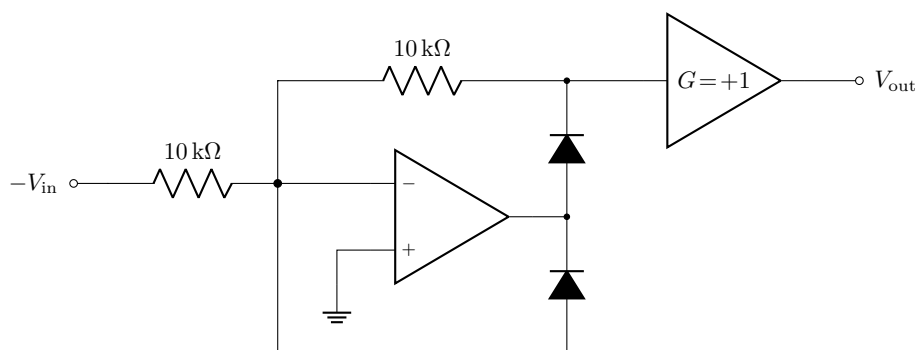
Below are two active rectifiers; that is, ideally, they realize the function

$$V_{\text{out}} = \begin{cases} V_{\text{in}}, & V_{\text{in}} \geq 0 \\ 0 & \text{otherwise.} \end{cases} \quad (7.132)$$

The first one is a “simple” rectifier,



and the second one is the “better” rectifier.<sup>14</sup>



Note that the second circuit takes the *inverted* signal  $-V_{\text{in}}$  instead of  $V_{\text{in}}$ , which could be implemented by another inverting amplifier that is not shown.

The questions are:

1. Why are these active rectifiers? (Note that unlike simple diodes, these circuits really make a transition at  $V_{\text{in}} = 0$ , rather than at one forward diode-drop above ground.)
2. Why is the “better” circuit better? (*Hint*: it has to do with the slew rate; what is the state of the op-amps when  $V_{\text{in}} < 0$ ?)

**Solution.** Tracing through the simple rectifier: Note that if  $V_{\text{in}} \geq 0$ , then the output of the op-amp can maintain the inverting-input voltage  $V_{\text{in-}} = V_{\text{in}}$  by keeping its output at one forward diode drop above  $V_{\text{in}}$ . However, if  $V_{\text{in}} < 0$ , then the op-amp can’t pull  $V_{\text{in-}}$  negative through the diode, so the op-amp rails negative. The output is taken from  $V_{\text{in-}}$ , so the forward diode-drop doesn’t matter for calculating  $V_{\text{out}}$ .

Tracing through the better rectifier: If  $-V_{\text{in}} \leq 0$  (i.e.,  $V_{\text{in}} \geq 0$ ), then the output of the op-amp can maintain the inverting-input voltage  $V_{\text{in-}} = 0$  by pulling its output one diode drop above  $V_{\text{in}}$  and conducting via the upper diode, and the op-amp acts like an inverting amplifier. Again, the diode drop doesn’t matter since the output is buffered at the correct point. For  $-V_{\text{in}} \geq 0$  (i.e.,  $V_{\text{in}} \leq 0$ ), the lower diode conducts instead, and the upper diode disconnects—this means that  $V_{\text{out}} = 0$  because the output buffer sees the virtual ground at the inverting input to the op-amp. The output of the op-amp is one forward diode-drop below ground.

The difference in these circuits is in the crossing through zero, because in the former case, the op-amp output swings from the negative rail to  $+0.7\text{ V}$ . In the latter case, the op-amp only swings from  $-0.7\text{ V}$  to  $+0.7\text{ V}$ , which reduces the tendency of the op-amp to glitch when fast input signals have zero-crossings.

<sup>14</sup>Paul Horowitz and Winfield Hill, *The Art of Electronics*, 2nd ed. (Cambridge, 1989), pp. 187-8 (ISBN: 0521370957).

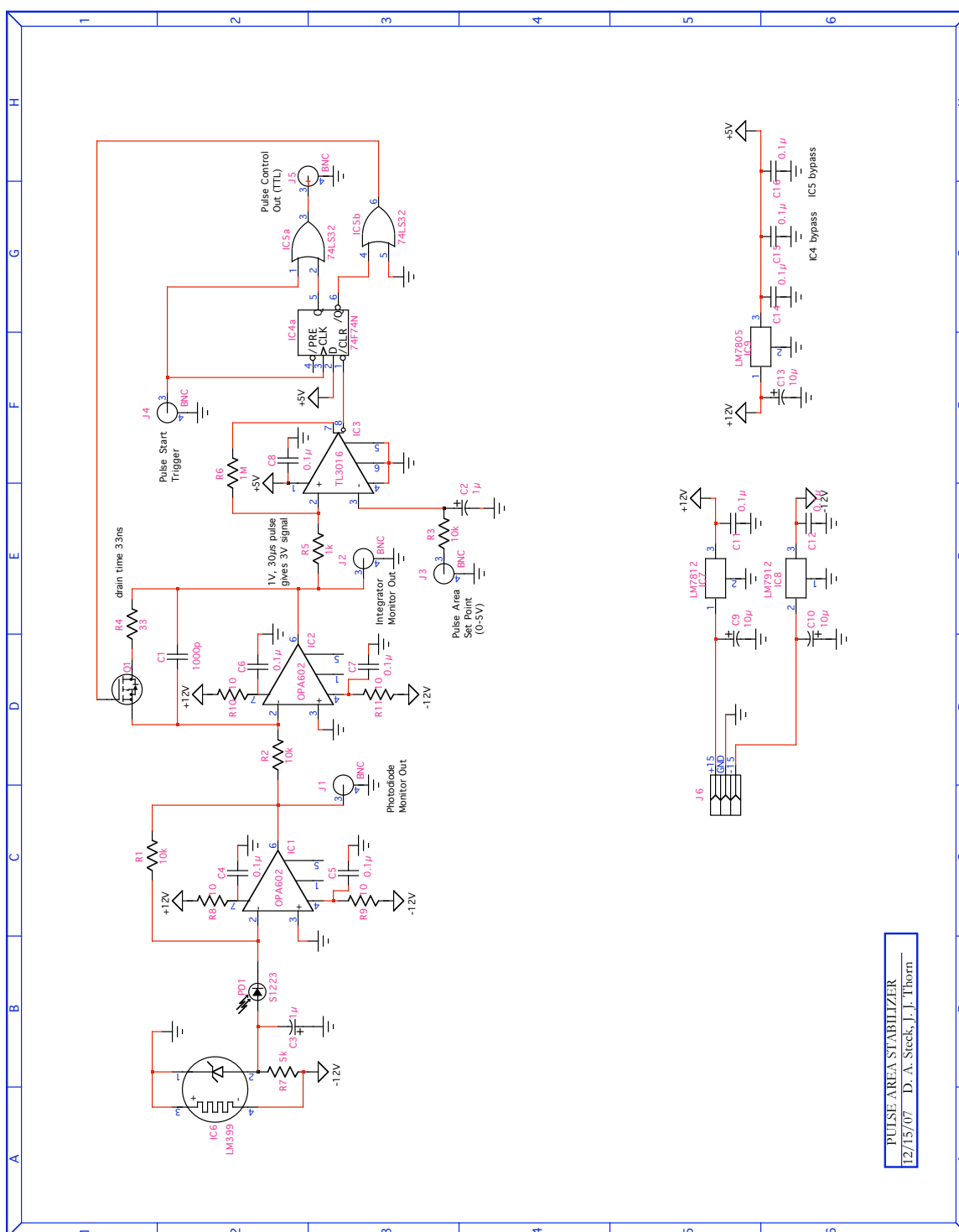


### 7.12.5 Pulse-Area Stabilizer

The circuit on the next page is designed for the following purpose. To take photographs with a laser pulse, it is desirable to have the same exposure from each pulse. But the intensity of the laser drifts. Rather than try to stabilize the intensity of the laser, we can compensate for the drift by changing the *duration* of each laser pulse to compensate. By making the **pulse area** or integrated energy of each pulse the same, the photographs have exactly the same exposure, independent of the laser intensity.

Try to trace through the following features in the circuit.

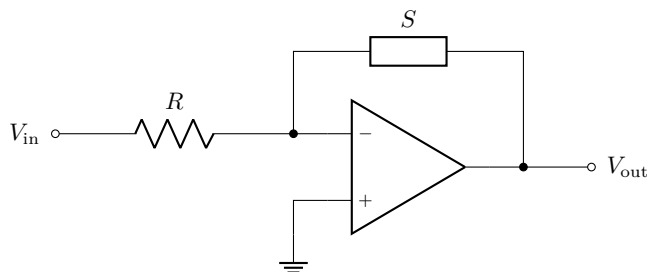
1. The reference IC6 provides a stable voltage to bias the photodiode. What is the voltage at pin 2?
2. You should then convince yourself that the photodiode PD1 is reverse-biased. This helps to improve the speed of the photodiode. The photodiode itself acts as a current source, with current flowing from cathode to anode.
3. What kind of op-amp circuit is IC1? How is the output related to the photodiode signal? (Answer: the op-amp output is positive and proportional to the photocurrent.)
4. IC2 is an integrator, with a MOSFET to reset the integrating capacitor.
5. IC3 is a comparator, connected as a Schmitt trigger. It detects when the integral of the laser pulse intensity reaches a set value from input jack J3. Its output drives digital logic circuitry, to which we will return later after we have studied digital electronics.



## 7.13 Exercises

### Problem 7.1

Consider the following op-amp circuit, which contains one resistor (of resistance  $R$ ), and one schmesistor.



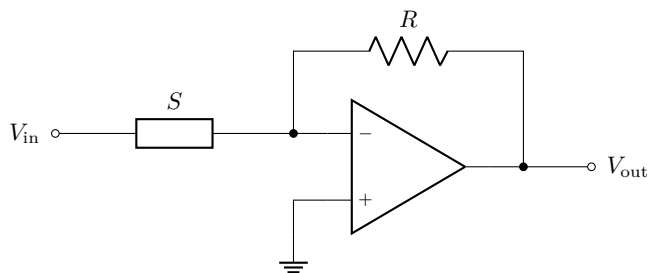
Recall (Problem 1.15) that a “schmesistor” is a device that obeys “Schmohm’s law,”

$$V = I^2 S, \quad (7.133)$$

where  $S$  is the “schmesistance.” Derive an expression for  $V_{\text{out}}$  in terms of  $V_{\text{in}}$ ,  $R$ , and  $S$ .

### Problem 7.2

Consider the following op-amp circuit, which contains one resistor (of resistance  $R$ ), and one “schmesistor.”



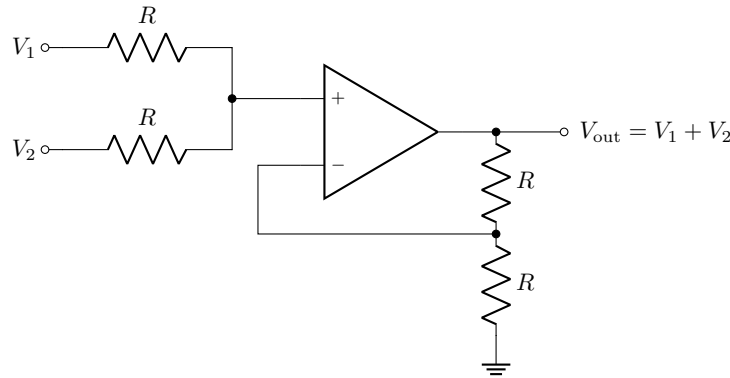
Recall that a “schmesistor” is a device that obeys “Shmohm’s law,”

$$V = I^2 S, \quad (7.134)$$

where  $S$  is the “schmesistance.” Derive an expression for  $V_{\text{out}}$  in terms of  $V_{\text{in}}$ ,  $R$ , and  $S$ .

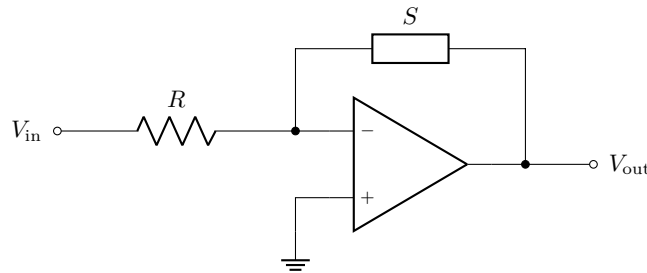
### Problem 7.3

Show that  $V_{\text{out}} = V_1 + V_2$  in the circuit below.



#### Problem 7.4

(a) Consider the following op-amp circuit, which contains one resistor (of resistance  $R$ ), and one schmapacitor (see Problem 2.9) of schmapacitance  $S$ .



Derive an expression for  $V_{\text{out}}(t)$  in terms of an arbitrary, time-dependent input  $V_{\text{in}}(t)$ ,  $R$ , and  $S$ . Assume an ideal op-amp.

Recall that the schmapacitor is defined by the relation

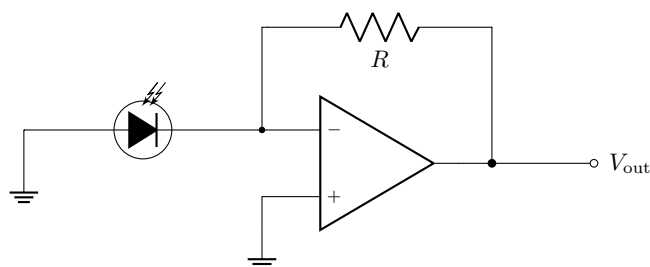
$$I = S \frac{d^3 V}{dt^3}, \quad (7.135)$$

where  $I$  is the schmapacitor current and  $V$  is the voltage drop across the schmapacitor.

(b) Of course, schmapacitors aren't real. Describe *briefly* and *qualitatively* how you could build an equivalent circuit using real-world components.

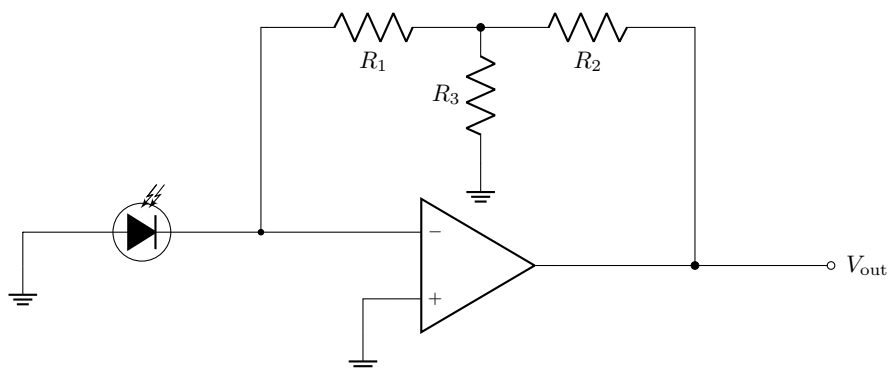
#### Problem 7.5

(a) A photodiode produces a *backwards* current (i.e., current flows from cathode to anode) when detecting light. (Think of this as the opposite of an LED, where a *forward* current causes light to be *emitted*. In fact, an LED can work as a photodiode, though not a particularly great one.) Consider the photodiode-amplifier (op-amp) circuit below, which acts as a **transimpedance amplifier** (current input, voltage output). Write down an expression for the output voltage in terms of the photodiode current and  $R$ . Is the output voltage positive or negative when you shine light on the photodiode?



(b) The Hamamatsu S1223 is a standard, medium-area ( $2.4\text{ mm} \times 2.8\text{ mm}$ ), general-purpose, silicon PIN photodiode. The sensitivity is specified at  $0.52\text{ A/W}$ . Assuming the photodiode collects all the power from a steady,  $1\text{-}\mu\text{W}$  laser beam, and the resistor is  $R = 10\text{ k}\Omega$ , what is the output voltage?

(c) For a very sensitive circuit (i.e., to register small input powers, of the order of  $\text{nW}$ ), it may be necessary to use a very large resistor. But also recall that op-amps aren't always happy with feedback resistances much over  $1\text{ M}\Omega$ , and very large resistors may be difficult to source. There is a nice trick to get around this, however. Consider the modified transimpedance amplifier below, with a "T network" in the feedback loop.



What is the effective feedback resistance  $R$  of this circuit (i.e., the resistance that makes this equivalent to the original circuit)? Find a combination of  $100\text{ k}\Omega$  or smaller resistors that gives an effective  $R$  of  $1\text{ G}\Omega$ .

### Problem 7.6

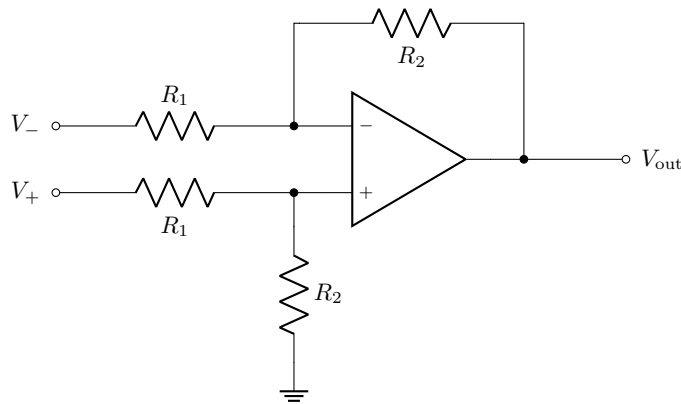
Design a current source (constant-current regulator), according to the following specs:

1. The only allowed parts are: one (ideal) op-amp, one resistor, and the load.
2. The load current is controlled by an input voltage, with a  $1\text{-V}$  input change corresponding to a  $1\text{-mA}$  change in output current.
3. The load need not be ground-referenced (i.e., the load need not have any direct connection to ground, or to any particular voltage).

*Hint:* think about how the standard inverting op-amp circuit works.

### Problem 7.7

Consider the (op-amp) differential-amplifier circuit shown below. Recall that for this to behave as a good differential amplifier (i.e., for perfect common-mode rejection), the two  $R_1$ – $R_2$  resistor pairs must be matched perfectly (in terms of ratio), assuming ideal op-amp behavior.

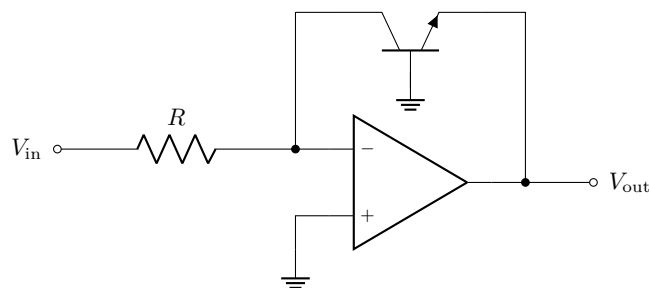


Of course, real resistors aren't perfectly matched. As a model for this, suppose that the feedback resistor has a resistance  $R_2 + \delta R$ , where  $\delta R$  is a small perturbation to this resistance.

- Rederive an expression for  $V_{\text{out}}$  in terms of the input voltages and resistances. Keep only first-order terms in  $\delta R$ .
- Write down an expression for the CMRR. Consider a unity-gain amplifier with  $R_1 = R_2 = 10 \text{ k}\Omega$ . Give a numerical estimate (in dB) for the expected CMRR if you use 1% resistors. Repeat for 0.01% resistors.

### Problem 7.8

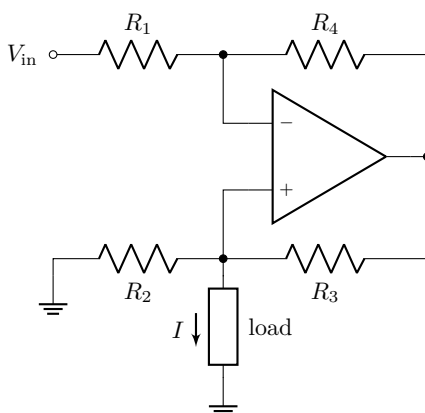
Show that the op-amp circuit below behaves approximately as a **logarithmic amplifier**. Under what conditions does this circuit really function logarithmically?



*Hint:* the function of the transistor, as in the inverting amplifier, is to convert a current into a voltage, so use an appropriate relation to describe the transistor.

### Problem 7.9

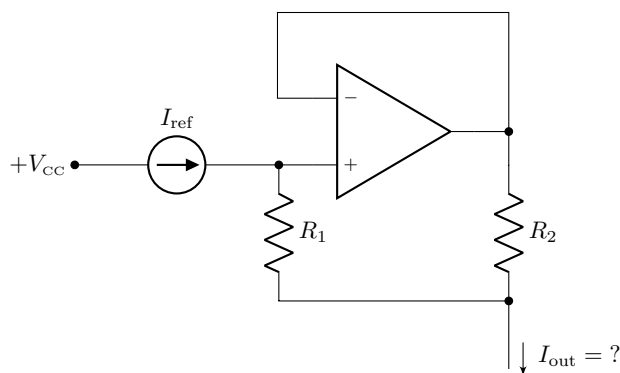
For the circuit below, the **Howland current source**, show, provided  $R_4/R_1 = R_3/R_2$ , that  $I = -V_{\text{in}}/R_2$ .



Assume an ideal op-amp.

### Problem 7.10

The circuit shown below is an op-amp current source, based on the reference current source  $I_{\text{ref}}$ .

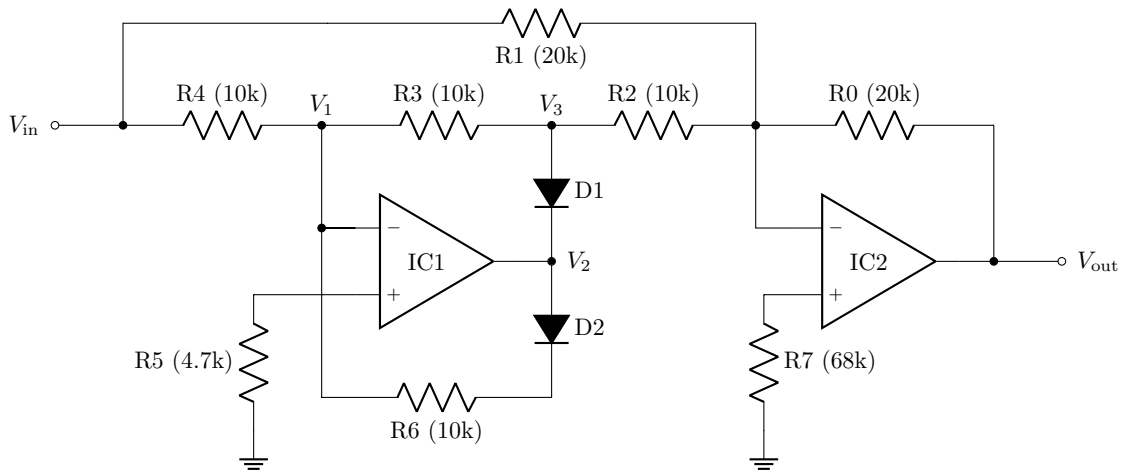


- Compute  $I_{\text{out}}$ , treating the op-amp as ideal.
- Treating the op-amp as having a **finite** open-loop gain  $A$ , derive an expression for  $I_{\text{out}}$  in terms of  $I_{\text{ref}}$ ,  $R_1$ ,  $R_2$ , and the (unknown) voltage  $V$ . In your solution you should make it clear that  $I_{\text{out}}$  depends **weakly** on  $V$  for large  $A$ .
- Take the  $A \rightarrow \infty$  limit to show that the dependence on  $V$  (and  $A$ ) disappears, and that you recover your result from (a).
- For finite  $A$ , the dependence of  $I_{\text{out}}$  on  $V$  implies a finite effective output impedance of this current-source circuit. Derive an expression for this effective output impedance (that appropriately relates a small change in voltage to a small change in current). Is it large or small for a decent op-amp?

### Problem 7.11

Consider the circuit below.<sup>15</sup> Show that for this circuit,  $V_{\text{out}} = |V_{\text{in}}|$ .

<sup>15</sup>Miles A. Smither, "Improved absolute-value circuit," in Bill Furlow, Ed., *Circuit Design Idea Handbook* (Cahners Books, 1974), p. 13 (ISBN: 0843602058).

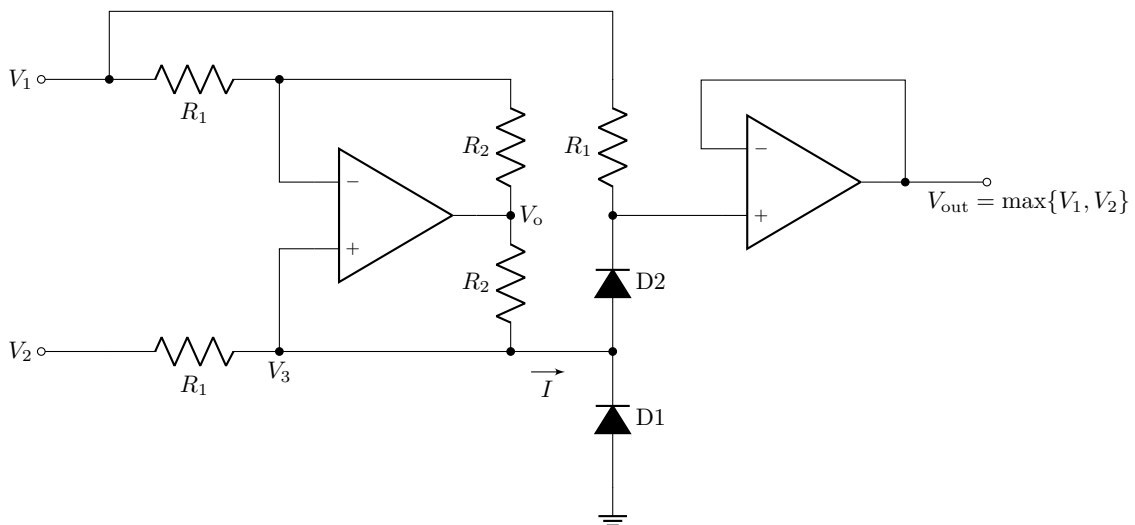


This circuit looks a little complicated, so here is some guidance. Note that you should be able to do this with very little math, provided you break the circuit down into manageable parts that you have already learned about.

- The IC2 op amp is connected in one of the basic op-amp circuits; what is it?
- Now the tricky part is understanding the diode network in the IC1 circuit. Begin by treating D1 and D2 as ideal diodes (i.e., no forward voltage drop). Now you should see that only one diode conducts at a time, and in either case IC1 is connected as one of the basic op-amp circuits (which one?). Work out the voltages  $V_1$ ,  $V_2$ , and  $V_3$ , and handle separately the cases where the output of IC1 is positive or negative.
- Finally, consider the original circuit with *real* diodes, and argue that the forward voltage drops don't matter.

### Problem 7.12

The purpose of this problem is to show that, in the circuit below, that  $V_{\text{out}} = \max\{V_1, V_2\}$  (i.e., it selects the larger of the input voltages). The diodes here are *real* diodes (i.e., you must account properly for any voltage drops across the diodes), and there are no restrictions on the signs of the input voltages.





(a) Begin by showing that the current  $I$  is given by

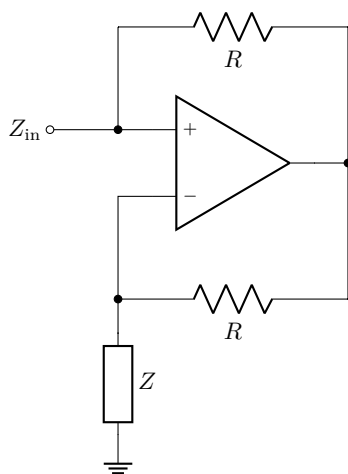
$$I = \frac{V_2 - V_1}{R_1}. \quad (7.136)$$

(The voltage labels  $V_3$  and  $V_o$  are there for a reason!)

(b) Then use the current  $I$  to show that  $V_{\text{out}} = \max\{V_1, V_2\}$ .

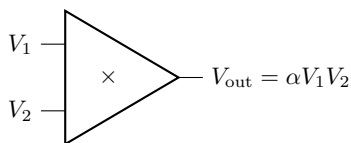
### Problem 7.13

The circuit below is a **negative-impedance converter**, composed of an op-amp (you may assume the golden rules apply here), two identical resistors of resistance  $R$ , and a generic circuit element of impedance  $Z$  (could be a resistor, capacitor, etc.). What is the impedance  $Z_{\text{in}}$  at the input terminal?



### Problem 7.14

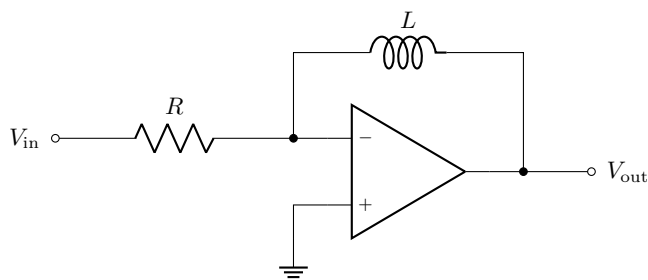
The amplifier below is a multiplying amplifier, with the output voltage proportional to the product of the input voltages.



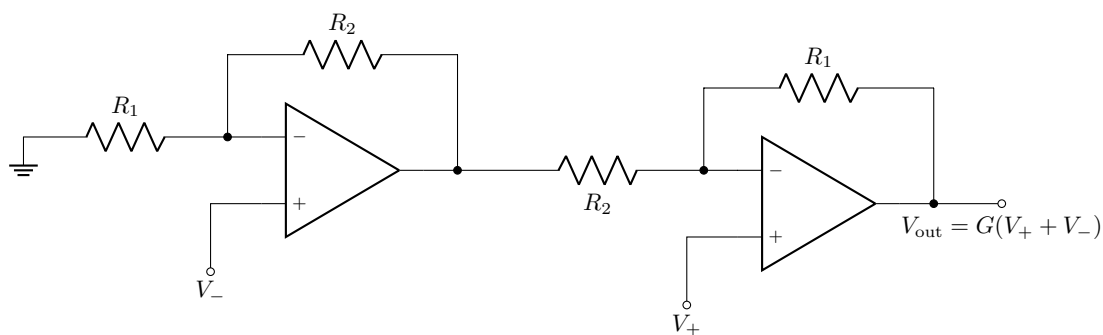
Using this multiplier and an op-amp, design a circuit that behaves as a **square-root amplifier** (i.e., with  $V_{\text{out}} \propto \sqrt{V_{\text{in}}}$ ). Show that your circuit behaves as advertised.

### Problem 7.15

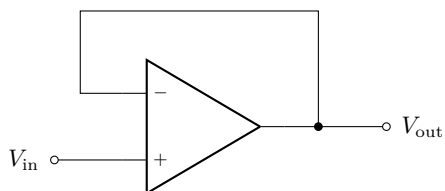
Consider the circuit below. Design an equivalent circuit using an op-amp, the **same** resistor  $R$ , and a **capacitor**, and give the value  $C$  of the capacitance in terms of  $R$  and  $L$ .

**Problem 7.16**

Show that the circuit below behaves as an instrumentation amplifier with high input impedance and gain  $G = 1 + R_1/R_2$ .

**Problem 7.17**

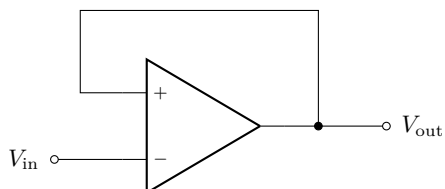
An op-amp unity-gain buffer is shown below.



Compute the closed-loop gain of this circuit, assuming a finite, open-loop gain  $A$ .

**Problem 7.18**

Consider the **flawed** op-amp circuit below.



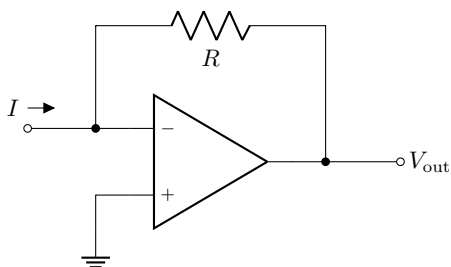
- What is  $V_{out}$ , if you (naïvely) apply the op-amp golden rules? Be *brief*.
- What is the flaw that makes the circuit not work according to the golden rules? Be *brief*.

(c) Show *mathematically* that the circuit doesn't work according to the golden-rule solution. [Use the op-amp relation  $V_{\text{out}} = A(V_{\text{in}+} - V_{\text{in}-})$ , where  $A$  is the open-loop gain, and consider what the op-amp "wants" to do in the presence of a small output error  $\epsilon$ .] What does the op-amp do instead?

(d) Go through the stability analysis of part (c) again, but take the error to be added to the finite- $A$  solution for  $V_{\text{out}}$ . Explain any difference you find in your solution here from that of part (c).

### Problem 7.19

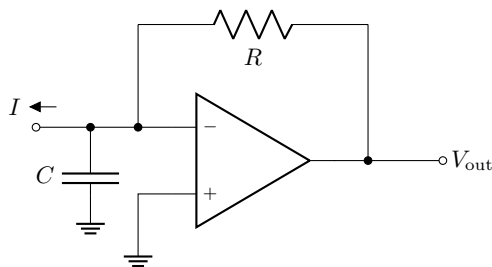
Consider the transimpedance amplifier below. The op amp has **finite** open-loop gain  $A$ .



- Derive an expression for  $V_{\text{out}}$  in terms of the input current  $I$ . Also take the limit as  $A \rightarrow \infty$ .
- Derive an expression for the input impedance of the circuit. Ignore any intrinsic input impedance  $R_i$  of the op-amp inputs.
- Derive an expression for the output impedance of the circuit.

### Problem 7.20

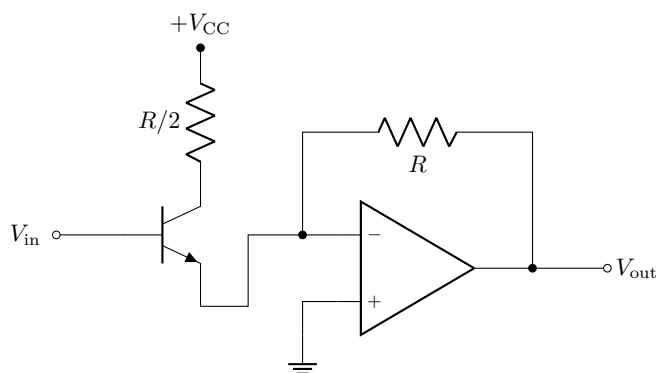
Consider the transimpedance amplifier below with input capacitance  $C$ . This circuit acts as a model for instability in a photodiode amplifier, where recall that a photodiode acts as a current source, and  $C$  models the photodiode capacitance.



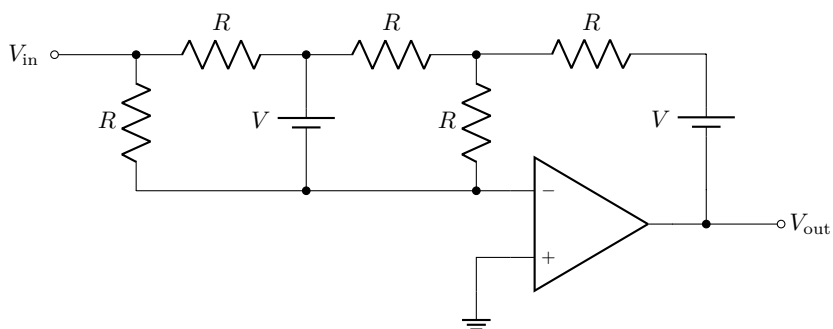
- Derive an expression for the frequency-dependent transimpedance  $Z(\omega)$  of the amplifier [i.e., such that  $V_{\text{out}} = IZ(\omega)$  for an input current of frequency  $\omega$ ], assuming a **finite** open-loop op-amp gain  $A$ . Also take the limit as  $A \rightarrow \infty$ .
- Show that, if the open-loop gain falls off like a low-pass filter,  $A(\omega) = A_0/(1 - i\omega/\omega_0)$ , that the magnitude of  $Z(\omega)$  is peaked at some frequency  $\omega > 0$  if  $A_0$  is sufficiently large.
- Find the peak frequency.

### Problem 7.21

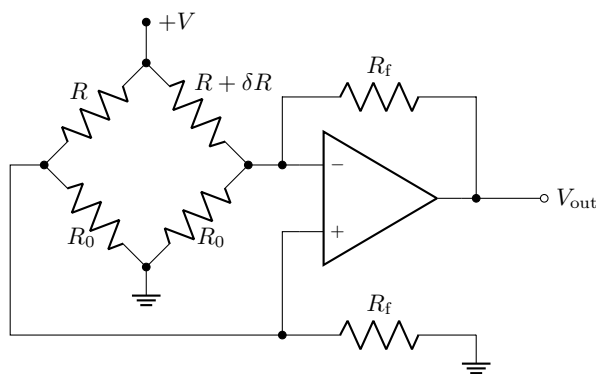
Derive an expression for  $V_{\text{out}}$  in the circuit below, in the regime where the output is not railed (assume that  $R$  is of the order of  $10\text{ k}\Omega$ ). Also, you may assume  $V_{\text{in}} \geq 0$  (but assume  $V_{\text{in}}$  is not so large that it damages the transistor) and an ideal op amp. Finally, the op amp is powered from  $\pm V_{\text{CC}}$ . Be clear about any assumptions you make.

**Problem 7.22**

In the circuit below, derive an expression for  $V_{\text{out}}$ , assuming an ideal op amp.

**Problem 7.23**

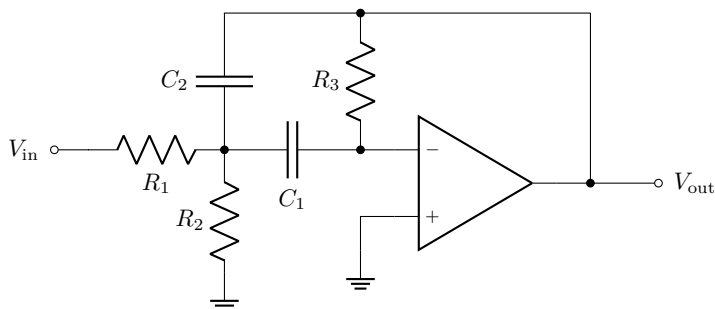
In the circuit below, the diamond-shaped arrangement of resistors is called a **Wheatstone bridge**, and can act as a sensitive measure of the mismatch of two resistors (here  $R$  and  $R + \delta R$ ). This is useful, for example, in sensing the value of a thermistor (semiconductor resistor whose resistance varies with temperature).



For this circuit, derive an expression for  $V_{\text{out}}$ . To keep the algebra under control, give the answer only to lowest order in the small perturbation  $\delta R$  [i.e., give the lowest-order term in the Taylor expansion about  $\delta R = 0$ ; you will need to use  $(a + \delta a)^{-1} \approx a^{-1} - \delta a/a^2$  for  $\delta a \ll a$ ].

**Problem 7.24**

Consider the following circuit, an active, inverting, band-pass filter. (This is the bandpass filter from the guitar preamplifier in Section 7.12.3.)



(a) Derive an expression for the gain function  $\tilde{G}(\omega)$ , and the amplitude gain function  $G(\omega)$ . Note that these are the same as the transfer function and amplitude transfer function, respectively, in the passive-filter case.

*Note:* an algebra program like *Mathematica* will make this problem considerably simpler. Contact me if you need help with learning how to use it for algebra or making plots.

(b) Derive asymptotic expressions for  $G(\omega)$  for large and small frequencies. From these expressions, argue that this is a band-pass filter.

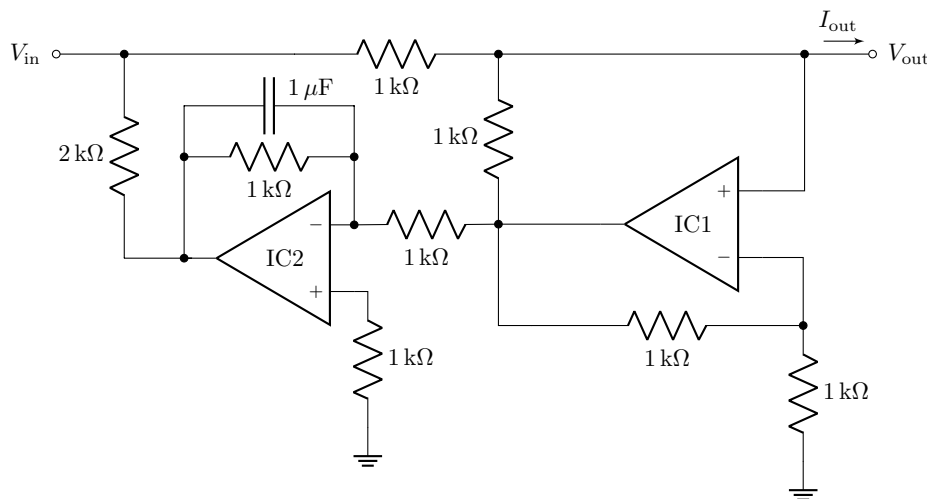
(c) Find an expression for the center frequency, where  $G(\omega)$  is maximum. Give a numerical value for this frequency (give the frequency in Hz, not rad/s), for the component values  $R_1 = 35.7\text{ k}\Omega$ ,  $R_2 = 28\text{ k}\Omega$ ,  $R_3 = 82\text{ k}\Omega$ ,  $C_1 = C_2 = 0.01\text{ }\mu\text{F}$ .

(d) Find an expression for the peak gain, i.e., the value of  $G(\omega)$  at the frequency you derived in part (c). Again, give a numerical value for the same component values.

(e) Make a log-log plot of  $G(\omega)$ , over a reasonable range of frequencies, for the component values above. (Keep in mind this is intended as an audio band-pass filter.)

**Problem 7.25**

Consider the circuit below, which is intended as a voltage-controlled current source.<sup>16</sup> Assume the “output voltage”  $V_{\text{out}}$  is held to a fixed voltage by an external source.

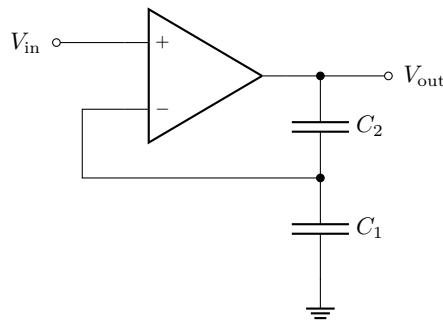


<sup>16</sup>This is the modulation input stage of the current-controller circuit in K. G. Libbrecht and J. L. Hall, “A low-noise high-speed diode laser current controller,” *Reviews of Scientific Instruments* **64**, 2133 (1993).

- (a) First, start by identifying the basic op-amp circuits for IC1 and IC2 (i.e., these are two standard op-amp circuits, connected via a network that includes three resistors). For IC2, first think about the circuit at dc (i.e., ignoring the capacitor).
- (b) Show that  $I_{\text{out}} = V_{\text{in}}/(1\text{ k}\Omega)$ , and is independent of  $V_{\text{out}}$ , for dc inputs. (You may find it useful to label the output voltages of IC1 and IC2 as  $V_1$  and  $V_2$ , respectively.)
- (c) Show that if  $V_{\text{in}}$  is disconnected (i.e., not held at any particular voltage), that  $I_{\text{out}} = 0$ , independent of  $V_{\text{out}}$  (for dc inputs).
- (d) The function of the capacitor is as follows. The above current regulation requires IC2 to generate a signal to cancel any currents drawn by IC1. However, due to propagation delays through IC2, the cancellation may not be accurate at high frequencies, and in the worst case, the circuit may even become unstable. Thus, the capacitor is there to roll off the gain of IC2, protecting against these effects. However, the cancellation no longer works, so redo (c), calculating  $I_{\text{out}}$  in terms of  $V_{\text{out}}$  at high frequencies (i.e., assume that  $I_{\text{out}}$  and  $V_{\text{out}}$  are the amplitudes of high-frequency, oscillating signals).

### Problem 7.26

Consider the circuit below.

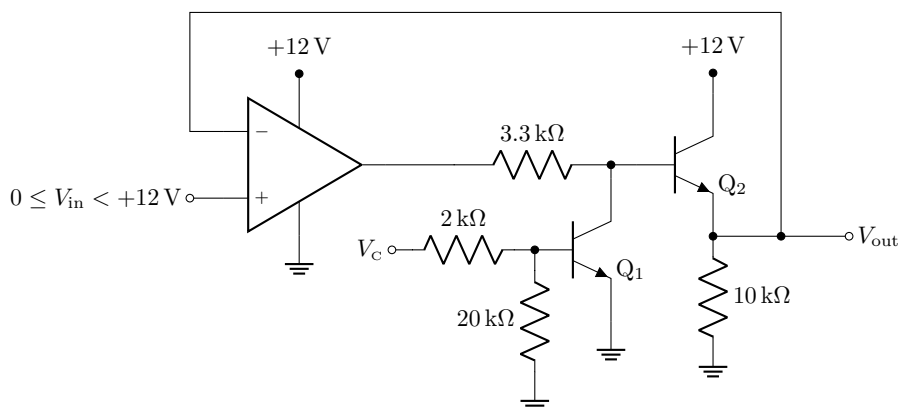


- (a) Assuming an *ideal* op-amp, what kind of amplifier is this? Compute  $V_{\text{out}}$  in terms of  $V_{\text{in}}$ .
- (b) For a *real* op-amp, what would the circuit do [i.e., why wouldn't it work as in (a)], and what *specifically* is it about the op-amp that causes the circuit to misbehave?

### Problem 7.27

The circuit below is an **analog switch**,<sup>17</sup> meaning that either  $V_{\text{out}} \approx 0$  or  $V_{\text{out}} = V_{\text{in}}$ , depending on the state of  $V_{\text{C}}$ . The point of this problem is to figure out what input corresponds to which state. ( $Q_1$  works as a switch; you may assume that the circuit is sanely designed, which means the op-amp is comfortable with the single-ended power supply, and the resistances are such that the  $Q_1$  is either completely OFF or saturated ON.)

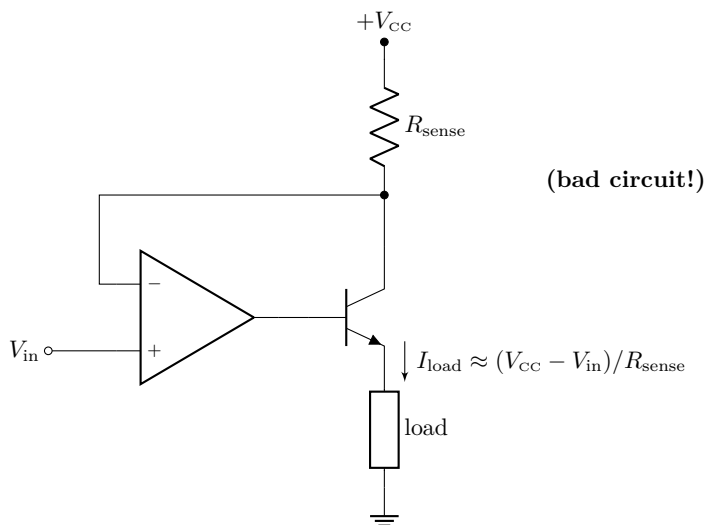
<sup>17</sup>This circuit from the Texas Instruments LF411 data sheet, <http://www.ti.com/lit/ds/symlink/lf411-n.pdf>.



- (a) What is  $V_{out}$  if  $V_C = 0$  V? What is the state of  $Q_1$ , and what does  $Q_2$  do? Is the op-amp railing or in negative-feedback mode?
- (b) What is  $V_{out}$  if  $V_C = +12$  V? What is the state of  $Q_1$ , and what does  $Q_2$  do? Is the op-amp railing or in negative-feedback mode?

### Problem 7.28

Consider this op-amp/BJT current-source circuit.



- (a) Assuming the op-amp golden rules apply to this circuit, show that the expression for the current is approximately correct, and give a *better* expression that includes the transistor  $\beta$ .
- (b) As drawn, the circuit does not work; rather, the op-amp just rails without regulating the current. Explain why.
- (c) Propose two fixes for this circuit; at least one of these fixes should not involve any component changes.

### Problem 7.29

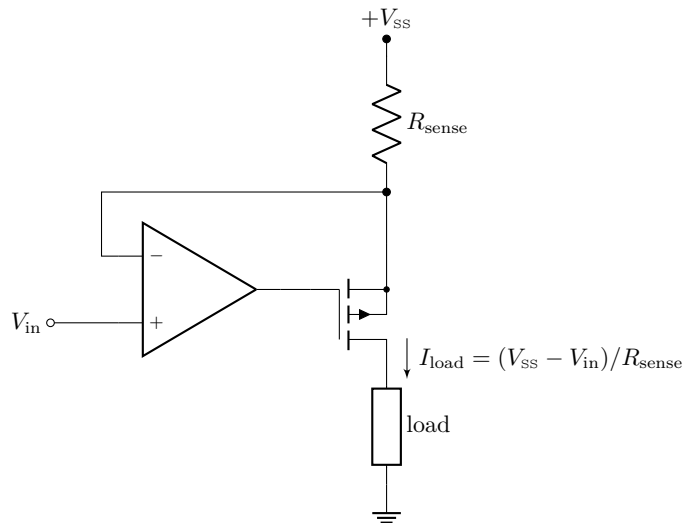
In this problem, you will analyze a schematic for a precision current source that powers a laser diode (see p. 6 in the referenced pdf file).<sup>18</sup> You will probably need to look up data sheets for various components in analyzing this circuit.

<sup>18</sup>Designed by Todd Meyrath, <http://george.ph.utexas.edu/~meyrath/informal/laser%20diode.pdf>

- (a) Suppose  $R_{\text{limit}} = 10\text{ k}\Omega$ . What is the maximum possible (dc) voltage at pin 3 of the AD820? Assume the 2N7000 is an open circuit. This is the control voltage that sets the current through the laser diode. Note that the 20-k $\Omega$  wiper of the (course adjust) pot is bypassed by a capacitor to ensure the control voltage has little high-frequency noise.
- (b) Now look at the relay on the left-hand side of the schematic. This is the “laser enable” part of the circuit, which either passes or overrides (i.e., sets to zero) the control voltage. Note the two *momentary* switches (i.e., they are push-buttons, only connected while you are actually pushing them but the circuit “remembers” the last one you pushed). Explain how this section of the circuit works to enable and disable the laser. Include the operation of the status LED (why doesn’t the LED burn out if it is powered by 15 V?), and explain why this circuit has a “soft-start” feature.
- (c) Now analyze the AD820, which regulates the diode current according to the control voltage at pin 3. For the purposes of this analysis, ignore C0, C1, and R1 (i.e., replace them by open connections) as well as R0, R2, and C2 (i.e., replace them by short circuits). These are needed for stability, but I haven’t discussed much about this yet. Thus, the feedback loop consists of the BUF634’s, resistor SR10, and the INA128 (assume the gain of the INA128 is 10). What is the purpose of the BUF634’s, why are there 2 of them, and what is with the 10  $\Omega$  resistors? What is the maximum current through the laser diode given your answer in (a)? *Explain.*
- (d) Look at the power-supply connections of the various amplifiers. Explain the differences in the bypass circuits, and *why* the chips are bypassed in different ways (i.e., why not bypass them all in the same way?).
- (e) Now go back to the AD820, and consider it separately from the other ICs, but this time *with* R0-R2, and C0-C2. Also assume pin 3 is grounded. What is the gain at low frequencies? High frequencies?

### Problem 7.30

Consider the op-amp current source shown below.



- (a) Show that

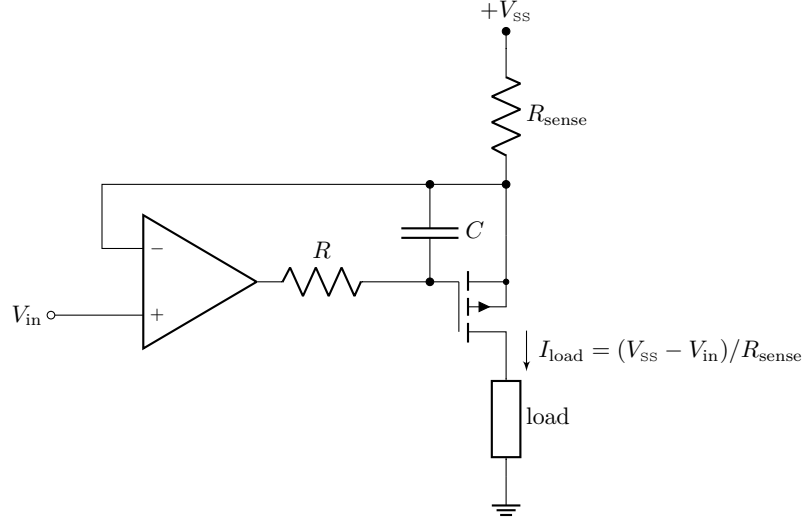
$$I_{\text{load}} = \frac{V_{\text{SS}} - V_{\text{in}}}{R_{\text{sense}}}, \quad (7.137)$$

as advertised.

- (b) Note that the MOSFET is a p-channel device. Explain why an n-channel device doesn’t work here (consider how the output changes as  $I_{\text{load}}$  increases).



(c) To analyze the stability issues in this circuit, consider the modified circuit below, with output resistor  $R$  (modeling the intrinsic output resistance of the op-amp), and capacitor  $C$ . The capacitor models the gate-source capacitance of the FET; for simplicity we are excluding the gate-drain capacitance, although it could be much more important due to the Miller effect, but this extra capacitance could be modeled in the following analysis with another capacitor (though some details of the load would have to be included in the analysis).



Consider a small ac signal  $v_{in}$  at frequency  $\omega$  at the input (and some dc bias  $V_0$ , so  $V_{in} = V_0 + v_{in}$ ). Treat the FET via its transconductance  $g_m$  (i.e.,  $i_{SD} = -g_m v_{GS}$ ) and the op-amp via the finite-gain formula (7.39)

$$V_{out} = \tilde{A}(V_{in+} - V_{in-}), \quad (7.138)$$

with response

$$\tilde{A}(\omega) = \frac{A_0}{1 - i\omega/\omega_c}, \quad (7.139)$$

where  $\omega_c$  is the cutoff frequency. Show that the ac load response  $i_{load}$  to  $v_{in}$  can be written

$$i_{load} = \frac{-g_m A_0}{(1 - i\omega RC)(1 - i\omega/\omega_c) + g_m(1 + A_0 - i\omega/\omega_c)R_{sense}} v_{in}. \quad (7.140)$$

(d) Note that the expression above can change sign, which we will see can lead to instability. Intuitively, the op-amp output is inductive, forming a destabilizing resonance with the capacitive gate input of the FET. Show that the circuit exhibits resonance behavior by finding the extreme value of  $|i_{load}|$ , as well as the frequency at which the extremum occurs. Discuss the dependence of the peak frequency on the FET parameters  $g_m$  and  $C$  (simplify your analysis by working in the regime  $A_0 \gg 1$  and  $\omega_c \ll 1/RC$ ). Also simplify your notation by defining the frequency  $\omega_0$  by

$$\omega_0^2 := \frac{\omega_c}{RC}. \quad (7.141)$$

(e) To more precisely consider stability problems in this circuit, we must consider the *voltage* feedback. To do this we should consider the voltage  $v_s$  at the inverting input, given by

$$v_s = -i_{load}R_{sense} =: G(\omega)v_{in}, \quad (7.142)$$

to be the circuit “output,” which then defines the voltage gain  $G$  of the circuit. Since this is fed back to the inverting input,  $G$  should be positive and preferably large. The problem comes when  $\text{Re}[G]$  is

*negative*, because it means the feedback is *positive*. In particular, the circuit becomes unstable when  $\text{Re}[G(\omega)] < -1$ , because perturbations at any frequency satisfying this condition will grow. Find the extreme value of  $\text{Re}[G]$  and show that  $\text{Re}[G] < -1$  when

$$\frac{A_0 \omega_c RC}{\sqrt{g_m R_{\text{sense}}(\sqrt{4A_0 \omega_c RC} + \sqrt{g_m R_{\text{sense}}})}} > 1, \quad (7.143)$$

again in the limits  $A_0 \gg 1$  and  $\omega_c \ll 1/RC$ . Using typical parameters for a high-speed op-amp and a power MOSFET ( $\omega_c/2\pi = 100$  Hz,  $R = 50 \Omega$ ,  $R_{\text{sense}} = 50 \Omega$ ,  $g_m = 1 \text{ U}$ ,  $A_0 = 5 \times 10^6$ ) solve the equation to find the range of gate-source capacitance that destabilizes the circuit. Since typical gate capacitances are on the order of 100 pF, you should find that the circuit is unstable (i.e., it oscillates). More components in the feedback loop are necessary to stabilize it.

Again, to keep your calculations organized, you should use the following definitions to simplify your calculation:

$$\begin{aligned} G_0 &:= g_m R_{\text{sense}} A_0 \omega_0^2 \\ \omega_1 &:= \omega_0 \sqrt{1 + g_m R_{\text{sense}}(1 + A_0)} \\ \gamma_1 &:= \frac{\omega_c^2 + (1 + g_m R_{\text{sense}}) \omega_0^2}{\omega_c}. \end{aligned} \quad (7.144)$$



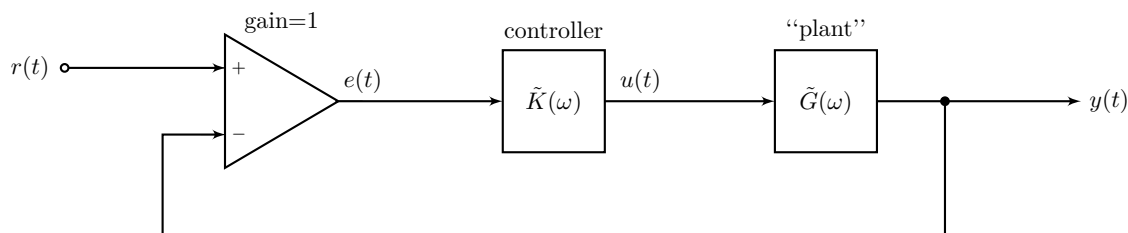
# Chapter 8

## PID Control

A major theme in our study of op-amps is that **feedback**, and *negative* feedback in particular, is a useful tool for improving the behavior of amplifiers. It is also useful in the realization of circuits that would otherwise be complex or difficult to implement (the logarithmic amplifier is a good example; see Problem 7.8). Of course, feedback is similarly useful beyond op-amp circuits, and we will consider feedback control more generally as a tool for maintaining systems in a desired “target” state. As it turns out, op-amp circuits are useful in realizing one of the popular, general-purpose control methods, **PID control**, which we will define shortly.<sup>1</sup>

### 8.1 Basics of Linear Control

Schematically, we can represent a feedback-control system as in the diagram below.



There are several important elements that interact here.

- The **plant** is the system to be controlled. (That’s plant as in “chemical plant,” not a shrub.) We assume the plant to have an input and an output. The output is some scalar quantity that we want to control, such as temperature, position, voltage, frequency, speed, etc. The input is some other scalar “knob” by which we can affect the plant. In linear control theory, we will assume that the plant is a linear filter with transfer function  $\tilde{G}(\omega)$ .
- The **controller** is a system that analyzes the state of the plant and implements a control procedure to the plant input. Again, we will treat this as a linear filter with transfer function  $\tilde{K}(\omega)$ .
- The **goal** of the control is to make the **output signal**  $y(t)$  follow the **input signal**  $r(t)$  as closely as possible. The control system should, however, be robust to environmental perturbations, which are something like random changes in  $r(t)$ . In temperature stabilization, for example, the temperature control should be robust to fluctuations in the surrounding temperature (e.g., due to the day/night cycle).

---

<sup>1</sup>For a more complete, readable introduction to control theory, see John Bechhoefer, “Feedback for physicists: A tutorial essay on control,” *Reviews of Modern Physics* **77**, 783 (2005) (doi: <http://dx.doi.org/10.1103/RevModPhys.77.783>), available at [http://www.sfu.ca/chaos/papers/2005/rmp\\_reprint05.pdf](http://www.sfu.ca/chaos/papers/2005/rmp_reprint05.pdf).

- The **error** is defined as the difference between the desired and actual states:

$$e(t) := r(t) - y(t). \quad (8.1)$$

Because of the negative sign on  $y(t)$  here, feeding this error signal into the controller amounts to negative feedback to the plant, assuming the low-frequency transfer characteristics of the controller and plant have no phase shift.

- The **feedback signal**  $u(t)$  is the error  $e(t)$ , modified by the controller in frequency space.

Of course, more realistic systems may not be linear, and may have vector inputs and outputs. The scalar case is still important both conceptually and practically, so we will focus on only this here; however, note that the vector case can sometimes be treated well enough as several, parallel scalar loops—the nature of feedback control is to correct for errors, and so it can often tolerate some slop in the model. The nonlinear case is more complex in theory, but often a simple, pragmatic approach is to approximate the nonlinear system by a linearized version, so that linear theory applies. Again, this can sometimes work even when the approximation is quite drastic.

## 8.2 Example: First-Order Plant, Proportional Control

As a simple example, suppose we take the plant to be a first-order, low-pass filter, with

$$\tilde{G}(\omega) = \frac{G_0}{1 - i\omega/\omega_0}, \quad (8.2)$$

where  $G_0$  is the dc gain, and  $\omega_0$  is the cutoff frequency of the filter. A more concrete example where this model applies is temperature control of a room, where the input is a simple, electric heater, and the output is the room temperature. The low-pass-filter nature of the room is apparent in the exponential settling of the room temperature when the input (power setting of the heater, *not* a thermostat) changes.

For the controller we will implement simple **proportional control**, which just means that the control signal is proportional to the error. That is, we have a constant transfer function

$$\tilde{K}(\omega) = K_P, \quad (\text{proportional controller transfer function}) \quad (8.3)$$

where  $K_P$  is the **proportional gain**.

### 8.2.1 General Result: Closed-Loop Transfer Function

To analyze our simple example, we will first examine a more generally useful result. To introduce some notation, for time domain quantities like  $y(t)$ , we will denote their frequency-domain counterparts by  $\tilde{y}(\omega)$ —that is  $\tilde{y}(\omega)$  is the amplitude of the frequency  $\omega$  that is present in  $y(t)$ . Given the connections in the circuit above, we have

$$\tilde{y}(\omega) = \tilde{K}(\omega) \tilde{G}(\omega) \tilde{e}(\omega), \quad (8.4)$$

where

$$\tilde{e}(\omega) = \tilde{r}(\omega) - \tilde{y}(\omega). \quad (8.5)$$

Eliminating the error  $\tilde{e}$ , we have

$$\tilde{y}(\omega) = \tilde{K}(\omega) \tilde{G}(\omega) [\tilde{r}(\omega) - \tilde{y}(\omega)]. \quad (8.6)$$

Then

$$\tilde{y}(\omega) [1 + \tilde{K}(\omega) \tilde{G}(\omega)] = \tilde{K}(\omega) \tilde{G}(\omega) \tilde{r}(\omega), \quad (8.7)$$

so

$$\tilde{y}(\omega) = \frac{\tilde{K}(\omega) \tilde{G}(\omega)}{1 + \tilde{K}(\omega) \tilde{G}(\omega)} \tilde{r}(\omega). \quad (8.8)$$

This gives the output response  $\tilde{y}$  in terms of the input  $\tilde{r}$ . Since the control system is linear, we have derived the transfer function for the entire control system, or the **closed-loop transfer function**

$$\tilde{T}(\omega) = \frac{\tilde{K}(\omega) \tilde{G}(\omega)}{1 + \tilde{K}(\omega) \tilde{G}(\omega)}. \quad (8.9)$$

(closed-loop transfer function)

Relating this back to op-amps, in the limit where the gain product  $\tilde{K}\tilde{G}$  becomes large, the transfer function approaches unity (and is otherwise less than unity).

### 8.2.2 Frequency-Domain Solution

Returning to our example, we have  $\tilde{K} = K_P$  and  $\tilde{G}$  defined by Eq. (8.2), so the closed-loop transfer function becomes

$$\tilde{T}(\omega) = \frac{K_P G_0}{K_P G_0 + 1 - i\omega/\omega_0} = \left( \frac{K_P G_0}{K_P G_0 + 1} \right) \frac{1}{1 - i\omega/\omega_0(K_P G_0 + 1)}. \quad (8.10)$$

This is *still* the transfer function for a low-pass filter, but now the dc gain is

$$\tilde{T}(\omega = 0) = \frac{K_P G_0}{K_P G_0 + 1}, \quad (8.11)$$

compared to the original dc gain of  $G_0$ , and the control becomes ideal as  $K_P \rightarrow \infty$  (at least in this simple model; this is not true in general for *real* control systems). Also, the new cutoff frequency is  $\omega_0(K_P G_0 + 1)$ , which is larger than the original  $\omega_0$ , particularly for large  $K_P$ . Since the cutoff frequency is inversely proportional to the decay time, we can see that a larger cutoff frequency is desirable, as it means the control system “settles” more quickly.

### 8.2.3 Time-Domain Solution

To examine this settling behavior more, we can also transfer the analysis for this example into the time domain. Returning to Eq. (8.6), and putting in Eq. (8.2),

$$\tilde{y}(\omega) = \tilde{K}(\omega) \tilde{G}(\omega) [\tilde{r}(\omega) - \tilde{y}(\omega)] = \frac{K_P G_0}{1 - i\omega/\omega_0} [\tilde{r}(\omega) - \tilde{y}(\omega)]. \quad (8.12)$$

Rearranging a bit, we find

$$\tilde{y}(\omega) (\omega_0 - i\omega) = \omega_0 K_P G_0 [\tilde{r}(\omega) - \tilde{y}(\omega)]. \quad (8.13)$$

We can change this to the time domain by identifying the time-domain counterparts to each variable, and using  $\partial/\partial t \equiv -i\omega$ , to find

$$\omega_0 y(t) + \dot{y}(t) = \omega_0 K_P G_0 [r(t) - y(t)]. \quad (8.14)$$

Solving for  $\dot{y}$ ,

$$\dot{y} = -\omega_0 y + \omega_0 K_P G_0 (r - y). \quad (8.15)$$

There are two terms here. The first is a simple damping term, again with a time constant of  $1/\omega_0$ . The second term is a forcing term, where the system is “driven” by the error  $e = r - y$ . The system always tries to eliminate the error. Since it does so via simple exponential relaxation, it is always *stable*—it never “runs away” from the zero-error point. Note that the drive is stronger for larger  $K_P$ , meaning that more control has more effect on the system, as we should expect.

### 8.2.4 Constant Input and Proportional Droop

As a simpler version of this example, let's try a constant input  $r(t) = r$ . What is the steady-state solution ( $\dot{y} = 0$ )? From Eq. (8.15), we have

$$\omega_0 y_{ss} = \omega_0 K_P G_0 (r - y_{ss}), \quad (8.16)$$

and solving for the steady-state output  $y_{ss}$ , we find

$$y_{ss} = \frac{K_P G_0}{1 + \omega_0 K_P G_0} r. \quad (8.17)$$

Then we see that the (proportional) control system only achieves the goal perfectly (eventually) in the limit  $K_P \rightarrow \infty$ . This is the fundamental problem with proportional control: the controller only acts if there is error, so there *must* be some steady-state error, or **droop**, for any finite proportional gain  $K_P$ .<sup>2</sup> In op-amp circuits, the gain is large for just this reason, but for more complex, real-world control systems (electronic, mechanical, etc.), there are usually limits on  $K_P$  to maintain loop stability.

### 8.3 Integral Control

One approach to fixing the problem is to introduce an infinite gain *only at dc*, where the time delays that usually cause feedback-loop stability problems won't matter much. This is precisely what an integrator does: recall that an op-amp integrator (Section 7.4.2) has a gain of the form  $-i/\omega RC$ . More generally, integral control has a transfer function of the form

$$\tilde{K}(\omega) = \frac{iK_I}{\omega\tau}, \quad (\text{integral controller transfer function}) \quad (8.18)$$

where  $\tau$  is a time constant and  $K_I$  is the (dimensionless) integral gain (note that  $\tau$  just acts like another gain parameter here). Then noting that  $1/(-i\omega)$  is an antiderivative, the controller output is

$$u(t) = \frac{K_I}{\tau} \int_0^t dt' e(t'). \quad (8.19)$$

That is, the controller has a built-in “memory” of past error in the feedback. This allows correction of the droop, because we no longer require an error at the *present* moment to have a nonzero control signal  $u(t)$ .

#### 8.3.1 Example: First-Order Plant, Integral Control

Now back to the example that we introduced in Section 8.2. In the time domain, the same steps leading up to (8.15) now give

$$\dot{y} = -\frac{y}{\tau} + \frac{K_P G_0}{\tau^2} \int_0^t dt' [r - y(t')], \quad (8.20)$$

where we are taking  $\tau = 1/\omega_0$  and we are still assuming a constant control input  $r$ . It is more convenient to handle an ordinary differential equation, rather than an integro-differential equation, so we can differentiate this equation to obtain

$$\ddot{y} = -\frac{\dot{y}}{\tau} + \frac{K_P G_0}{\tau^2} [r - y(t)]. \quad (8.21)$$

In steady state,  $\ddot{y} = \dot{y} = 0$ , and so we have

$$y_{ss} = r, \quad (8.22)$$

which means that we obtain exactly the target in steady state: there is no droop with (ideal) integral control.

#### 8.3.2 Frequency Domain

In the frequency domain, for this example with general closed-loop transfer function (8.9) and example plant function (8.2), we have

$$\tilde{T}(\omega) = \frac{\tilde{K}(\omega) \tilde{G}(\omega)}{1 + \tilde{K}(\omega) \tilde{G}(\omega)} = \frac{1}{1 - i\omega/\omega_0 K_I - \omega^2/\omega_0^2 K_I}. \quad (8.23)$$

<sup>2</sup>For a good story of proportional droop, see the introduction to David Sellars, “An Overview of Proportional plus Integral plus Derivative Control and Suggestions for Its Successful Application and Implementation,” <http://hdl.handle.net/1969.1/5215>.

This transfer function is second order in the denominator because of the frequency dependence of the integrator, and is qualitatively different than the first-order low-pass filter that we obtained for proportional control.

To see this, compare this to a damped, forced harmonic oscillator, which has the form

$$\ddot{y} + \gamma\dot{y} + \omega_0^2 y = f(t), \quad (8.24)$$

where  $\gamma$  is a damping rate, and  $f(t)$  is a forcing function (we have set the mass to 1). In the frequency domain, this becomes

$$(-\omega^2 - i\gamma\omega + \omega_0^2)\tilde{y} = \tilde{f}(\omega), \quad (8.25)$$

and solving for  $\tilde{y}$  gives

$$\tilde{y} = \frac{\tilde{f}(\omega)/\omega_0^2}{1 - i\gamma\omega/\omega_0^2 - \omega^2/\omega_0^2}. \quad (8.26)$$

The transfer function here has the same form as for the integrator control of a low-pass filter in Eq. (8.23). In the integrator control problem, a large integral gain  $K_I$  is equivalent to a large oscillation frequency  $\omega_0$  relative to the damping rate  $\gamma$  in the harmonic-oscillator problem. This leads to underdamped oscillations, which means the controller is overshooting the target state.

Again, note that  $\tilde{T}(\omega) \rightarrow 1$  as  $\omega \rightarrow 0$ . This means that there is no steady-state droop in this integral-control example.

## 8.4 Proportional-Integral (PI) Control

It is, of course, possible to combine the benefits of proportional and integral control by using a controller with both features. The simplest way to combine these is a simple linear combination:

$$\tilde{K}(\omega) = K_P + \frac{iK_I}{\omega\tau}. \quad (8.27)$$

(PI-controller transfer function)

This is the transfer function for **proportional-integral (PI) control**. The second term gives integral control, which eliminates droop issues. The first term is a proportional term, which gives more high-frequency response, and thus faster setting. In the example of the single-pole plant from Section 8.2, if we work out the closed-loop transfer function, we obtain

$$\tilde{T}(\omega) = \frac{1 - i(K_P/K_I)(\omega/\omega_0)}{1 - i(1 + K_P)\omega/\omega_0 K_I - \omega^2/\omega_0^2 K_I}. \quad (8.28)$$

(Again  $\omega_0 = 1/\tau$  here.) Note that in the dc limit,  $\tilde{T}(\omega \rightarrow 0) = 1$ , which means there is no steady-state droop, as in the integral-control case. In the high-frequency limit,

$$\tilde{T}(\omega \rightarrow \infty) \sim \frac{-i(K_P/K_I)}{(1 + K_P)/K_I - i\omega/\omega_0 K_I} = \frac{-iK_P}{(1 + K_P) - i\omega/\omega_0}. \quad (8.29)$$

In the high-frequency limit, the transfer function reduces to a first-order transfer function, in which case we no longer expect overshoot behavior, as we did in the integral case.

## 8.5 Proportional-Integral-Derivative (PID) Control

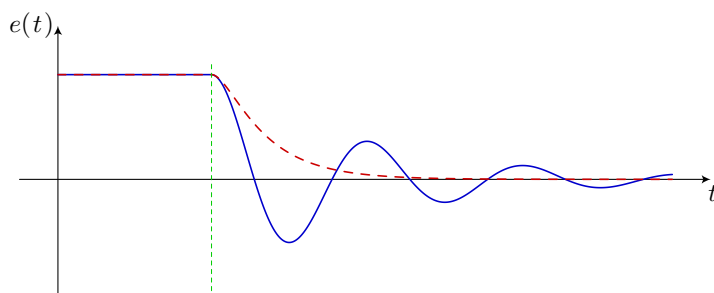
In **proportional-integral-derivative (PID) control**, the idea is to add a derivative term to PI control, so the transfer function is

$$\tilde{K}(\omega) = K_P + \frac{iK_I}{\omega\tau} - i\omega\tau K_D. \quad (8.30)$$

(PID-controller transfer function)

Here,  $K_D$  is the derivative gain. Intuitively, this can help in cases where overshoot and ringing is a problem. Qualitatively, consider the overshooting case below in blue, where the input is suddenly changed at the time marked by the dashed green line.





Qualitatively, the overshoot occurs because the slope of the error signal is too steep, due to the action of the controller. By putting in a term proportional to the derivative, the controller “senses” the steep slope corresponding to an impending overshoot, and reduces the control action. This can result in better settling, as in the dashed red line. In terms of the closed-loop transfer function, the effect of the derivative gain is to modify the damping coefficient of the feedback system, which can eliminate the ringing and promote better settling.

Setting the parameters for a PI or PID loop is something of an art. We won’t get into this here, but one reasonably simple method for setting the gain parameters is the **Ziegler-Nichols method**.<sup>3</sup>

---

<sup>3</sup>J. G. Ziegler and N. B. Nichols, “Optimum Settings for Automatic Controllers,” *Transactions of the ASME* **64**, 759 (1942), copies available at <http://chem.engr.utc.edu/Student-files/x2008-Fa/435-Blue/1942-paper.pdf> and <http://www.driedger.ca/Z-N/Z-n.pdf>. See also Allard Mosk, “Tutorial on Experimental Physics of Ultracold Gases,” in *Interactions in Ultracold Gases: From Atoms to Molecules*, Matthias Weidemüller and Claus Zimmerman, Eds. (Wiley-VCH, 2003), p. 215 (doi: 10.1002/3527603417.ch5).

## **Part II**

# **Digital Electronics**



## Chapter 9

# Binary Logic and Logic Gates

### 9.1 Binary Logic

The idea behind *binary logic* is to represent information using only two states. You can call these states **TRUE** and **FALSE**, or you can use the corresponding numerical values 1 and 0. We will explore in much more detail how to represent and use information in this form, but for now, note that we call the fundamental (abstract) element that carries these states a **bit**. That is, a single bit can have either the values 0 or the value 1.

The idea behind **digital logic** and **digital electronics** is to represent the binary states by two different electronic states, usually different voltages or voltage ranges, but sometimes different currents. For example, the standard for **transistor-transistor logic (TTL)** is to use nominal voltages of 0 V for **FALSE**, and +5 V for **TRUE**. We will get into more detailed specifics later.

Changing information into a digital representation has advantages and disadvantages. The main *disadvantage* of this approach is that it is necessary to **sample** analog signals (i.e., change continuous signals into discrete representations). The main *advantage* is in robustness to noise, as long as the noise amplitude is far below the physical separation between the logic states (e.g., TTL logic is robust to noise interference provided the noise is smaller than 5 V). Of course, for sophisticated logic systems (computers), often the advantages far outweigh the disadvantages.

We will treat binary logic as an abstract concept for now, and learn how to manipulate binary information. Then we will come back later to the physical implementation of binary logic.

### 9.2 Binary Arithmetic

In *binary arithmetic*—the binary analogue of the more usual arithmetic—the first thing to deal with is how to *represent* numbers in binary. To keep things relatively simple, we will stick to representing *integers* in binary (as opposed to rational approximations to real numbers, which are represented in either **fixed-point** or **floating-point** notation, the latter of which is more complicated).

#### 9.2.1 Unsigned Integers

The most basic form of a binary integer is an **unsigned integer**. Unsigned integers are just like decimal integers, but instead of counting from 0–9 and then carrying a 1 to the next place, you just count from 0–1 and then carry instead. (So the counting is 0, 1, 10, 11, 100, 101, 110, 111, ...) You can understand converting between binary and decimal best via an example. Suppose that we have the unsigned integer  $1011_2$  (the subscript “2” denotes binary, or base-2 arithmetic). There are four digits, which represent, from right to left, the “ones,” “twos,” “fours,” and “eights” places (just like the ones, tens, etc. in decimal counting). Then proceeding from the ones (rightmost) place, or the **least-significant bit (LSB)**, to the eights (leftmost)

place, or the **most-significant bit (MSB)**.

$$1011_2 = 1 \times 2^0 + 1 \times 2^1 + 0 \times 2^2 + 1 \times 2^3 = 1 + 2 + 8 = 11. \quad (9.1)$$

Note that with  $N$  digits, we can represent 2 values with each bit, for a total of  $2^N$  numbers (i.e., ranging from 0 to  $2^N - 1$ ).

### 9.2.1.1 Binary-Coded Decimal

Another representation of unsigned integers comes in **binary-coded decimal (BCD)**, where the idea is to convert each digit in a decimal number to binary, using 4 bits per decimal number. For example,  $11_{10} = 1011_2$ , but in BCD, this would be written 00010001. This representation is “wasteful” in that there are 16 states for each digit but only 10 decimal digits, but this representation is very convenient for implementations of digital numeric displays.

### 9.2.1.2 Hexadecimal

**Hexadecimal** arithmetic is just base-16 arithmetic. The 16 states are represented by 0–9 as usual, and the “extras” by A–F for the values 10–15. Since 4 bits, or a **nybble** (8 bits is a **byte**) has 16 states, a single hexadecimal digit is a convenient and compact representation for a binary nybble. Thus, for example,  $10111010_2 = \text{BA}_{16}$ .

## 9.2.2 Negative Values and Sign Conventions

Besides unsigned integers, it is useful to represent *negative* integers in binary. There are multiple conventions for this, however.

### 9.2.2.1 Sign-Magnitude Convention

The simplest convention, the **sign-magnitude convention**, is to tack on an extra bit (as a new MSB) to represent the sign, and the rest of the digits are just like an unsigned integer. For example, one nybble ranges from 0–15 as an unsigned integer, but as a signed value, it ranges from  $-7$  to  $+7$  as a signed integer (the three LSB’s range from 0–7, and the MSB gives the sign). Note that one value (1000) is “wasted” in this convention, because it is not different from (0000). The main advantage is the simplicity of the scheme. You can see the relatively serious disadvantage, however, by considering a couple of example numbers,

$$0001_2 = 1_{10}, \quad 1001_2 = -1_{10}. \quad (9.2)$$

Unfortunately, adding these two numbers gives  $1010_2 = -2_{10}$ , but really we’d like these to add to *zero*.

### 9.2.2.2 Two’s Complement

Preserving this additive-inverse property of negative numbers is the idea behind the **2’s-complement** representation: if  $n$  is a positive integer, just define the number  $(-n)$  such that it satisfies  $n + (-n) = 0$  in binary addition. For example, suppose we have

$$n = 0001_2 = 1_{10}. \quad (9.3)$$

Then  $-1_{10}$  in 2’s-complement notation is

$$-n = 1111_2 = -1_{10}. \quad (9.4)$$

To see this, first note that

$$n + (-n) = 10000_2, \quad (9.5)$$

but the important point is that we *drop* the MSB, because we regard addition in 4-bit binary as being modulo 16.

The advantage, of course, is that addition works as expected with positive and negative numbers in 2's-complement notation, and so does multiplication.

There are a couple of useful procedures for finding 2's-complement values (i.e., for finding the negative counterpart of a positive number):

1. Start by exchanging  $0 \longleftrightarrow 1$  on each digit, then add 1 to the result. For example, in the example above,

$$0001_2 \longrightarrow 1000_2 + 1 = 1111_2. \quad (9.6)$$

2. Note also that in  $N$ -bit arithmetic,  $2^{N-1} = -2^{N-1}$ . So we just need to figure out what number to add to  $-2^{N-1}$  to get the number we want. For example, in the above example in 4-bit arithmetic,  $2^3 = 1000_2 = -2^3 = -8$ . We want  $-1$ , which means we have to add 7 to  $-8$ . Since  $7 = 0111_2$ , we just say

$$-1 = -8 + 7 = 1000_2 + 0111_2 = 1111_2. \quad (9.7)$$

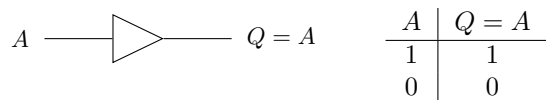
Alternately, we are just saying that for any negative number that we want in  $N$ -bit arithmetic, add  $2^N$  and then find the *unsigned* binary value. (In the example, add 16 to  $-1$  to get 15, or  $1111_2$ . In this convention, we're just taking the unsigned range of  $2^{N-1} + 1$  to  $2^N$ , and shifting the whole block to below zero.

## 9.3 Logic Gates

So far, we have discussed only binary-logic values and how to use them to represent numbers. But we also need to implement *transformations* on logic values, which are accomplished via **logic gates**, which is basically a logic-valued function of logic variables. We will talk about the simplest logic gates now, and just mention that more complicated gates can be represented in terms of the simpler ones.

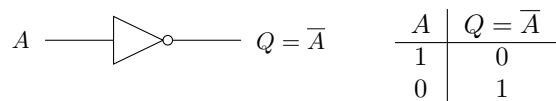
### 9.3.1 One-Input Gates

The simplest logic gates are the one-input gate, which takes only one logic value as input. A simple example is the **buffer gate**, which simply copies its input  $A$  to the output  $Q$ . The symbol for the buffer is below.

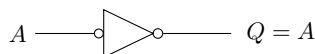


Above on the right is the **truth table**, a table enumerating all inputs and the corresponding output values.

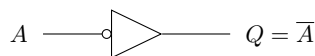
The other main one-input gate is the **inverter** or NOT gate, which changes the state of the input. The symbol and truth table are below.



Note the circle “o” in the diagram represents a NOT operation, which is denoted symbolically by a bar (that is, if  $A = 1$ , the  $\bar{A} = 0$ ). This same NOT operation may be applied to inputs as well. For example, this is a buffer gate,



and this is another NOT gate.

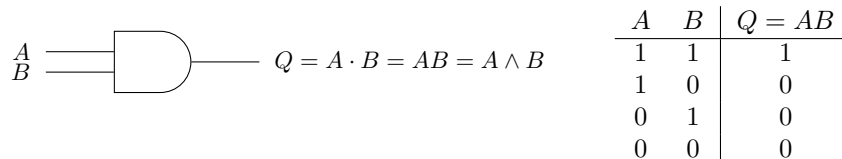


The buffer and NOT gates are the only one-input gates. The only other possibilities (in terms of truth-table content) have a fixed output for any input, which is usually not drawn as a gate with an input.

## 9.3.2 Two-Input Gates

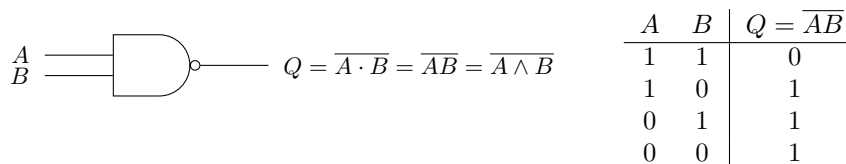
### 9.3.2.1 AND and NAND

Two-input gates are important, and easily available as electronic components. The first gate is the **AND** gate, whose symbol and truth table are shown below.



There are several notations for the **AND** operation in the diagram. Note that the output is only **TRUE** if *both* inputs are **TRUE**.

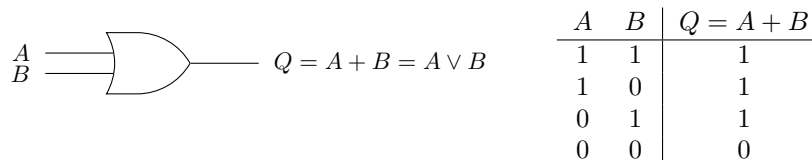
Adding a **NOT** to the output of the **AND** gate gives a **NAND** gate (i.e., **NOT AND**), which is just the negation of the **AND** gate.



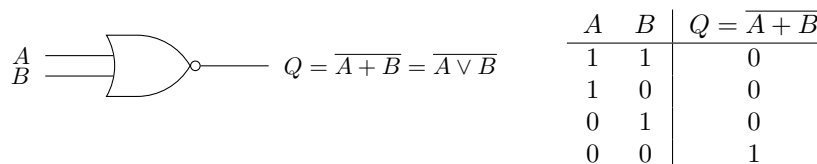
This gate is more important than it may seem at first glance, as we'll return to below.

### 9.3.2.2 OR and NOR

The next gate is the **OR** gate, whose operation is symbolically represented by “+.”



The output here is **TRUE** if either input is **TRUE** (or both inputs are **TRUE**). Of course, we can add a **NOT** to the output.



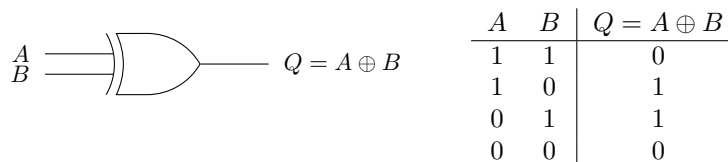
Thus, we obtain the **NOR** gate (**NOT OR**).

### 9.3.2.3 Universal Gates

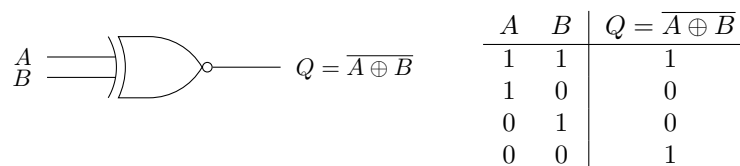
The **NAND** and **NOR** gates are special, because they are **universal gates**. That is, *any* logic operation can be realized by connecting a bunch of **NAND** gates, or by connecting a bunch of **NOR** gates.

### 9.3.2.4 XOR and XNOR

We'll briefly also mention the **XOR** gate (“exclusive **OR**”),



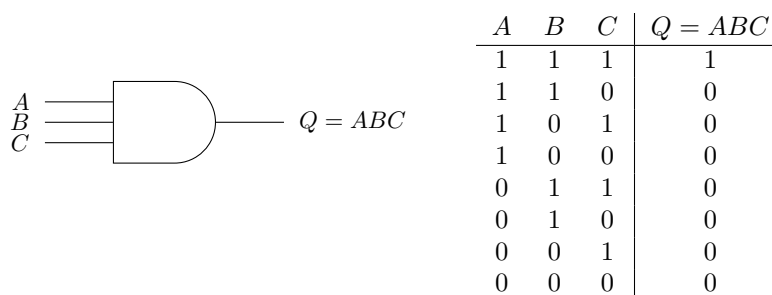
which is just like the **OR**, but the output is **FALSE** if *both* inputs are **TRUE**. The complement of the **XOR** gate is the **XNOR** gate (“exclusive **NOR**”), which is again like the **NOR** gate except for the case of two **TRUE** inputs.



In mathematical logic, the **XNOR** is the same as “if and only if.”

### 9.3.3 More Complex Gates

More complex gates are possible; for example, consider the 3-input **AND** gate below.

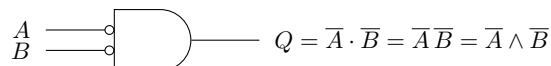


The idea is a reasonable obvious generalization of the 2-input **AND**: the output is **TRUE** only when *all* inputs are **TRUE**.

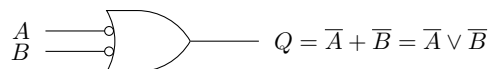
## 9.4 Circuit Practice

Here are a couple of gates with negated inputs.

- (a) Work out the truth table and find which 2-input gate that we introduced above is equivalent.



- (b) Do the same for this gate.



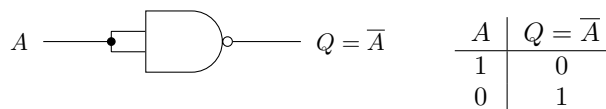
- (c) How do you make an inverter from a **NAND** gate?

**Solution.**

- (a) **NOR** gate.

- (b) **NAND** gate.

- (c) Tie the inputs together.





## 9.5 Exercises

### Problem 9.1

Convert:

- (a)  $89_{10}$  to 8-bit, unsigned binary
- (b)  $-89_{10}$  to 8-bit, signed (sign-magnitude) binary
- (c)  $-89_{10}$  to 8-bit, signed (2's complement) binary
- (d)  $89_{10}$  to hexadecimal
- (e)  $ABCD_{16}$  to decimal
- (f)  $0110011001100110_2$  to hex
- (g)  $0110011001100110_2$  to decimal
- (h)  $10011001_2$  (2's complement binary) to decimal
- (i)  $10011001_2$  (sign-magnitude binary) to decimal
- (j)  $11111111_2$  (unsigned) to decimal

### Problem 9.2

Convert:

- (a)  $75_{10}$  to 8-bit, unsigned binary
- (b)  $-75_{10}$  to 8-bit, signed (sign-magnitude) binary
- (c)  $-75_{10}$  to 8-bit, signed (2's complement) binary
- (d)  $75_{10}$  to hexadecimal
- (e)  $ABBA_{16}$  to decimal
- (f)  $1010101010101010_2$  to hex
- (g)  $1010101010101010_2$  to decimal
- (h)  $11011101_2$  (2's complement binary) to decimal
- (i)  $11011101_2$  (sign-magnitude binary) to decimal
- (j)  $11111101_2$  (unsigned) to decimal

### Problem 9.3

(a) Suppose  $x$  is a power of 2 (i.e.,  $x = 2^n$  for some positive integer  $n$ ,  $n \in \mathbb{Z}^+$ ). What does  $x$  “look” like when written out in (unsigned) binary? (That is, how can you recognize powers of two, when written in binary, just by looking at them?)

(b) In writing computer programs it is sometimes useful to check whether an integer is a power of 2. (One example is in computing numerical Fourier transforms, where the most common algorithms operate only on arrays whose lengths are powers of 2.)

A nice trick for checking if  $x$  is a power of 2 is to compute  $x \wedge (x - 1)$ . That is, subtract 1, and compute the bitwise AND with the original. How do you tell from the result if  $x$  is a power of 2? (You might try this on some examples to see the pattern.)

Note: “bitwise AND” means to compute the AND of corresponding binary digits. For example,  $1100_2 \wedge 1010_2 \equiv 1100_2 \cdot 1010_2 = 1000_2$ .

**Problem 9.4**

Short-question potpourri:

- (a) For any integer expressed in decimal, suppose you add the digits. The result is divisible by 3 if and only if the original number is also divisible by 3. Does this statement also hold for binary numbers? If so, explain why. If not, provide a counterexample.
- (b) What is  $1000\ 0000$  in decimal? Interpret the given number as an 8-bit, signed (2's complement) binary number.
- (c) What is  $111_{10}$  in unsigned binary?
- (d) What is  $111_{10}$  in hex?
- (e) Suppose  $1010_2 \oplus B = 1100_2$ , where the operation is bitwise. What is  $B$  in decimal?

**Problem 9.5**

Short-question potpourri:

- (a) What is  $2^{38\ 457} - 3$  in (unsigned) binary? (I suggest describing *how* to write down the binary expression, not actually writing it out.)
- (b) What is  $1010\ 1010_2 - 0101\ 0101_2$ ? (Do the calculation in binary.)
- (c) What is  $1101_2$  in hex?
- (d) What is  $1101_2 + 3_{16}$ ? Interpret both numbers as 4-bit, 2's complement binary numbers, and give your answer in 4-bit, 2's complement binary.
- (e) Suppose you generalize decimal-fraction notation (e.g.,  $0.9 = 9/10$ ) to binary fractions in the obvious way (e.g.,  $0.1_2 = 1/2$ ). What is the value of  $0.\overline{1}_2 \equiv 0.111\dots_2$  (i.e., the overbar here means a repeating digit) in decimal?

**Problem 9.6**

In each expression below, assume the operation to be a **bitwise logic operator**. That is, it works on corresponding bits in the binary representations of the numbers; for example, a bitwise **OR** on two binary numbers would look like:  $00001010_2 + 00000110_2 = 00001110_2$ .

Compute the result in each case, assuming 8-bit binary arithmetic.

- (a)  $25_{10} + 72_{10} = ?$  (Give a decimal result.)
- (b)  $32_{10} \cdot 223_{10} = ?$  (Give a decimal result.)
- (c)  $23_{10} \oplus -23_{10} = ?$  (Give a decimal result; assume signed binary, 2's complement.)
- (d)  $\overline{57}_{10} \oplus \overline{-123}_{10} = ?$  (Give a decimal result; assume signed binary, 2's complement.)

**Problem 9.7**

The binary **bit-shift operators**  $\ll$  and  $\gg$  are defined as follows:

- $n \ll k$  (" $n$  left-shifted by  $k$  places") is the binary representation of  $n$ , with all bits shifted to the left by  $k$  places (with 0's filling in any vacated spots). Thus  $0011_2 \ll 2 = 1100_2$ .
- $n \gg k$  (" $n$  right-shifted by  $k$  places") is the binary representation of  $n$ , with all bits shifted to the **right** by  $k$  places (removing any digits that "fall off the end"). Thus  $0110_2 \gg 2 = 0001_2$ .

For parts (a)–(c), compute the result in each case, assuming binary arithmetic with arbitrarily many digits allowed (i.e., no fixed number of bits like we will have in the last two parts). *Don't* write out all the digits of the answer, but give a decimal expression that *doesn't* involve a bit shift operator.

- (a)  $2^{207} \ll 5 = ?$

(b)  $(2^{4051} - 1) \gg 5 = ?$

(c)  $(2^{303} + 1) \ll 2 = ?$

For parts (d)–(e), assume fixed, 64-bit (unsigned) binary arithmetic (i.e., a fixed 64 binary “digits,” where adding 2 to the 64-bit binary representation of  $2^{63}$  is  $0_2$ ).

(d) What is the smallest integer  $m$  for which  $288 \ll m = 0$ ?

(e) What is the smallest integer  $n$  for which  $288 \gg n = 0$ ?

# Chapter 10

## Boolean Algebra

### 10.1 Algebras and Boolean Algebra

Intuitively, a **Boolean algebra** is an abstract, compact notation for logic (in which we will see that  $1+1 = 1$ ). A Boolean algebra is defined on the set  $\{0, 1\}$  (these are the “values” for Boolean variables), and has two **binary operations** “+” and “ $\cdot$ ” defined, though not the usual addition and multiplication. The + operation is defined by the truth table for the OR gate, while the  $\cdot$  operation is defined by the truth table for the AND gate. Recall that the truth tables for the AND and OR operations on Boolean variables  $A$  and  $B$  are invariant under the exchange of  $A$  and  $B$ , so both + and  $\cdot$  are **commutative** (i.e., the order of the variables don’t matter):

$$A + B = B + A, \quad AB = BA. \quad (10.1)$$

(We’re not bothering to write the “ $\cdot$ ” explicitly here.) These operations are also **associative**, which means that the order of two successive operations does not matter:

$$A + (B + C) = (A + B) + C, \quad A(BC) = (AB)C. \quad (10.2)$$

The other usual algebraic property that holds here is the **distributive** property:

$$A(B + C) = (AB) + (AC). \quad (10.3)$$

Finally, a number of simple identities hold for the Boolean binary operators:

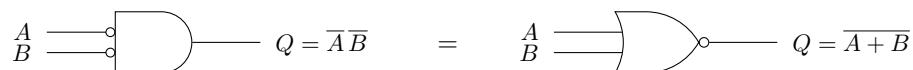
1.  $\overline{\overline{A}} = A$
2.  $A \cdot 0 = 0$
3.  $A \cdot 1 = A$
4.  $A \cdot A = A$
5.  $A \cdot \overline{A} = 0$
6.  $A + 0 = A$
7.  $A + 1 = 1$
8.  $A + A = A$
9.  $A + \overline{A} = 1$

Roughly speaking, the NOT operation (bar) is something like a minus sign, in which case some of these identities seem familiar, but some seem less so (like  $A + 1 = 1$ ).

## 10.2 Boolean-Algebraic Theorems and Manipulations

### 10.2.1 De Morgan's Theorems

Recall the circuit practice from Section 9.4, where we examined **AND** and **OR** gates, where both inputs are negated. For example, the negated-input **AND** gate is equivalent to **NOR** gate, as shown schematically below.



In algebraic notation, this is

$$\overline{A} \cdot \overline{B} = \overline{A + B}. \quad (10.4)$$

(De Morgan theorem)

Similarly

$$\overline{A + B} = \overline{A} \cdot \overline{B}. \quad (10.5)$$

(De Morgan theorem)

These are extremely useful in transforming negated expressions, as we will see.

### 10.2.2 Absorption Theorems

Two other useful theorems are called **absorption theorems**:

$$\begin{aligned} A + (A \cdot B) &= A \\ A \cdot (A + B) &= A. \end{aligned} \quad (10.6)$$

(absorption theorems)

We will leave the proofs as exercises (in circuit practice).

### 10.2.3 Another Theorem

Here is another theorem that is often useful:

$$A + \overline{A}B = A + B. \quad (10.7)$$

(negated AND theorem)

We will again leave the proof as an exercise, but essentially this is saying that because of the **OR** with  $A$ , the  $\overline{A}$  never really matters.

### 10.2.4 Example: XOR Gate

As an example of Boolean algebra and implementation of algebraic expressions in gates, consider the **XOR** operation, where we would like to show that

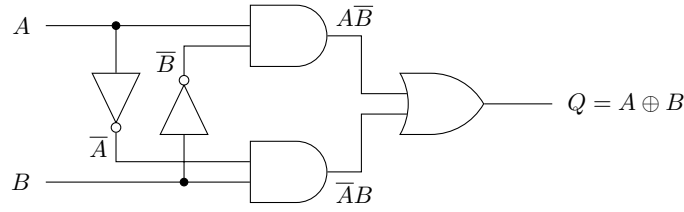
$$A \oplus B = \overline{A}B + A\overline{B}. \quad (10.8)$$

(XOR expression)

We can first do this by working through the truth table for the right-hand side, and verifying that it matches the truth-table results for  $A \oplus B$ .

$A$	$B$	$\overline{A}B$	$A\overline{B}$	$A \oplus B$
1	1	0	0	0
1	0	0	1	1
0	1	1	0	1
0	0	0	0	0

Now using this expression, we can show how to implement an **XOR** gate, in terms of regular gates.



To trace through this, the negations  $\bar{A}$  and  $\bar{B}$  are realized with NOT gates, and finally two AND gates and an OR gate to generate the correct combination.

#### 10.2.4.1 NAND-Gate Realization

As we alluded to before, NAND and NOT gates are universal, and can be used to realize any gate. So how can we realize an XOR gate out of only NAND gates, for example? Let's do some algebraic transformations to see how to do this. First, starting with the expression (10.8),

$$A \oplus B = \bar{A}B + A\bar{B}, \quad (10.9)$$

we can add in  $A\bar{A} = 0$  and  $B\bar{B} = 0$  to obtain

$$A \oplus B = B(\bar{A} + \bar{B}) + A(\bar{A} + \bar{B}). \quad (10.10)$$

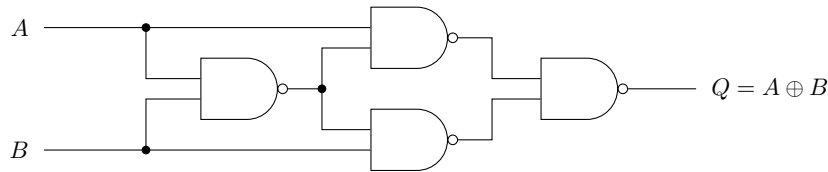
Then using the second De Morgan theorem (10.5),  $\bar{A} + \bar{B} = \overline{A \cdot B}$ ,

$$A \oplus B = B(\overline{AB}) + A(\overline{AB}). \quad (10.11)$$

Using the same theorem once more,

$$A \oplus B = \overline{[B(\overline{AB})][A(\overline{AB})]}. \quad (10.12)$$

Now notice that every operation here is a NAND, and we need one operation to generate  $\overline{AB}$ , two more to combine it with  $A$  and  $B$ , and one more for the final combination. The circuit realizing this expression is shown below.



Note that we can obtain a simpler but less-efficient expression by applying the De Morgan theorem only once as follows:

$$A \oplus B = \bar{A}B + A\bar{B} = \overline{(\bar{A}B)(A\bar{B})}. \quad (10.13)$$

This is less efficient in terms of NAND gates: two NANDs are needed to make  $\bar{A}$  and  $\bar{B}$ , two more to make the combinations  $\bar{A}B$  and  $A\bar{B}$ , and one more to make the final combination (a total of 5).

#### 10.2.5 Example: Algebraic Simplification

As another example of simplifying an expression, consider the three-variable expression  $(A + B) \cdot (A + C)$ . Starting out, we can distribute twice,

$$\begin{aligned} (A + B) \cdot (A + C) &= A \cdot (A + C) + B \cdot (A + C) \\ &= AA + AC + AB + BC. \end{aligned} \quad (10.14)$$

Then  $AA = A$ , and using the first absorption theorem in Eqs. (10.6) to write  $A + AC = A$ ,

$$\begin{aligned}(A + B) \cdot (A + C) &= (A + AC) + AB + BC \\ &= A + AB + BC.\end{aligned}\tag{10.15}$$

Then also  $A + AB = A$ , so

$$(A + B) \cdot (A + C) = A + BC,\tag{10.16}$$

a somewhat simpler expression (two operations vs. three).

## 10.3 Karnaugh Maps

It can often be difficult to see how to realize a particular logic function in terms of logic gates, just via algebraic manipulations. One tool that makes this more intuitive, at least for small numbers of inputs (3 or 4), is the **Karnaugh map**. (The cases of 1 and 2 inputs we've already mostly covered with standard gates, and it's easy to do these exhaustively.) These days, computers can exhaustively search logical expressions to find the simplest ones, so you might ask, why bother with techniques like the Karnaugh map for doing algebra? Well, algebraic simplicity is a subjective thing (depending somewhat on the goals of the particular problem at hand), and symbolic manipulation systems on computers have notoriously poor "taste" when it comes to "simple." The Karnaugh map is also such an elegant and downright *cool* tool that it's really worth studying.

The first idea behind a Karnaugh map is to make a 2D table of logic inputs and truth-table values. The second twist is to order the inputs using a 2-bit **Gray code**, which means that we count as 00, 01, 11, 10, instead of the usual binary-counting order. The point is when counting this way we change only 1 bit at a time (in regular binary, this doesn't happen when we count from 01 to 10). The motivation for Gray codes comes from mechanical implementations of logic, where you may get spurious transitional states if bits don't change synchronously (this happens in some fast logic circuits as well). That is, when counting from 01 to 10, the actual sequence may be 01 to 00 to 10 if the LSB changes before the MSB. In terms of the Karnaugh map, the idea is to keep "related" input states grouped together.

The process of hunting for simplifications in a Karnaugh map is hard to explain, but it's easy to get the idea by studying a few examples.

### 10.3.1 Three-Input Example

Before, as a Boolean-algebraic example, we showed in Eq. (10.16) that

$$(A + B) \cdot (A + C) = A + BC.\tag{10.17}$$

We will show how to obtain this and other transformations via the Karnaugh map. The first task is to write out the diagram as a table. Notice that the four values of  $AB$  are along the top, in Gray-coded order, and the two  $C$  values are along the side. We are also writing out each output value for each set of possible input values, so this is just a truth table in 2D form, here for  $(A + B) \cdot (A + C)$ .

	$AB$	00	01	11	10
$C$					
0		0	0	1	1
1		0	1	1	1

Now to analyze this, the idea is to look for blocks of 1's in square or rectangular shapes ( $2 \times 1$ ,  $2 \times 2$ , etc.). One example is below.

		AB			
		00	01	11	10
C	0	0	0	1	1
	1	0	1	1	1

$BC$        $A$

$A + BC$

The  $2 \times 2$  block here corresponds to every input where  $A = 1$ , hence we have labeled it “ $A$ .” Similarly, the  $2 \times 1$  block corresponds to inputs where  $B = 1$  and  $C = 1$  (hence  $BC = 1$ ), so we label it “ $BC$ .” The entire group of 1’s is the union of these two, so the logical expression is the OR of these two blocks, hence an equivalent expression is  $A + BC$ .

Generally speaking, bigger blocks correspond to simpler expressions, so the best simplifications occur by covering the 1’s with large blocks. Also, usually it is best to look for blocks with dimensions of 2, 4, 8, etc. As an example, note that we could have done the last covering without any overlap if we kept the  $4 \times 4$  block and then introduced a  $1 \times 1$  block, as below.

		AB			
		00	01	11	10
C	0	0	0	1	1
	1	0	1	1	1

$\overline{A}BC$        $A$

$A + \overline{A}BC$

The small block corresponds to  $\overline{A}BC$ , since we have to restrict all three variables, and this leads to the more complicated (but equivalent) expression  $A + \overline{A}BC$ .

We can also look at some other attempts to simplify with a Karnaugh map that will yield less compact results, just to illustrate the technique. For example, we can “overcover” the 1’s by using two  $4 \times 4$  blocks,  $A$  and  $B$ , and then combine them. However, we must exclude one location that has a zero; the location is  $\overline{A}B\overline{C}$ .

		AB			
		00	01	11	10
C	0	0	0	1	1
	1	0	1	1	1

$B$        $A$

$A + B\overline{A}\overline{C}$

Thus, to combine these, we *negate* the null block and AND the result with the  $B$  block to obtain  $B\overline{A}\overline{C}$ . We then OR this with  $A$  to obtain  $A + B\overline{A}\overline{C}$ .

Another possibility is to use a similar technique, but focusing on the 0’s. For example, we can identify a block of mostly zeros,  $\overline{A}$ . However, we must exclude the 1, which is located at  $\overline{A}B\overline{C}$ .

		AB			
		00	01	11	10
C	0	0	0	1	1
	1	0	1	1	1

$\overline{A}$        $\overline{A}B\overline{C}$

$\overline{\overline{\overline{A}B\overline{C}}}$



So to find the region of 0's, we have to negate the 1 block and **AND** this with the 0 block, to obtain  $\overline{A}\overline{A}BC$ . Then to get the block of 1's we must negate the overall result, to obtain  $\overline{\overline{A}\overline{A}BC}$ . Note that by De Morgan's theorems this is equivalent to  $A + \overline{A}BC$ , as we showed two diagrams ago.

### 10.3.2 Four-Input Example

In addition to three-input problems, it is not much harder to extend the analysis to four-input problems. For example, consider the following truth table.

		$AB$				
		00	01	11	10	
$CD$	00	1	1	1	1	$\overline{B}\overline{C}$
	01	1	0	0	1	
	11	0	0	0	0	
	10	0	1	1	0	
		$B\overline{D}$				

$\overline{B}\overline{C} + B\overline{D}$

One thing to notice is that we have extended the vertical direction to cover the two variables  $C$  and  $D$  together, and the other is that we have put in  $4 \times 4$  blocks that “wrap” around from top to bottom or right to left. That is, the Karnaugh map has periodic boundary conditions for the purposes of finding blocks.

### 10.3.3 XOR Example

In searching for blocks, it is somewhat harder to see XOR and XNOR blocks, but it is possible. An example is below.

		$AB$				
		00	01	11	10	
$C$	0	0	1	0	1	$\overline{C}$
	1	0	0	0	0	

$A \oplus B$

$(A \oplus B)\overline{C}$

Due to the ordering of the horizontal axis, the  $A \oplus B$  block is split, but we can combine it with the  $\overline{C}$  via an **AND** (to intersect the blocks) to obtain a relatively simple expression.

Note that some flexibility is usually beneficial when using a Karnaugh map: it is not necessarily a good tool for finding solutions in terms of a *particular* gate (e.g., all **NAND** gates).

### 10.3.4 Race Hazards

A **race condition** is a spurious output of a circuit if the inputs don't change state simultaneously (i.e., a “glitch”). This can be a big problem if this output is the input to a latch or a memory circuit that will “trigger” on the glitch.

Intuitively, in a Karnaugh map, a glitch is possible if the changing inputs cross between disjoint blocks of 1's, because the output state is being controlled by transitions of two gates feeding into the same final gate. For example, returning to the  $(A + B)(A + C)$  example, suppose we make a transition between 111 to 011 in  $ABC$ . In this logic realization,

		AB			
		00	01	11	10
C	0	0	0	1	1
	1	0	1	1	1

$BC$        $A$

$A + BC$

we stay inside the  $BC$  block, so we don't expect any glitches: the output stays at 1 during the transition. However, in *this* realization,

		AB			
		00	01	11	10
C	0	0	0	1	1
	1	0	1	1	1

$\overline{A}BC$        $A$

$A + \overline{A}BC$

we must cross between blocks, so a glitch is possible. Specifically, when  $A$  goes from  $1 \rightarrow 0$ , a slight delay in  $\overline{A}$  going from  $0 \rightarrow 1$  results in the output going momentarily to 0 during the input transition, even though it should remain as 1.

As another example, let's return to the four-input example.

		AB			
		00	01	11	10
CD	00	1	1	1	1
	01	1	0	0	1
	11	0	0	0	0
	10	0	1	1	0

$B\overline{D}$

$\overline{B}\overline{C}$

$\overline{B}\overline{C} + B\overline{D}$

There is a similar problem here when  $ABCD$  goes from  $1100 \rightarrow 1000$ , because we cross in between blocks. However, by adding another block, we can "protect" the circuit from glitches in this transition. Here, we add  $\overline{C}\overline{D}$ , and combine it with an OR operation.

		AB			
		00	01	11	10
CD	00	1	1	1	1
	01	1	0	0	1
	11	0	0	0	0
	10	0	1	1	0

$B\overline{D}$

$\overline{B}\overline{C}$

$\overline{C}\overline{D}$

$\overline{B}\overline{C} + B\overline{D} + \overline{C}\overline{D}$

Of course, the price for robustness is a more complicated expression.

## 10.4 Circuit Practice

### 10.4.1 Boolean-Algebra Theorems

Here, you should prove two things that we only introduced earlier.

(a) Prove the first absorption theorem in Eqs. (10.6):

$$A + (A \cdot B) = A. \quad (10.18)$$

Use a truth table or algebra.

(b) Prove the second absorption theorem in Eqs. (10.6):

$$A \cdot (A + B) = A. \quad (10.19)$$

**Solution.**

(a) First suppose  $B = 0$ . Then

$$A \cdot B = A \cdot 0 = 0, \quad (10.20)$$

and so

$$A + (A \cdot B) = A + 0 = A. \quad (10.21)$$

Now take the other case, where  $B = 1$ . Then

$$A \cdot B = A \cdot 1 = A, \quad (10.22)$$

and so

$$A + (A \cdot B) = A + A = A. \quad (10.23)$$

(b) Using the same method, first suppose  $B = 0$ . Then

$$A + B = A + 0 = A, \quad (10.24)$$

and so

$$A \cdot (A + B) = A \cdot A = A. \quad (10.25)$$

Now taking the other case  $B = 1$ ,

$$A + B = A + 1 = 1, \quad (10.26)$$

and so

$$A \cdot (A + B) = A \cdot 1 = A. \quad (10.27)$$

### 10.4.2 Karnaugh Map

Write down the Karnaugh map and a logic circuit for the following function: the output is 1 if and only if the input, a 3-bit unsigned integer, is prime. (Don't count 0, 1, or 2 as prime integers.)

**Solution.** The primes are 3, 5, and 7. In binary, these are 011, 101, and 111. Hence the Karnaugh map:

		AB			
		00	01	11	10
C	0	0	0	0	0
	1	0	1	1	1

(A + B)C

The simplest solution for a logic gate is to make a  $3 \times 1$  block, noting that the  $AB$  part is specified by  $(A + B)$ .

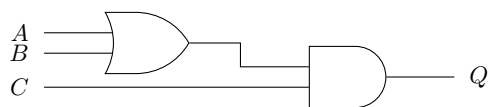
$C$	$AB$			
	00	01	11	10
0	0	0	0	0
1	0	1	1	1

$\overline{A}\overline{B}$

$\overline{\overline{A}\overline{B}C} = (A + B)C$

An alternate, and equivalent solution (via De Morgan's theorems) is shown below, by starting with the  $4 \times 1$  block  $C$ , and then excluding the 0.

The circuit to realize this function is shown below.



## 10.5 Exercises

### Problem 10.1

Simplify the expression

$$Q = \overline{A} \overline{B} \overline{C} + \overline{A} \overline{B} C + \overline{A} B \overline{C} + A \overline{B} \overline{C}, \quad (10.28)$$

and draw a logic circuit that realizes it. (This can be done with only 3 2-input gates and 1 3-input gate; try to at least reduce this somewhat, and it's best if your solution reflects the symmetry of the original expression. Also, try to use algebraic transformations rather than writing out truth tables.)

### Problem 10.2

Simplify the expression

$$Q = \overline{A} \overline{B} \overline{C} + \overline{A} \overline{B} C + \overline{A} B \overline{C} + A \overline{B} \overline{C}, \quad (10.29)$$

and draw a logic circuit that realizes it. (This is possible with only 1 3-input gate; try to at least reduce this somewhat, and it's best if your solution reflects the symmetry of the original expression. Also, try to use algebraic transformations rather than writing out truth tables.)

### Problem 10.3

(a) Simplify the following Boolean expression:

$$Q = (A + B)(\overline{B} + A) + \overline{\overline{A} \overline{B} + A + \overline{B}} + B. \quad (10.30)$$

(b) Sketch a realization of this expression (after simplifying!) using only 2-input NAND gates.

### Problem 10.4

(a) Simplify the following Boolean expression:

$$Q = \overline{(\overline{A} + \overline{B})(\overline{A} + \overline{B})} + (\overline{A} + \overline{B})(\overline{A} + \overline{B}). \quad (10.31)$$

(b) Sketch a realization of this expression (after simplifying!) in terms of only XNOR gates.

### Problem 10.5

(a) Simplify the following Boolean expression:

$$Q = \overline{[(\overline{A} \overline{B} \overline{C})(\overline{B} \overline{A} \overline{C})][(\overline{A} \overline{B})(\overline{A} \overline{B} \overline{C})]}. \quad (10.32)$$

(b) Sketch a realization of this expression (after simplifying!) in terms of only NOR gates.

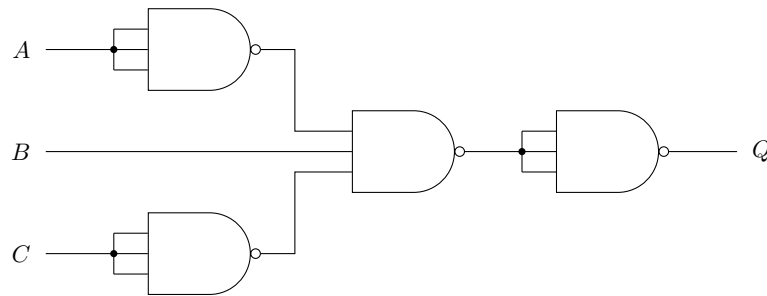
### Problem 10.6

Simplify the following Boolean expression:

$$\begin{aligned} Q = & A + B + C + D + E \\ & + \overline{A} \overline{B} + \overline{B} \overline{C} + \overline{C} \overline{D} + \overline{D} \overline{E} + \overline{E} \overline{A} \\ & + \overline{A} B C + \overline{B} C D + \overline{C} D E + \overline{D} E A + \overline{E} A B \\ & + \overline{A} B C D + \overline{B} C D E + \overline{C} D E A + \overline{D} E A B \\ & + \overline{E} A B C. \end{aligned} \quad (10.33)$$

**Problem 10.7**

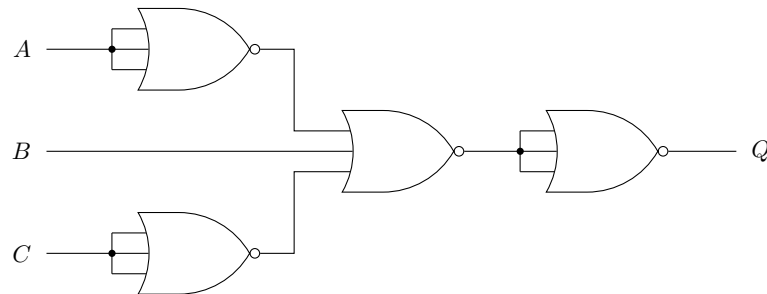
Consider the following circuit, based on 3-input NAND gates.



- Write down the logic (Boolean) expression for the circuit.
- Write down the truth table.

**Problem 10.8**

Consider the following circuit, based on 3-input NOR gates.



- Write down the logic (Boolean) expression for the circuit.
- Write down the truth table.

**Problem 10.9**

Show how you can realize an XOR gate ( $A \oplus B = \overline{A} \cdot B + A \cdot \overline{B}$ ) using only NOR gates.

**Problem 10.10**

Show how you can realize an XNOR gate ( $\overline{A \oplus B} = \overline{\overline{A} \cdot B + A \cdot \overline{B}}$ ) using only NAND gates.

**Problem 10.11**

One possible generalization of boolean logic is to three states: the usual TRUE (1) and FALSE (0), plus a third UNKNOWN (?) state. The (still associative and distributive) operations  $+$  and  $\cdot$  can then be defined by the following:

- if  $A$  and  $B$  are both Boolean (i.e., either 0 or 1), the result is just the Boolean result.
- $A + 1 = 1$  for any  $A$  (including ?)
- $A \cdot 0 = 0$  for any  $A$  (including ?)
- $? + 0 = ?$

- $? + ? = ?$
- $? \cdot 1 = ?$
- $? \cdot ? = ?$
- $\overline{A}$  is the same as in Boolean algebra if  $A$  is 0 or 1, and ? otherwise.

(Think through how these rules sensibly generalize Boolean algebra, where UNKNOWN means “could be either TRUE or FALSE.”)

Prove that both De Morgan theorems still hold in this generalized logic algebra. You can do this by truth table, but try to be more clever, breaking the proof down into cases you can work out algebraically.

### Problem 10.12

- Show that the XOR operation is both commutative and associative. (You may assume  $+$  and  $\cdot$  are commutative and associative.)
- Simplify the following Boolean expression. (Eliminate any constant values, 0 or 1, from your final expression.)

$$Q = \left[ \left[ \left[ (A \oplus B) \oplus (\overline{A} \oplus B) \oplus (\overline{A} \oplus \overline{C}) \right] \oplus B \right] \oplus C \right]. \quad (10.34)$$

### Problem 10.13

Write out the Karnaugh map for a circuit where the output is true if the input (a 3-bit, unsigned integer, 0–7) is in the Fibonacci sequence. Give a circuit implementation in terms of 2-input gates.

### Problem 10.14

Write out the Karnaugh map for a circuit where the output is true if the input (a 3-bit, unsigned integer, 0–7) is one of the first 6 digits of  $\pi$ . Give a circuit implementation in terms of 2-input gates.

### Problem 10.15

A **semiprime number** is a positive integer that is the product of two prime numbers. The prime numbers need not be distinct, 1 doesn’t count as one of the primes. For example, 0–3 are not semiprime, but 4 is.

- Write down the Karnaugh map for the function of the boolean variables  $A$ ,  $B$ ,  $C$ , and  $D$ , which is true when the concatenation  $ABCD$  (when converted to decimal as an unsigned integer) is semiprime.
- Find a (reasonably simple) boolean expression for this 4-bit semiprime function you diagrammed in (a).
- Sketch a logic implementation of this function in terms of logic gates.

### Problem 10.16

- Write down the Karnaugh map for the function of the Boolean variables  $A$ ,  $B$ ,  $C$ , and  $D$ , which is true when the concatenation  $ABCD$  (when converted to decimal as an unsigned integer) is greater than or equal to 6.
- Find a (reasonably simple) boolean expression for the logic function you diagrammed in (a).
- Sketch a logic implementation of this function in terms of only 2-input NAND gates.

**Problem 10.17**

Find logic to perform multiplication of two 2-bit (unsigned) integers (i.e., 0–3), with a 4-bit output.

*Hint:* use a separate Karnaugh map for each output bit.<sup>1</sup>

**Problem 10.18**

Find logic to perform addition of two 2-bit (unsigned) integers (i.e., 0–3), with a 3-bit output.

*Hint:* use a separate Karnaugh map for each output bit.

**Problem 10.19**

- (a) Write down the Karnaugh map for the function of the Boolean variables  $A$ ,  $B$ ,  $C$ , and  $D$ , which is true when  $ABCD$  (when converted to decimal as an unsigned integer,  $D$  being the LSB) is odd *and* greater than 4.
- (b) Find a (reasonably simple) boolean expression for the logic function you diagrammed in (a).
- (c) Sketch a logic implementation of this function in terms of only 2-input NAND gates.

**Problem 10.20**

- (a) Write down the Karnaugh map for the function of the Boolean variables  $A$ ,  $B$ ,  $C$ , and  $D$ , which is true when  $ABCD$  (when interpreted as an unsigned integer,  $D$  being the LSB) is valid as a BCD integer.
- (b) Find a (reasonably simple) boolean expression for the logic function you diagrammed in (a).
- (c) Sketch a logic implementation of this function in terms of only 2-input NOR gates.

**Problem 10.21**

Consider the Boolean expression

$$Q = (\overline{A} + B) + AC + \overline{C}(A + B). \quad (10.35)$$

- (a) Write down the Karnaugh map for this expression, and find a simpler (equivalent) expression. Indicate the blocks corresponding to the three terms in the logical expression above.
- (b) This expression has a race hazard; find it and explain the problem (i.e., give the initial and final states of  $ABC$ , and explain what needs to happen to produce a spurious output state).

<sup>1</sup>Paul Horowitz and Winfield Hill, *The Art of Electronics*, 2nd ed. (Cambridge, 1989), Exercise 8.14 (ISBN: 0521370957).





# Chapter 11


## Physical Implementation of Logic Gates

So far, we have studied logic and logic gates, but logic is much more useful if we can implement logic gates *physically*. Generally speaking you can buy these as prepackaged integrated circuits, but it is still useful to understand how to implement these, for (1) extra intuition and (2) to understand the limits and quirks of commonly available electronic logic gates. We will start with simple examples of logic realizations and progress to realistic (but more complicated) cases.


The material here in this chapter relies on previous material on diodes from Chapter 3 and transistors from Chapter 4. However, we will *briefly* review some of the relevant material here.

### 11.1 Simple Mechanical Switches

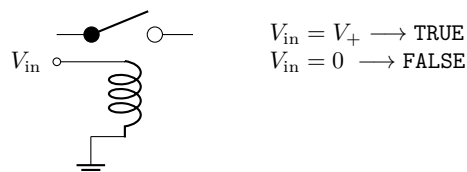
Fundamentally, electronic logic gates work by involving switches. Typically these are some form of *electronic* switches, but of course these can be ordinary mechanical switches (equivalent to connecting two points by a wire or breaking the wire connection). One simple convention for a **single-pole, single-throw (SPST)** switch (“single pole” = single circuit to break or connect; “single throw” = single possible connection to make or break), as shown below, is that the closed (conducting) or ON state is **TRUE**,

closed = TRUE  


and the open or OFF state is **FALSE**.

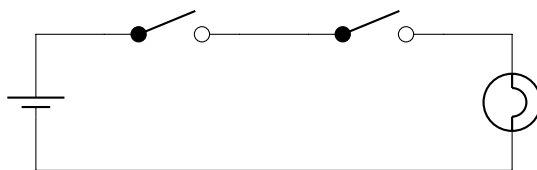
open = FALSE  


This convention agrees with a common convention for logic in terms of voltage levels, where **HIGH** voltage is **TRUE** and **LOW** voltage is **FALSE**, if we consider a **relay** (magnetically controlled switch), as shown below.

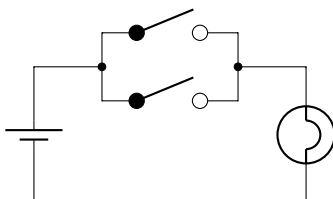


The relay is pulled closed when the voltage is **HIGH** (at some voltage  $V_+$ ), due to the magnetic field of the coil; when the voltage is zero, there is no field and the switch pops open (due to the action of a spring).

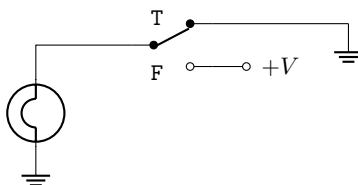
Using switches it is easy to see how to construct an **AND** gate, if two switches are in series, since both switches must close to light the light bulb (the logical “output” here).



For an **OR** gate, the two switches are in parallel, so only one switch needs to be closed to light the bulb.



As a final example, an inverter is shown below.



The switch here is a **single-pole, double-throw (SPDT)** switch (“double throw” = two alternative contacts for the switch), where “up” on the switch is **TRUE** and “down” on the switch is **FALSE**.

## 11.2 Diode Logic (DL)

The simplest “purely electronic” examples of logic come in the form of **diode logic (DL)**. Before examining some DL gates, first let’s review how diodes work.

### 11.2.1 Diode Review

A diode is a two-terminal device, as shown below, and it acts as a one-way valve for current: current can only flow from the anode to the cathode (in the direction of the “diode arrow” in the schematic symbol).



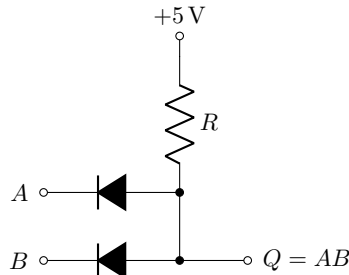
That is, as in the diagram below, if the anode voltage  $V_A$  is greater than the cathode voltage  $V_B$ , then current flows; otherwise, no current flows.

$$A \circ \begin{array}{c} I \rightarrow \\ \blacktriangleright \end{array} \circ B = \begin{array}{c} I \rightarrow \\ \text{if } V_A > V_B \end{array} A \circ \circ B \quad \text{or} \quad \begin{array}{c} (I = 0) \\ \text{if } V_A < V_B \end{array} A \circ \circ B$$

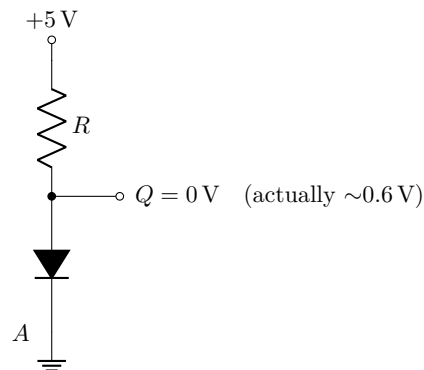
You can think of the diode as being a short circuit in the first case (“forward-biased”), and an open circuit in the second (“reverse-biased”). However, the real situation is a bit more complicated: a slightly better model is that there is a forward voltage drop of around 0.6 to 0.7 V when the diode is conducting current (at least for a small-signal silicon diode).

### 11.2.2 DL AND Gate

Now to see how to realize gates in DL. Below is a realization of an AND gate. The DL convention here is that 0 V is **FALSE**, and +5 V is **TRUE**.



If both inputs are at +5 V then all points are at the same voltage, including the output  $Q$ . If one input is low, say  $A$ , then the situation is as shown below.

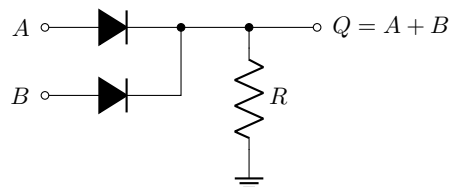


The diode is forward-biased, and thus shorts the output  $Q$  to ground. The power-supply voltage (+5 V) is dropped across the resistor because the diode causes sufficient current to flow through the resistor to ground to ensure this. The state of the other input ( $B$ ) is irrelevant here, because either it “agrees” with  $A$ , or if it is **HIGH**,  $B$ ’s diode is reverse-biased, so it is disconnected from the circuit.

Actually, the output voltage is not quite 0 V in the latter case; because the diode has a forward-voltage drop, the output **FALSE** state is more like 0.6 V.

### 11.2.3 DL OR Gate

Another DL gate, the OR gate, is shown below.



Here, if either input is at +5 V, then the corresponding diode is forward-biased, so output is at +5 V [actually, (+5 – 0.6) V if we account for the diode’s voltage drop]. If both inputs are at 0 V, then the whole circuit, including the output, is also at 0 V.

The main problem in the DL circuits is that one of the signal-voltage states “degrades” by 0.6 V on each gate, so not many gates can be cascaded while keeping the signal levels distinguishable. This motivates the use of *active* devices in logic circuits that can maintain the proper voltage levels.

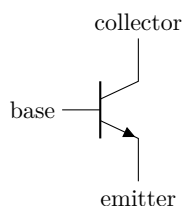
## 11.3 Resistor-Transistor Logic (RTL)

A step up in terms of sophistication is **resistor-transistor logic (RTL)**, which is obsolete but relatively easy to understand. However, bear in mind these circuits are still very much practical. If you need an inverting logic buffer in a circuit, for example, it may well be easier to reach for a couple of resistors and a transistor (and fit them onto your circuit board!) rather than a whole 14-pin chip.

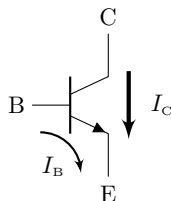
Again, we first have to review how a transistor—specifically, the NPN bipolar junction transistor (BJT)—behaves, in particular as a switch.

### 11.3.1 BJT Review

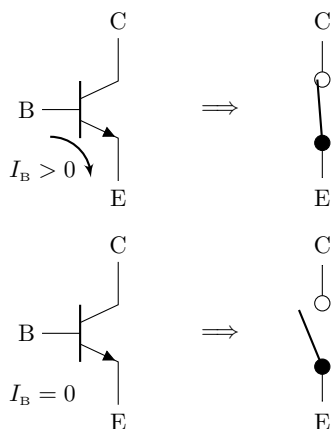
Recall that a BJT (bipolar junction transistor) is a three-terminal device, with terminals labeled as in the diagram below.



The transistor acts as a switch for current, based on another current. We will consider two currents,  $I_B$  from the base to the emitter, and  $I_C$  from the collector to the emitter, as shown below.



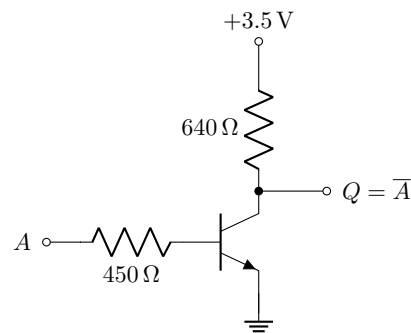
Then  $I_B$  acts as the control current, and  $I_C$  is the current to be switched. Simplistically, if there is some current  $I_B$ , then  $I_C$  can flow, so the C–E path acts as a closed switch.



However, if  $I_B = 0$ , then the C–E path acts as an open switch. There are some extra voltage drops to consider here, but this simple model suffices to understand RTL-gate operation.

### 11.3.2 RTL NOT Gate

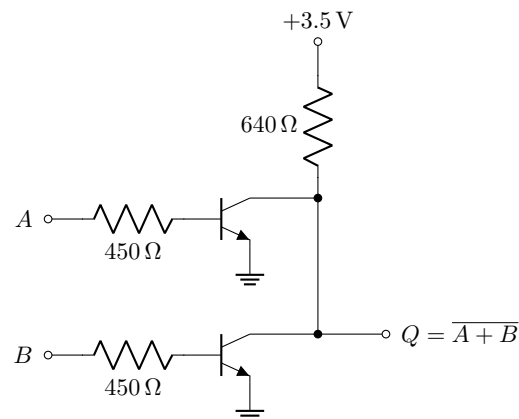
The RTL convention is that +3.5 V is **TRUE**, with 0 V **FALSE**. The simplest RTL gate is an inverter or NOT gate, shown below.



If the input is **TRUE**, then  $I_B > 0$ , and the C–E path conducts. This pulls the output down near ground, or **FALSE**. If the input is **FALSE**, then  $I_B = 0$ , and the C–E path is broken. The resistor pulls the output up to the supply voltage, or **TRUE**.

### 11.3.3 RTL NOR Gate

A slightly more complicated example is the **NOR** gate, shown below.

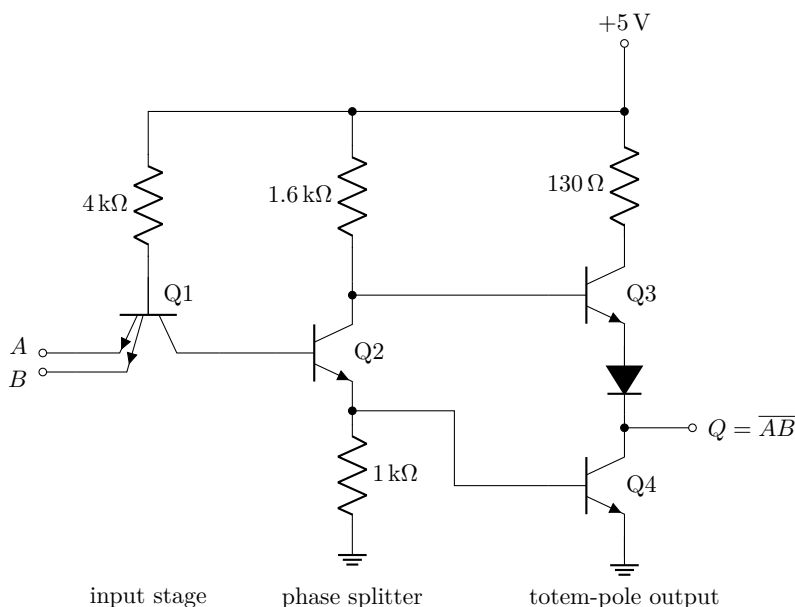


The operation is the same as the **NOT** gate, but here either input can pull the output to ground; the output is only pulled up **HIGH** in voltage if both inputs are **FALSE** or **LOW**.

RTL works reasonable well and doesn't suffer from the (cumulative) degradation problems of DL, because the output levels are set by the supply levels, not the inputs. However, DL is obsolete because the “return” to the high state when the transistors stop conducting is via the pull-up resistor. This transition is slow if there is a significant capacitive load on the output. (The **LOW** transitions when the transistors conduct are fast because the BJT collectors have effectively a very low impedance.)

## 11.4 The Real Thing: Transistor-Transistor Logic (TTL)

A common standard still in modern use is **transistor–transistor logic (TTL)**. The nominal convention is that +5 V is **TRUE**, and 0 V is **FALSE**. The circuitry is somewhat more complicated, and we'll go through the classic TTL **NAND** gate, shown below, as an example.



There are three different stages: the input stage (Q1), the phase splitter (Q2), and the totem-pole output (Q3, Q4, and the diode). An unusual feature is the double-emitter input transistor Q1. It works just like a regular transistor, except that a base current to *either* emitter will switch the collector current. Let's trace the voltages through the circuit for two cases.

1. Suppose *A* or *B* is LOW. Then:

- Q1 is ON (collector conducts to grounded input).
- Q2's base is LOW, thus Q2 is OFF.
- Q3's base is HIGH (pulled up by 1.6-kΩ resistor), thus Q3 is ON.
- Q4's base is LOW (pulled down by 1-kΩ resistor), thus Q4 is OFF.
- The output is HIGH since it is pulled up via Q3 and the diode. The output is  $5\text{ V} - \text{Q3's voltage drop} - \text{the diode drop}$ , which works out to around 3.5 V.

2. Suppose *A* and *B* are both HIGH. Then:

- Q1 is OFF (collector disconnected from inputs).
- Q2's base is HIGH (pulled up via the B–C path of Q1, which acts like a diode), thus Q2 is ON.
- Q2's emitter is pulled LOW by Q4, which is ON.
- Q2's collector is pulled LOW since it is ON; so Q4's base is LOW, and Q3 is OFF.
- The output is LOW since it is pulled down via Q4. The output is  $0\text{ V} + \text{Q4's voltage drop}$ , which works out to around 0.4 V.

The point of all this is to generate a few useful and general observations.

1. The inputs “want” to be high, because they tend to be pulled up to the power-supply voltage via the 4-kΩ resistor and the B–E paths of Q1. Thus, the inputs *source* current when they are pulled LOW. In particular, open inputs are HIGH by default in TTL, and less current flows (less power is dissipated) when the inputs are HIGH. In particular, if you have unused inputs in TTL circuits, it is best to tie them HIGH (i.e., connect them to +5 V).
2. The output, when driving another TTL input, must *sink* current, roughly  $(5\text{ V})/(4\text{ k}\Omega) = 1.25\text{ mA}$ . One output can drive multiple inputs, but there is a limit to this, because the output has a limited current capacity. This limit is called **fanout**, which is typically  $\sim 10$  inputs for a standard TTL output.

3. The output voltages don't quite match the nominal values of 0 V and +5 V, so the TTL standard defines precisely the tolerance limits on signal voltages.
  - TTL circuits must recognize anything from +2.0 V to +5 V as HIGH.
  - TTL circuits must recognize anything from 0 V to +0.8 V as LOW.
  - The intermediate range of +0.8 V to +2.0 V is *indeterminate*: TTL circuits can do anything with inputs in this range and still conform to the TTL standard.

### 11.4.1 TTL Nomenclature

Standard TTL chips are most famously grouped into the 74XX (or 74XXX) family. For example, there are:

- 7400: quad, 2-input NAND (i.e., 4 NAND's per package)
- 7402: quad, 2-input NOR
- 7404: hex inverter (i.e., 6 NOT gates)
- 7408: quad, 2-input AND

and there are hundreds more, though many are now becoming obsolete. Note that these are also labeled as equivalent 54XX circuits, which are the military-grade versions.

These “classic” TTL circuits are now obsolete, but they still come in many popular “flavors.” These variations are labeled by a tag between the 74 and XX, for example 74LS00, 74F00, and 74HCT00 are all basically the same as the original 7400. The common flavors are:

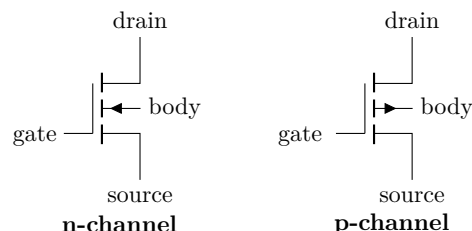
- L: low-power (slow, obsolete)
- H: high-speed (high-speed, obsolete)
- S: high-speed Schottky (high-power, obsolete)
- LS: low-power Schottky (common, modern-standard chip)
- AS, ALS: “advanced” S, LS
- F: fast (gates have ~4-ns propagation delay vs. ~10 ns for standard gates)

## 11.5 The Modern Thing: CMOS Logic

### 11.5.1 MOSFET Review

We have seen MOSFETs before in Section 5.2 while studying the analog side of electronics. Before seeing how they're used in making digital gates, we'll briefly review their operation, in a relatively simplistic way that suffices to explain their digital utility.

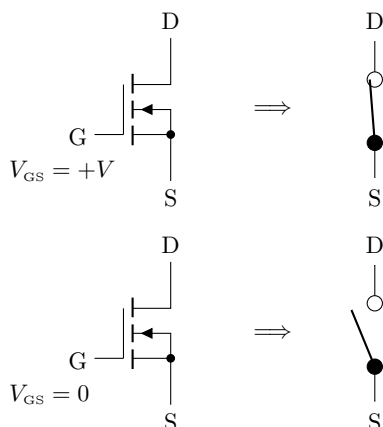
Recall that the MOSFET is a four-terminal device, with terminals labeled as in the diagram below.





To keep things simple, think of the body connection relatively ignorable compared to the others (it's usually shorted to the source lead), so this is more or less a three-terminal device. In this case, the current to be controlled ( $I_D$  for “drain current”) runs from drain to source in the n-channel MOSFET, while current runs from source to drain in the p-channel flavor (for this reason, the p-channel version is typically drawn in an inverted sense compared to what appears above).

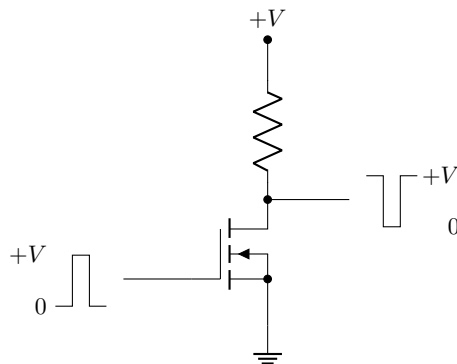
Like the BJT, the MOSFET acts as a switch for current. However, while the BJT was controlled by another current (the base current), the MOSFET is controlled by a *voltage*—namely, the gate–source voltage  $V_{GS}$ . First, consider the n-channel MOSFET. Speaking simplistically, if there is a sufficiently *positive* gate voltage (referenced to the source), then  $I_D$  can flow. If there is *zero* gate voltage, then no current flows. Thus, the MOSFET acts as a voltage-controlled switch (recall the MOSFET analog switch from Section 5.4.5).



For the moment, we won't worry too much about *how* positive  $V_{GS}$  must be to close the switch, but typically in the ON state,  $V_{GS}$  is at least as large as any voltage that would be presented to the drain terminal. And don't worry, we'll come back to how the p-channel MOSFET works, after we take the n-channel version for a test drive. But one more thing to note before carrying on is that the gate is electrically isolated from the rest of the transistor by an oxide insulating layer (the “O” in “MOSFET” is for “oxide”). The consequence of this is that there is some capacitance between the gate and the other terminals, which is important in considering the speed of MOSFET switching.

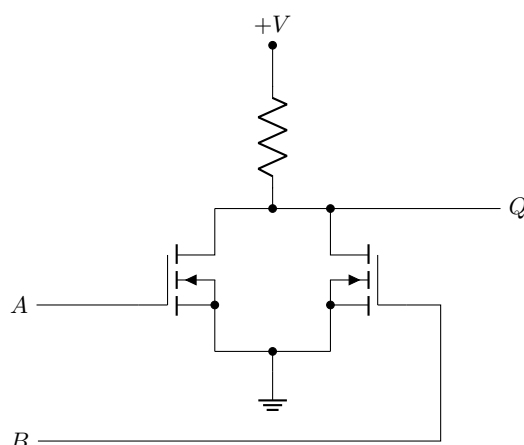
### 11.5.2 NMOS

Other than being under voltage control, the n-channel MOSFET can be used in much the same way as an npn transistor to make gates in the RTL style of Section 11.3, to make **NMOS logic gates** (the “N” here referring to n-channel FETs). For example, a NOT gate is shown below.



This behaves in much the same way as the RTL NOT gate of Section 11.3.2. With zero input voltage ( $V_{GS} = 0$ ), the MOSFET behaves as an open switch, so the output gets pulled up via the resistor to +V. With a positive input voltage, the MOSFET acts as a closed switch, so the output is actively pulled close to ground.

Other NMOS gates are of course possible. By putting two MOSFETs in parallel the circuit below realizes a NOR gate, in the same manner as the RTL NOR gate of Section 11.3.3.

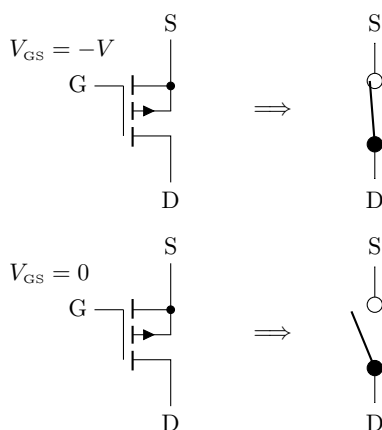


Again, either transistor can pull the output to ground, which means that either input being **HIGH** leads to a **LOW** output. A NAND gate can be made in a similar way by putting the two transistors in *series* with the resistor, so that both transistors must be activated to bring the output low. We'll see an example like this in the next section.

The main disadvantages of NMOS are in speed and power consumption. Like RTL, the resistor has a limited current available to pull the output up to  $+V$  in an upward transition. Also, the input gate capacitance is a limitation to speed—especially when the resistor of one logic gate drives the input capacitance of the next gate. In terms of power consumption, note that (like RTL and TTL) the logic gate draws quiescent current whenever the output is **LOW** (i.e., even if the gate isn't driving any load) because of the nature of the pull-up resistor.

### 11.5.3 PMOS

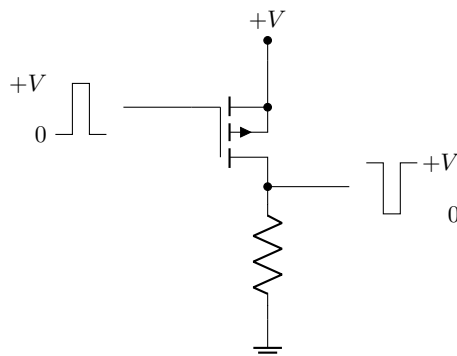
In contrast to the relatively asymmetric RTL and TTL standards, where the npn BJT was king, the p-channel MOSFET is an important counterpart to the n-channel version. The p-channel FET works just like the n-channel, but with the currents and voltages reversed. The diagram below summarizes this.



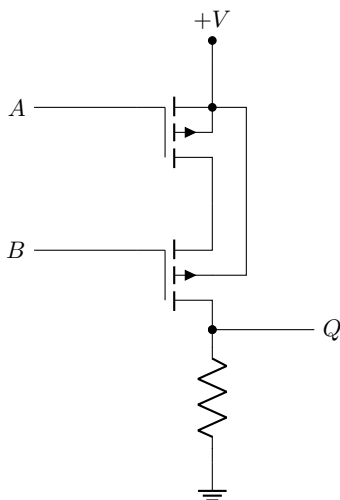
Again the current to be switched flows from source to drain. A zero gate voltage (with respect to the source, which is now drawn as the upper terminal) still causes the “switch” to be open. To close the MOSFET switch, a *negative* gate voltage (again, with respect to the source, so  $V_{GS} < 0$ ) is required to close the gate.

The p-channel MOSFET can be used in just the same way as the n-channel to realize its own logic

family, the **PMOS logic** family. The NOT gate is shown below, and is basically the mirror image of the NMOS NOT.



Another good example is the PMOS NOR gate shown below.



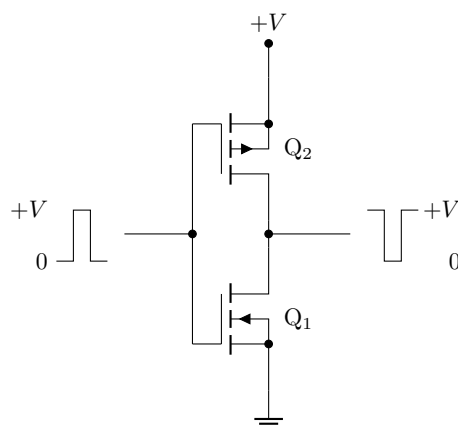
Here, the output is pulled LOW unless both transistors are ON. But the input voltages should be LOW to turn on the transistors. Since the output is LOW unless both inputs are LOW, this is a NOR gate. Note that both body connections are referenced to the same point, as the transistors would be fabricated with common bodies (i.e., common substrates); however, this body-terminal arrangement is not critical (for example, the body of the *B*-input MOSFET could be shorted to its source, in the way that MOSFETs usually come when packages as discrete devices).

Clearly, PMOS suffers from the same limitations as NMOS. PMOS even tends to be inferior to NMOS, but was nevertheless historically important because p-channel MOSFETs are more easily fabricated than n-channels.

#### 11.5.4 CMOS

The **CMOS logic** (“complementary MOS”) family is more complicated than NMOS or PMOS because it integrates both n-channel and p-channel MOSFETs (hence the “complementary”). The NMOS and PMOS families were predecessors because they were easier to fabricate, with CMOS being a particular challenge because the transistors do not share a common body (substrate) type—what usually happens now is that n-channel MOSFETs are fabricated on the p-type substrate, and the p-channel MOSFETs must be fabricated in n-type “wells” in the overall p-type substrate.

A good basic example is the NOT gate shown below. Note that this is just the inverter circuit from Section 5.6.2.



With a zero input,  $Q_1$  is OFF and  $Q_2$  is ON, so the output is pulled (actively) up to  $+V$ . With a positive input, the transistors exchange states, so the output is pulled (actively) down to ground.

Note that the complementary arrangement eliminates the resistor of NMOS/PMOS, along with its disadvantages.

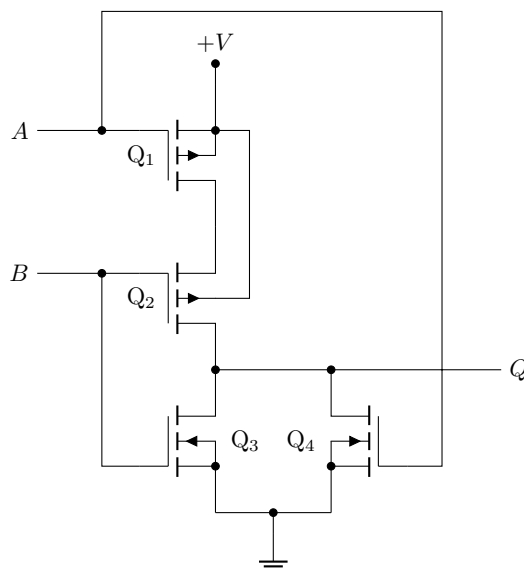
Some important observations we can make from this circuit that hold true more generally for CMOS gates are:

1. Because there is only one path of significant conductance from the output, there is no quiescent current.
2. The output is actively pulled to the appropriate state via a low-impedance path, so that the transitions are relatively fast (and relatively symmetric in terms of speed).
3. Indeed, an important characteristic of CMOS is low power consumption while the gates are static, punctuated by brief periods of current consumption during logic transitions (while gate capacitances are charged or discharged).
4. CMOS logic is also relatively insensitive to noise pickup (compared to NMOS, PMOS, RTL and even TTL) because all voltages are fixed via a low-impedance path to a power supply (or ground).
5. Because the input to the logic gate only drives MOSFET gates, the inputs are not naturally pulled HIGH like TTL gates. So while a TTL input can be operated simply by a switch to ground, a CMOS gate requires a pull-up resistor in addition to the switch.
6. The logic voltages are relatively flexible compared to TTL. The swing in the input level must be above some minimum (threshold) level to make the MOSFETs really turn ON. Typically this means that CMOS gates can operate over a supply range of 3–15 V. However, note that the inputs need to be close to either 0 or  $+V$  to make sure the appropriate transistors are really ON (i.e., to “saturate” the transistors, which really means to put them in the linear region, as in Section 5.3).
7. Relative to TTL, CMOS tends to be slow because of the requirement of charging/discharging gate capacitance. This is the flip side of power efficiency.
8. Because no steady-state current flows into CMOS-gate inputs, CMOS gates tend to have greater fanout compared to TTL gates (at least when driving other CMOS gates). But with increased fanout comes slower speed (due to increased load capacitance).
9. MOSFETs tend to have lower current-drive capability compared to BJTs, and so CMOS gates tend not to be as capable in terms of output drive, compared to TTL gates. Hybrid circuits like BiCMOS use MOSFETs as input stages with BJT output stages, combining the advantages of both types of transistor (in a way similar to JFET-input op-amps that use BJTs for the gain and output stages).

Overall, because of the power efficiency and advanced in fabrication technology (which allow high-density integrated fabrication of MOSFETs and that ameliorate some of the input-capacitance issues), CMOS logic has become overwhelmingly popular, and is by far the most common type of gate in integrated circuits (especially in complex circuits such as modern processors).

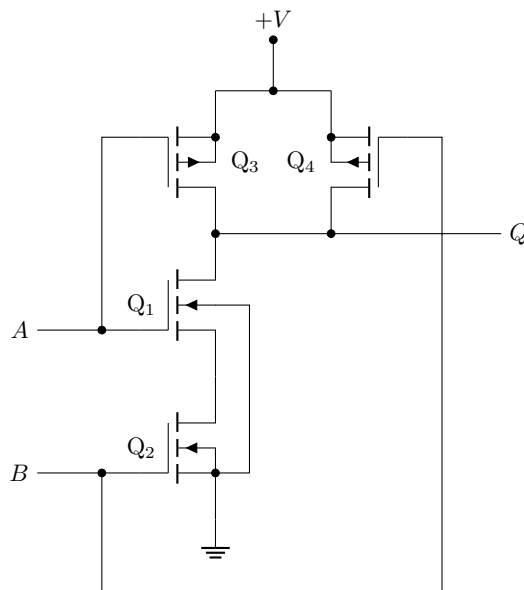
### 11.5.5 CMOS NAND and NOR

To see an example of a more complex CMOS gate, consider the NOR gate below.



This is a kind of combination of the NMOS and PMOS NOR gates above.  $Q_1$  and  $Q_2$  act as the NMOS NOR gate from above, and actively pull the output LOW if either input is HIGH.  $Q_3$  and  $Q_4$  act as the PMOS NOR gate from above, actively pulling the output HIGH only if both inputs are LOW.

The CMOS NAND gate below is something like an upside-down NOR.



Now  $Q_1$  and  $Q_2$  work together as an NMOS NAND gate, pulling the output LOW only if both inputs are HIGH.  $Q_3$  and  $Q_4$  work together as a PMOS NAND gate, pulling the output HIGH if either input is LOW.

The designs here are relatively simple, reference versions of the gates. The “production” gates discussed in the next section typically had a more complex design to improve speed and output driving capability.

### 11.5.6 CMOS Nomenclature

To continue the discussion of TTL nomenclature from Section 11.4.1, we can now discuss some commonly available CMOS devices. Again, an important side effect of CMOS designs is that they dissipate current while switching states, but not when “holding.” This feature makes CMOS generally power-efficient compared to TTL, and CMOS circuits take negligible steady-state input current (but of course are more susceptible to static discharge, which damages the oxide insulating layers of MOSFETs). CMOS devices are more flexible in terms of logic level, and can operate at HIGH voltages other than +5 V. There are separate CMOS logic families, but there are also CMOS variations of standard 74XX devices. For example, some flavors are:

- C: (e.g., 74CXX) operates from +3 to +15 V (compared to +5 V for standard TTL, which typically translates to an acceptable range of +4.5 to +5.5), with a nominal “trigger” range arranged symmetrically about  $\frac{1}{2}$  of the supply voltage. More specifically, for example, at a supply of +10 V, an input below 2 V is recognized (guaranteed) as LOW, and an input above 8 V is recognized as HIGH. The output current at a supply voltage of 5 V is  $\pm 1.75$  mA for the 74C00. (Compare this to the ability of the 7400 to sink 16 mA with a LOW output, but a much weaker current driving ability with a HIGH output.)
- HC: high-speed CMOS, operates from a supply of +2 to +6 V; at a supply of +6 V, an input below 1.8 V is recognized as LOW, and an input above 4.2 V is recognized as HIGH. The output current drive is equivalent to the 74CXX series at +5 V supply voltage.
- HCT: high-speed, compatible levels with TTL (+5 V, and a trigger-voltage range that matches the TTL specification).
- AC, ACT: advanced (i.e., fast) CMOS, ACT is the advanced HCT.

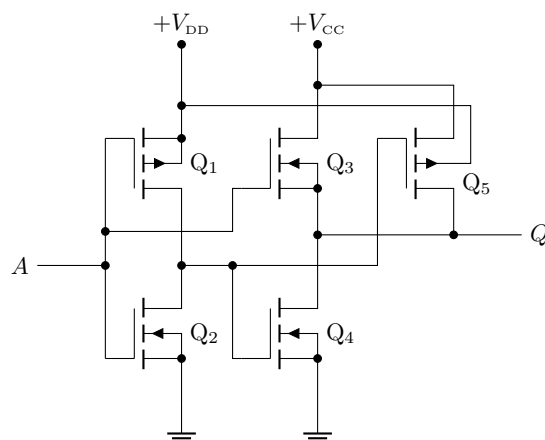
Another classic CMOS line of ICs is the 40XX-series family, which was historically important and still available now. This family is named more simply, so that for example a 4001 is a quad 2-input NOR, a 4011 is a quad 2-input NAND, and so on. The 4-digit part number is preceded by a short, manufacturer specific code (so that the CD4011 was manufactured by Intersil, whereas the MC4011 was by Motorola, or now ON Semiconductor). The 40XX series had the weak output-drive capability that we noted before (at a 5-V supply the output could drive about  $\pm 1$  mA, but increasing to  $\pm 7$  mA at a 15-V supply). For example, one of these CMOS chips can drive a low-power TTL gate (74LXX), but driving a regular TTL gate (74XX) is more complicated, and may or may not be feasible depending on the exact chips involved.

## 11.6 Circuit Practice

### 11.6.1 Mystery RTL/DL Gate

To practice BJT/diode-based logic gates, consider the circuit below. What kind of gate is this?





**Solution.** Start with the logic of the first circuit, input  $A$  at 0 V.  $Q_1$  and  $Q_2$  form a CMOS inverter, so  $Q_3$  is fed with  $\bar{A}$  ( $Q_3$  ON), and  $Q_4$  is fed with  $A$  ( $Q_4$  OFF). This means  $Q$  is pulled HIGH, so this is an inverter. Meanwhile  $Q_5$  is fed with  $A$  ( $Q_5$  ON). The case with  $A$  at  $+V_{DD}$  is similar, but now  $Q_4$  pulls  $Q$  LOW, while  $Q_5$  is disconnected.

The question of the level-shifting action is more subtle, and more detail about the MOSFET is required than what we reviewed here (we will need some more information about the body connection from Section 5.2). Since we know  $V_{DD} \geq V_{CC}$ ,  $Q_5$  acts as a “regular” MOSFET switch (see Section 5.4.5), and  $Q$  is pulled HIGH via both  $Q_3$  and  $Q_5$ . Also, note that  $Q$  cannot be at higher voltage than  $+V_{CC}$ , because  $Q_3$ ’s body–source connection forms a diode from  $Q$  to  $+V_{CC}$ , clamping the output to  $+V_{CC}$ .

The second circuit operates in a similar way. Again suppose input  $A$  is at 0 V.  $Q_1$  and  $Q_2$  again form a CMOS inverter, feeding  $Q_3$  with  $A$  ( $Q_3$  OFF, and  $Q_4$  with  $\bar{A}$  ( $Q_4$  ON). Thus,  $Q$  is pulled LOW via  $Q_4$ , so this is a buffer. The case with HIGH input pulls  $Q$  up via both  $Q_3$  and  $Q_5$  as before.

Note the asymmetry of the output stages.  $Q_4$  can effectively pull  $Q$  LOW, particularly when  $V_{DD}$  is much larger than  $V_{CC}$  (e.g., when interfacing 15-V logic to TTL), because the large gate voltage leads to a low drain–source impedance.  $Q_5$  is the complementary partner of  $Q_4$ , and as a p-channel MOSFET it tends to be weaker than the n-channel  $Q_4$ . Also when  $V_{DD}$  is large, in a downward transition of  $Q$  the gate of  $Q_5$  must drop from  $V_{DD}$  to below  $V_{CC}$  before it even starts to conduct. These effects lead to weaker current drive when  $Q$  is HIGH and a slower LOW-to-HIGH transition.  $Q_3$  is thus here as a “helper,” especially in the transition, since it can begin reacting earlier than  $Q_5$ . (Note that when  $Q$  is HIGH,  $Q_3$  won’t help much if  $V_{CC} = V_{DD}$ , but can still help if  $V_{DD}$  is larger by at least a few volts than  $V_{CC}$ .)

These circuits are, respectively, the internals of the CD4009 hex inverting converter and CD4010 hex buffer converter.<sup>1</sup>

<sup>1</sup><http://www.ti.com/lit/ds/schs020c/schs020c.pdf>



## 11.7 Exercises

### Problem 11.1

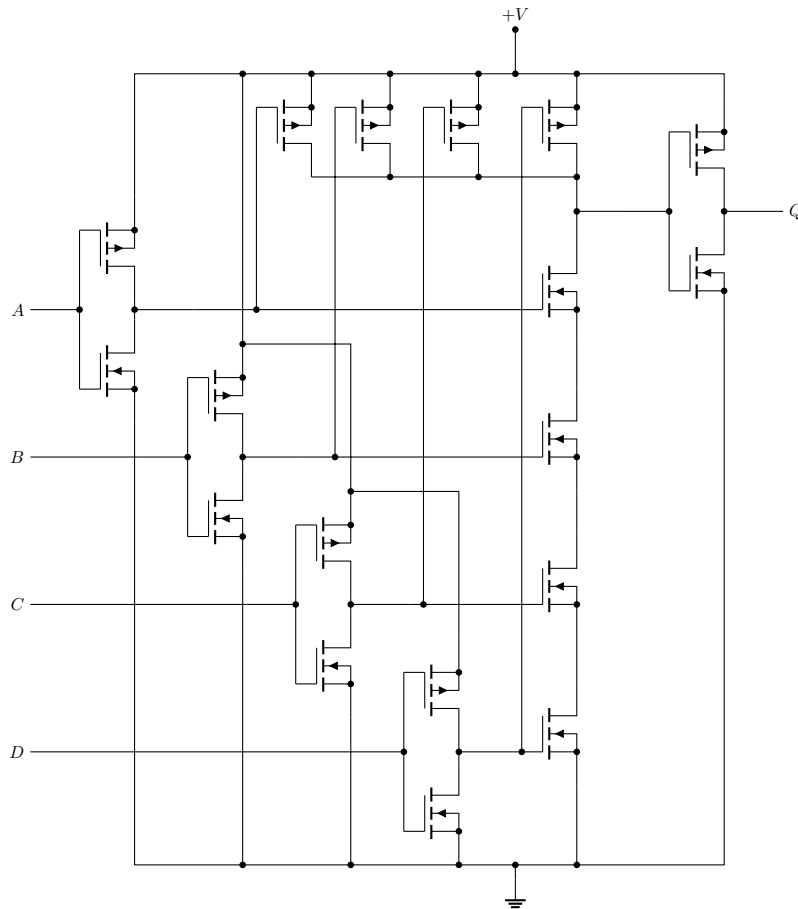
You have two switches (two-position switches; review how switches work if you need to!), a battery, a light bulb, and an infinite supply of wire. Devise a way to realize an XOR gate, where the switch positions are the inputs and the light bulb is the “output.” How about an XNOR gate?

### Problem 11.2

You have three switches, a battery, a light bulb, and an infinite supply of wire. Devise a way to realize the logic expression  $A \cdot (B + C)$ , where the switch positions are the inputs and the light bulb is the “output.”

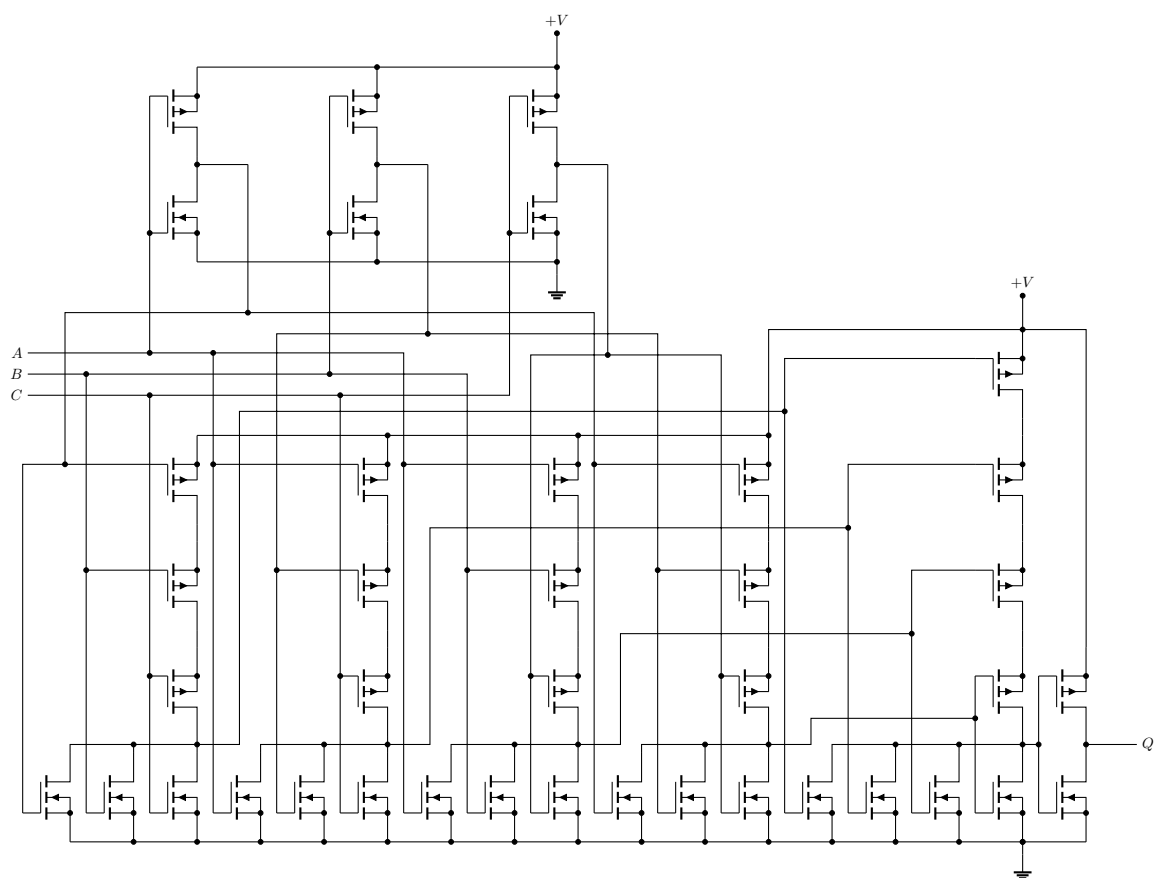
### Problem 11.3

The diagram below shows the internals of a classic CMOS logic device. Identify the gate (i.e., what logical expression does it realize?), and *briefly* describe its operation.



### Problem 11.4

Shown below is a CMOS realization of a logic gate. What kind? **Explain.**





## Chapter 12

# Multiplexers and Demultiplexers

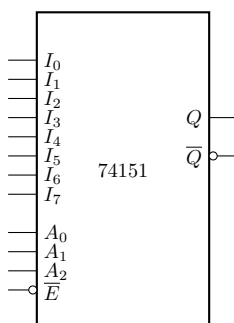
### 12.1 Multiplexers

Simply put, a **digital multiplexer** (or **MUX** for short) is a logic device that maps one of many (digital) inputs to one (digital) output. You select which input to connect to the output using the “address” inputs. The multiplexer is the logic analog of a many-to-one mechanical rotary switch.

Multiplexers are useful devices. For example, you can use them to “pack” data from multiple sources (“parallel data”) onto a single “serial” transmission line (e.g., for phone or computer networks). They can also be used to sample or “poll” data from multiple sources, and ultimately allow scaling of many digital devices into modern computers. Multiplexers are examples of **MSI (medium-scale integration)** devices, “medium-scale” here meaning dozens of gates on 1 chip.

#### 12.1.1 Example: 74151

An example of a multiplexer is the 74151 (which for example, with manufacturer and TTL-flavor codes would be something more like DM74LS151), an 8-input MUX, shown schematically below.



There are a number of features here:

- $I_0$ – $I_7$  are the 8 inputs.
- $A_0$ – $A_2$  are the 3 address lines, to select among the  $2^3 = 8$  inputs; the idea is to select input  $n$  by setting  $A_2A_1A_0$  to  $n$  in binary.
- $Q$  is the output: the selected input is copied to the output.
- $\overline{Q}$  is an inverted copy of the output.
- $\overline{E}$  or  $\overline{\text{ENABLE}}$  (read this as “enable LOW”; this input is also called “STROBE” in the data sheet<sup>1</sup>) is a

---

<sup>1</sup><http://www.ti.com/lit/ds/symlink/sn74ls151.pdf>

“chip enable” line. If  $\overline{E}$  is LOW, the chip works as we have described; if  $\overline{E}$  is HIGH, then  $Q = \text{LOW}$ ,  $\overline{Q} = \text{HIGH}$ , independent of the states of  $I_{0-7}$  and  $A_{0-2}$ .

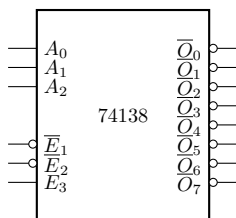
Another example of a common MUX is the 74150, a 16-input MUX (with 4-bit address).

## 12.2 Demultiplexers

A **demultiplexer** (or **DEMUX** for short) is the “opposite” of the MUX, in the sense that a single input is copied to a selected one of many possible outputs. Again, these are useful in, for example, packing and unpacking data to and from a transmission line via a MUX–DEMUX pair. Also a variation on the DEMUX is a **decoder**, which is the same as a DEMUX, but only *selects* the output, without copying any input (the “data” is effectively constant HIGH).

### 12.2.1 Example: 74138

A good DEMUX example is the 74138, a 1-to-8 DEMUX, as shown below.



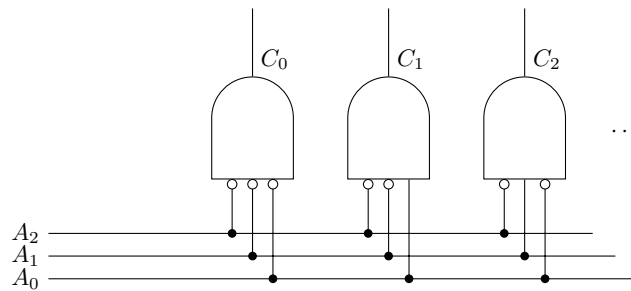
To go over the features here:

- $\overline{O}_0\text{--}\overline{O}_7$  are the 8 outputs, Note that they are inverted (i.e., their “normal,” unselected state is HIGH).
- $A_0\text{--}A_2$  are the 3 address lines, again to select among the  $2^3 = 8$  inputs in the same way as the MUX.
- $\overline{E}_1$ ,  $\overline{E}_2$ , and  $E_3$  are chip-enable inputs. The chip is enabled if  $\overline{E}_1 = \overline{E}_2 = \text{LOW}$  and  $E_3 = \text{HIGH}$ . Then the operation is as follows.
  - If the chip is enabled, then the selected output  $\overline{O}_j$  is LOW. (The others are HIGH.) In this case, the chip acts as a decoder.
  - If the chip is *not* enabled, then all outputs  $\overline{O}_j$  are HIGH.
  - To operate this chip as a DEMUX instead of just a decoder, use  $\overline{E}_1$  or  $\overline{E}_2$  as a data input. In this case, the selected output copies  $\overline{E}_1$  or  $\overline{E}_2$ , while the others remain HIGH. Alternately,  $E_3$  can work as a data input, in which case the selected output copies  $\overline{E}_3$  (with the others still HIGH).

Another example of a common MUX is the 74154, a 16-output decoder/DEMUX (with 4-bit address).

## 12.3 Making a MUX

The logic underlying a multiplexer is not difficult to understand. There are two basic elements: a decoder and “routing” logic. As an example, let’s consider an 8-input multiplexer. The decoder, as we described, takes address inputs  $A_0\text{--}A_2$ , and sets the corresponding one of 8 outputs  $C_0\text{--}C_7$  HIGH, with the others LOW. We can simply use AND gates to set each output when matching the correct address combination, as below.



To be more efficient in terms of gates, we would only use one NOT gate for each of  $\overline{A_0}$ ,  $\overline{A_1}$ , and  $\overline{A_2}$ , rather than NOT gates as shown, but the input NOT operations simplify the diagram.

The routing logic then has the algebraic form

$$Q = (C_0 I_0 + C_1 I_2 + \dots) \overline{\overline{E}}, \quad (12.1)$$

if we include an **ENABLE** input. Notice that the  $E$  is negated twice, which means we're negating the " $\overline{E}$ " input—the chip should be active when  $\overline{E}$  is LOW, so we should AND the output with the *negation* of this (so we do an AND with 1, copying the input to the output). Of course we could just write this as an  $E$ , but this emphasizes the enable-LOW function of this variable.

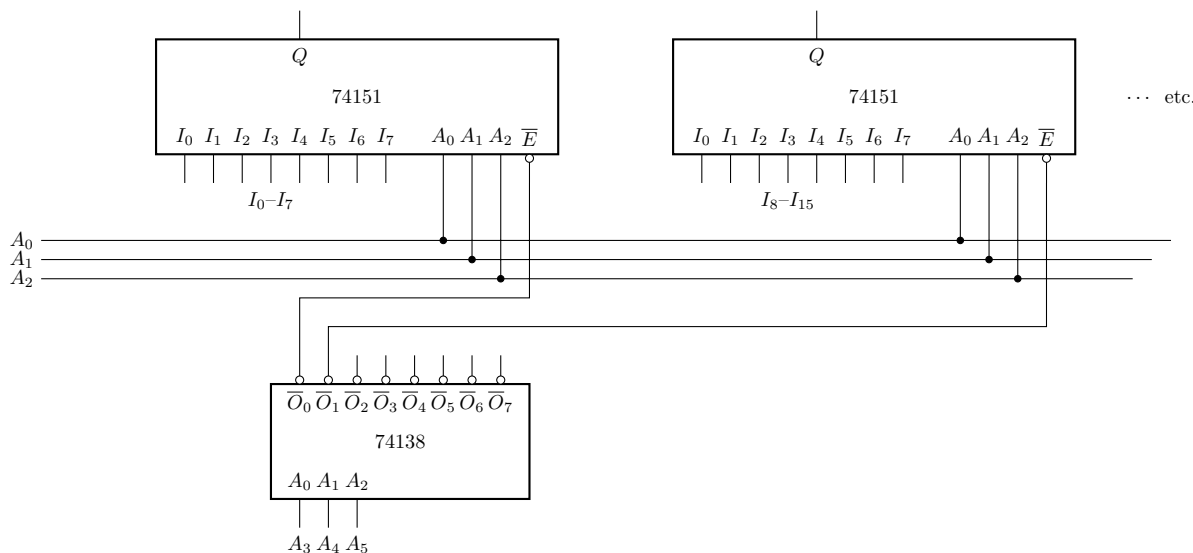
The above expression means that the disabled output state is LOW; this is how the 74151 works. You could also imagine a MUX in which the disabled state is HIGH, in which case

$$Q = (C_0 I_0 + C_1 I_2 + \dots) + \overline{E} \quad (12.2)$$

is the appropriate expression. In the disabled state  $\overline{E}$  is HIGH, forcing  $Q$  to be the same.

## 12.4 Expanding a MUX (or DEMUX)

A useful technique is to combine MUXs into larger MUXs. For example given,  $8 \times 8$ -input MUXs, how do we make a 64-input MUX? This shows the real idea behind having chip-enable inputs: we will use the  $\overline{E}$  inputs and an 8-output decoder, as shown below.

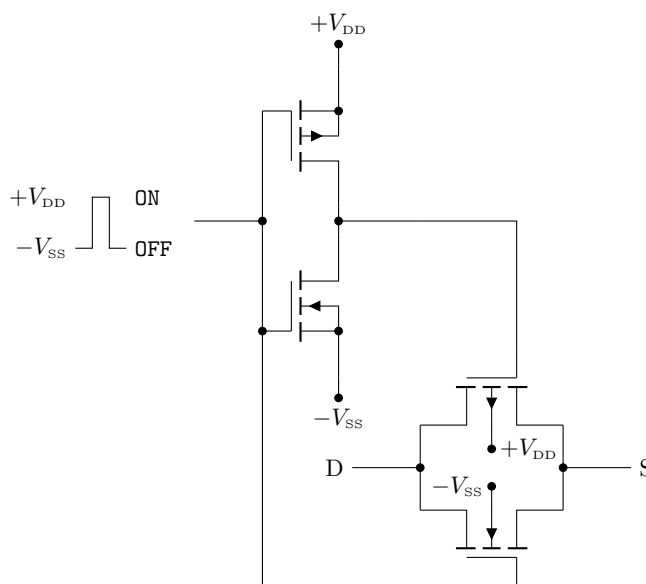


One detail that we have left out is that the  $Q$  outputs must be combined by an OR gate. An important alternative is to use chips with **three-state logic**. For these chips, the output is *disconnected* (high-impedance) when the chip is not enabled. In this case, you can just connect the chip outputs directly together, since only one chip will be enabled at a time. In this example, we can use the three-state alternative 74251 instead of the 74151.

## 12.5 Analog MUX/DEMUX

The same ideas behind *digital* MUX/DEMUX can apply to analog signals as well. An **analog switch** (or **CMOS switch**) is an electronic switch for analog signals, controlled by a digital input. An example is the DG412 quad SPST (normally open) analog switch. The switches are switchable electronically, e.g., from a computer-interface output. They are even good compared to mechanical switches in terms of noise: for example, you can run a digital control wire to the front panel from a circuit board rather than a signal line, and the noise pickup is not critical for a digital control line (it is easier to keep low-noise, critical signals on a well-grounded circuit board than to carry the signals on wires away from the board).

An example of a CMOS realization of an analog switch is shown below. The analog signal can travel between the D and S terminals (in either direction). These are switched by a complementary MOSFET pair (hence “CMOS switch”); the gates of the CMOS pair are driven by a control signal and its inversion from a CMOS inverting gate.



This CMOS-switch circuit can switch analog signals in the range between  $-V_{SS}$  to  $+V_{DD}$ . Note that packaged CMOS switches like the DG412 are usually controllable by a different logic level (typically something like TTL levels). Thus, there is typically also a level-translation layer to change the TTL input into the  $-V_{SS}/+V_{DD}$  logic levels to drive the switching portion. As compared to the previous analog switch (Section 5.4.5), the use of two switching MOSFETs yields a more uniform resistance in the analog path as a function of the analog signal voltage.<sup>2</sup>

An analog MUX/DEMUX (in the analog case, these devices act as bidirectional devices, so there is no distinction between MUX and DEMUX) is an array of analog switches, controlled by address lines. (All that is required is a sufficient quantity of CMOS switches, driven by a decoder.) This is really the analogue of a mechanical rotary-select switch. A good application here, for example, is to connect many sensors to a single microcontroller. An example of an analog MUX/DEMUX is the DG407, a dual 8-channel MUX.

## 12.6 Circuit Practice: Multiplexed Thermocouple Monitor

On the following pages, look over the schematic for a web-enabled thermocouple monitor.<sup>3</sup> A few things to look over:

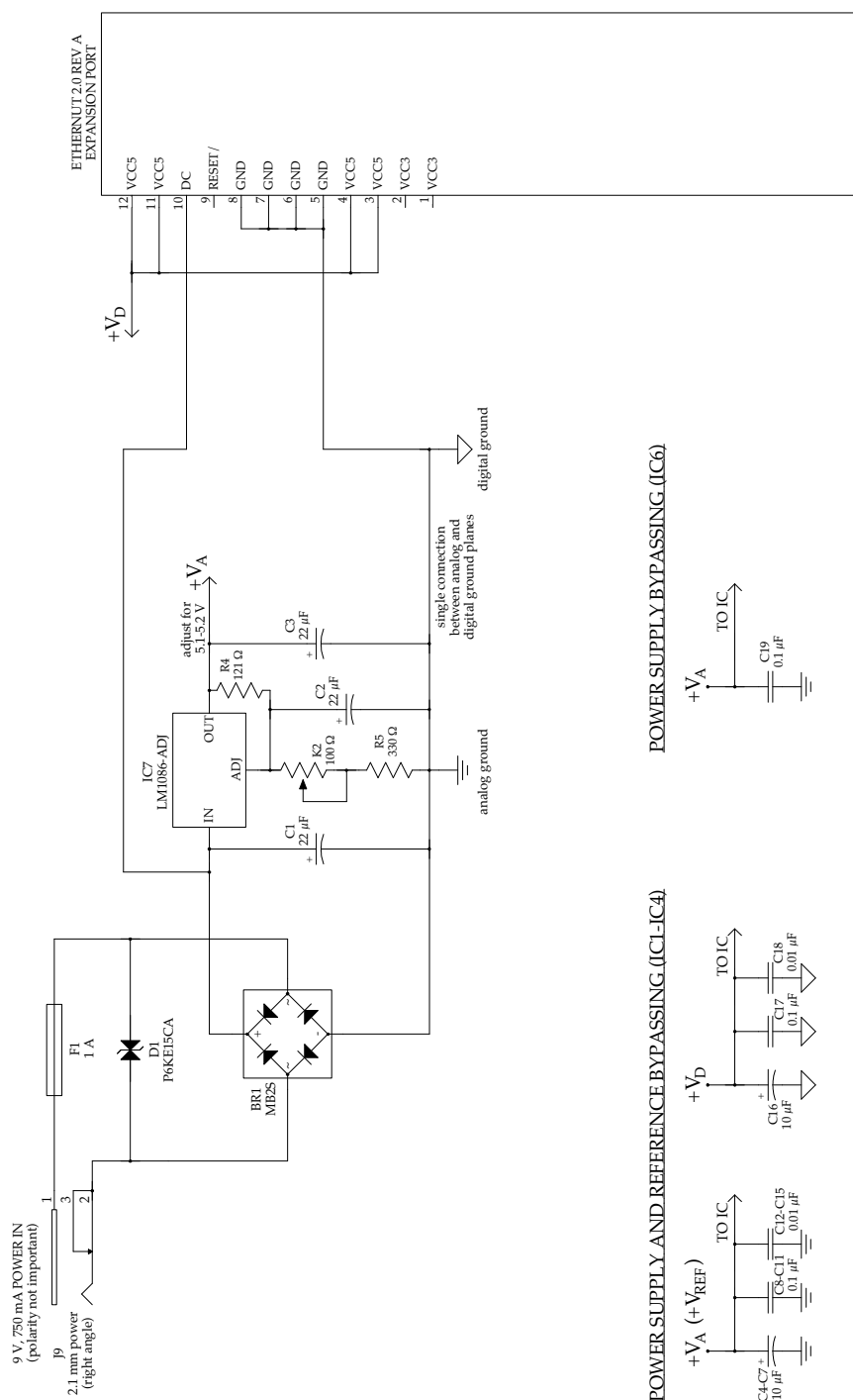
<sup>2</sup>For more details, a good tutorial is Analog Devices tutorial MT-088: <http://www.analog.com/media/en/training-seminars/tutorials/MT-088.pdf>.

<sup>3</sup><http://atomoptics-nas.uoregon.edu/~zoinks/#WebTC>

- There are provisions for 8 thermocouple inputs.
- The thermocouples are monitored by the AD594, which provides “ice point” compensation and buffers the thermocouple signal. However, this is a relatively expensive chip, so rather than having one for each thermocouple, we just use an analog multiplexer (IC1).
- The output is converted to digital via an analog-to-digital converter (IC3).
- The multiplexer address is controlled by an Ethernet microcontroller, which also reads out the ADC. The Ethernet has an ethernet port, and thus has firmware to make a web site to display the temperatures. It also controls an LCD display on the actual box.







## 12.7 Exercises

### Problem 12.1

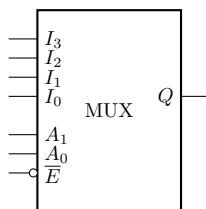
Show how to make a 4-input (digital) multiplexer from ordinary logic gates.<sup>4</sup>

### Problem 12.2

Look up the 74139. What is it? Show how to hook it up, using at most an extra inverting buffer, to make an 8-output decoder. Show how (by adding extra logic gates) to implement an 8-output DEMUX.

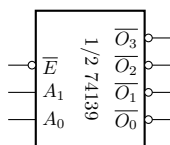
### Problem 12.3

Suppose that you end up stuck on an isolated desert island, surrounded by sand, coconut trees, and (oddly) a near-infinite supply of 4-bit MUXs, shown schematically below. Explain how you can, *in principle*, use your treasure cache to create *any* logic circuit you can dream up, to use to call for help. Make *specific* wiring diagrams to support your argument as necessary (you should have at least one diagram showing a wired-up MUX.) Assume you have found some MacGyverish way to adapt coconuts and coconut fibers to create whatever power supplies and wiring you need. However, you may *not* take the MUXs apart to obtain the individual logic gates inside.

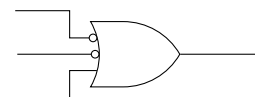
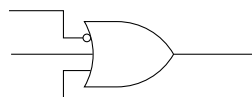
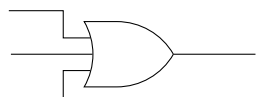
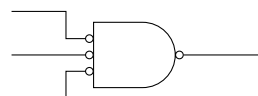
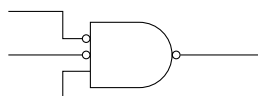
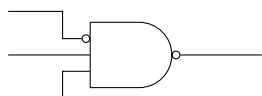


### Problem 12.4

(a) Recall that the 74139 is a 2-bit, 4-output decoder/DEMUX, 2 per package, with 1 enable-LOW input per decoder. One of the DEMUXs is shown below.

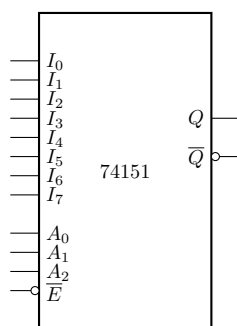


Show how to wire this DEMUX to produce **each** of the following gates (i.e., assign the gate inputs/outputs with the appropriate DEMUX labels).



(b) Recall that the 74151, shown below, is an 8-input MUX. Show how to wire this up to realize a 4-input OR (and simultaneously NOR) gate.

<sup>4</sup>Paul Horowitz and Winfield Hill, *The Art of Electronics*, 2nd ed. (Cambridge, 1989), Exercise 8.17a (ISBN: 0521370957).





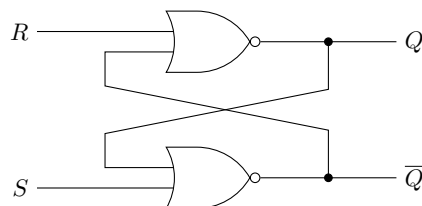
# Chapter 13

## Flip Flops

### 13.1 Flip-Flop Construction: SR Flip Flop

A **flip-flop** is relatively simple logic circuit that involves feedback (i.e., such that the output of a gate drives its own input, generally via other gates). Flip-flops are useful devices, and as we will see, they are the basis of digital memory.

The basic flip-flop is the **SR flip-flop** (“SR” for “set–reset”). A realization in terms of NOR gates is shown below.



To analyze this, let’s work out the truth table.

$S$	$R$	$Q$	$\overline{Q}$
1	0	1	0
0	1	0	1
0	0	1	0
0	0	0	1
1	1	0	0

A few things to notice here: First, there are *two* rows with inputs  $SR = 00$ , with different outputs. You should convince yourself that both are consistent with the circuit. In the first two rows, the fact that one input is **HIGH** fixes the state of the corresponding **NOR** gate, which then fixes the state of the other one. But in this multivalued, or **bistable** state, the inputs don’t fix the state of either gate. Rather, we have to *assume* that  $Q$  is in some state (i.e., it was set in this state in the past), which then fixes  $\overline{Q}$ . **This bistable state is the defining characteristic of a flip-flop:** it means there is **hysteresis** in the circuit, so that the state of the circuit “remembers” the past state. It is in this sense that a flip-flop can act as memory.

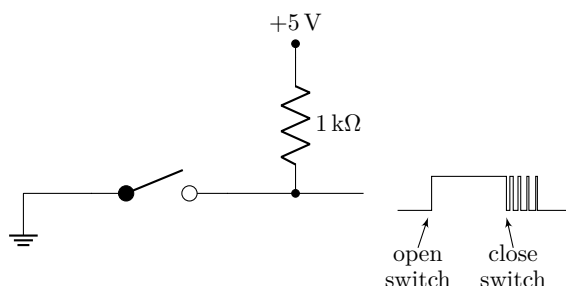
A second feature in the truth table is that the state  $SR = 11$  is a “bad” state, since  $Q = \overline{Q}$ , which means our output notation is in some sense itself bad. However, having the complementary outputs is convenient, even if nonessential. The more important problem with this state, however, is that the outputs don’t match either of the two “hysteresis states,” which we want to use as memory. So if we take the inputs from the bad state to  $SR = 00$ , it will collapse into one or the other hysteresis state in an ill-defined way, which is not very useful. Generally speaking, the bad state is to be avoided when using a flip-flop for its intended purpose.

Then this is how you use a flip-flop:

- The inputs  $R$  and  $S$  are normally 0 (i.e., the flip-flop is in one of the memory states).
- Bringing  $S$  to 1 and back to 0 (the “set” operation) changes  $Q = 1$  and  $\overline{Q} = 0$ . This state is “remembered” when  $R = S = 0$ .
- Bringing  $R$  to 1 and back to 0 (the “reset” operation) changes  $Q = 0$  and  $\overline{Q} = 1$ . This state is “remembered” when  $R = S = 0$ .

### 13.1.1 Application: Debounced Switch

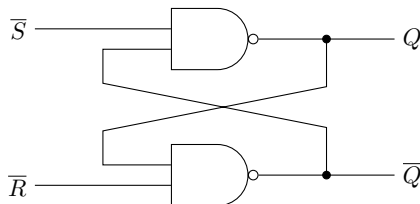
A simple application of the SR flip-flop is to make a **debounced switch**. Recall that switches are mechanical devices that make and break electrical connections. We can use a switch as in the schematic below to toggle between TTL HIGH and LOW.



That is, if the switch is open, the output is pulled up by the resistor to +5 V, while a closed switch corresponds to a 0-V output.

The problem is that the output will really look like the output shown for one open/close cycle. When we open the switch, the output goes **HIGH** with no problem, because the switch cleanly breaks the connection. However, when closing the switch, there is a problem. The contacts must close, and normally they are held together by some spring pressure. But when they close, one contact smacks into the other and “bounces” off of it, just like dropping a chunk of metal on a hard floor. The spring action pushes the contacts together again, and the result is a few extra, short pulses due to the switch bounce, typically on ms time scales.<sup>1</sup> This is a real problem, for example, if the pulse is to drive the input of a counter. For example, the switch could be actuated by items on a manufacturing line, to count the number of items produced; it would obviously not be a very good count if there were several extra bounces for each item to count.

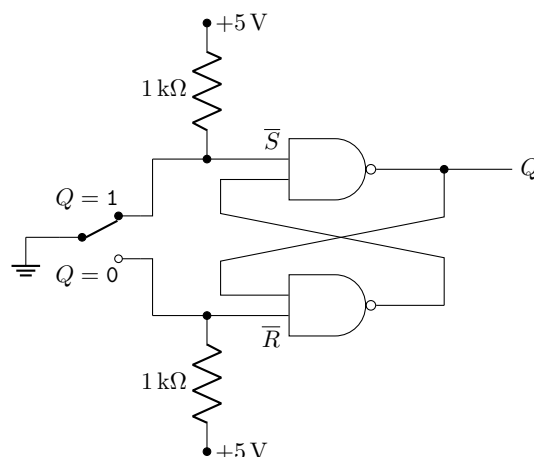
A simple solution to this uses a flip-flop and a slightly more complicated switch. Before getting to that, let’s introduce a functionally equivalent variant of the above RS flip-flop, now based on NAND gates.



The operation is the same as before, but note the inputs are  $\overline{R}$  and  $\overline{S}$ , so their senses are inverted. That is, the “usual” input state should be  $\overline{R} = \overline{S} = 1$ . Then you bring  $\overline{S}$  momentarily to 0 to set the flip-flop (i.e.,  $Q = 1$ ), and you bring  $\overline{R}$  momentarily to 0 to reset it ( $Q = 0$ ). We will leave the analysis of this flip-flop as a circuit-practice exercise.

Now the debounced switch uses an SPDT switch (the “bouncy” switch used an SPST switch). The “up” switch state sets  $Q = 1$ , and the “down” switch state sets  $Q = 0$ .

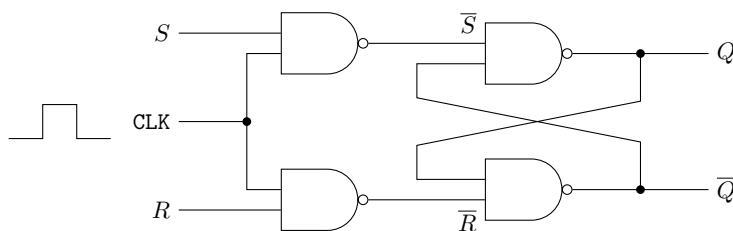
<sup>1</sup>example: <http://www.maximintegrated.com/en/app-notes/index.mvp/id/287>



During a bounce, the switch makes no connection to either contact, so both  $\overline{S}$  and  $\overline{R}$  are 1. This is the memory state, so the flip-flop holds the last switch state, which persists through the duration of the bouncing.

## 13.2 Clocked Flip-Flops

An important class of flip-flops, one step up in sophistication from the basic flip-flops above, is that of **clocked flip-flops**. A **clock** is an external, typically periodic logic signal that synchronizes signals in complex circuits. We will see some examples later, but in complex circuits, this synchronization is important in avoiding problems with race conditions. The idea in a clocked flip-flop is that the input datum is only accepted during a particular phase of the clock cycle, for example when the clock is **HIGH**. The clock then functions as a “gate” for the input data. An example of a clocked SR flip-flop is shown below.



Note that when the CLK signal is **LOW**, this guarantees that the flip-flop is in the memory state. We can write the truth table for this circuit as follows.

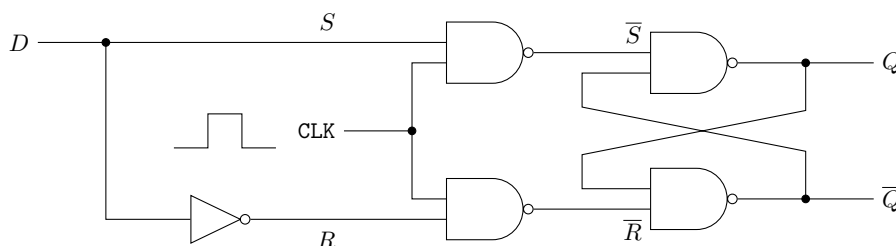
$S$	$R$	$Q_{n+1}$
0	0	$Q_n$
1	0	1
0	1	0
1	1	“bad”

Here,  $Q_n$  is the output state after the  $n$ th pulse. If  $S$  and  $R$  are **LOW**, then the flip-flop stays in the memory state, and the state  $Q_n$  persists to the next clock cycle. Otherwise, the clocking **NAND** gates act as inverters for the  $S$  and  $R$  inputs, so the inputs set and reset as usual (with momentary **HIGH** action), but only when the clock is **HIGH**.

### 13.2.1 D-Type Flip-Flop

An important class of clocked flip-flop is the **D-type flip-flop**, which is basically the clocked SR flip-flop, but where the two inputs are always in opposite states, as shown below. Here, the  $D$  or “data” input drives the  $S$  and  $R$  inputs oppositely.

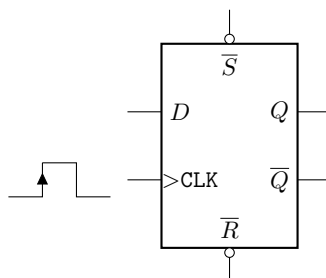




The idea is that there is now only one input, and the flip-flop latches the value of the data while CLK is HIGH. This flip-flop is often called a **data latch**.

### 13.2.2 Edge-Triggered, D-Type Flip-Flop

A somewhat more sophisticated and realistic D-type flip-flop is the **edge-triggered, D-type flip-flop**. Here, “realistic” means you can buy these prepackaged (e.g., the 7474 gives you two of these per chip). The main difference is that the data-latching action happens on the *rising edge* of the CLK pulse. Since the edge has a short duration compared to the HIGH phase of the clock, the timing is more precise in this convention. Schematically, this flip-flop is shown below.



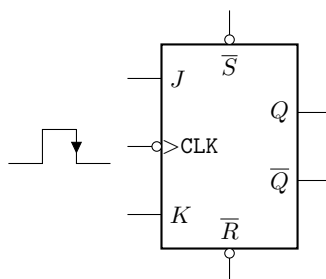
Some things to notice:

- $D$  is the data input, as in the regular D-type flip-flop.
- CLK is the clock input. Again, the datum is latched on the rising edge. Schematically, it is common to indicate this by drawing an arrowhead on the edge of a sample clock pulse, as shown, and also to include a “>” next to the CLK pin.
- $Q$  is the output as usual, and  $\bar{Q}$  is an inverted output copy.
- $\bar{S}$  and  $\bar{R}$  are “jam” set and reset inputs. These override the output, independent of the CLK state (so they work just like the inputs to the SR flip-flop). These are often called  $\overline{\text{PRE}}$  and  $\overline{\text{CLR}}$  (for preset/clear).

This flip-flop is good, for example, for storing data until they are “passed on” to a computer (e.g., in data-acquisition systems, when data arrive with timing determined by a physical system, but need to be loaded into a computer with its own timing).

### 13.2.3 JK Flip-Flop (Edge-Triggered)

A slightly more complicated variation on the edge-triggered, D flip-flop is the **edge-triggered, JK flip-flop**. This is like the D version, but there are *two* data inputs ( $J$  and  $K$ ), with no indeterminate states for  $J$  and  $K$ . This is available in the 74112/74112A (2 per chip), and the now-obsolete 7476/7476A (also 2 per chip). The flip-flop is shown schematically below.

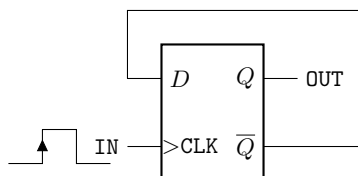


The operation with the two new inputs is as follows.

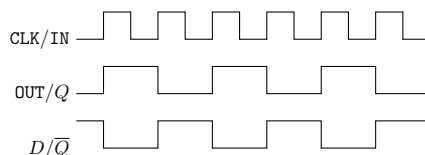
- The CLK on this flip-flop, as drawn, triggers on the *falling* edge of the clock. Note the NOT circle on the clock input, and the sample clock pulse. This is how the 74112/74112A and 7476A work (the plain 7476 triggers on the positive edge). That is, this device is **negative-edge triggered**.
- If  $J = 0$  and  $K = 0$ , then  $Q$  persists on the next CLK pulse.
- If  $J = 1$  and  $K = 0$ , then  $Q = 1$  on the next CLK pulse.
- If  $J = 0$  and  $K = 1$ , then  $Q = 0$  on the next CLK pulse.
- If  $J = 1$  and  $K = 1$ , then  $Q$  **inverts** on the next CLK pulse (i.e., it “toggles”).
- $\bar{S}$  and  $\bar{R}$  are still jam set and reset inputs.

## 13.3 Counters

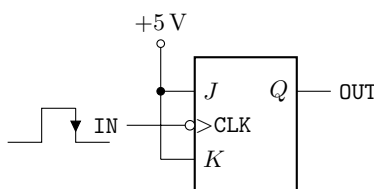
One useful application of flip-flops is in realizing **counters**, which count input pulses by incrementing a binary output. The basic building block of a counter is the **divide-by-2 counter**, shown below in terms of a D-type flip-flop.



The timing diagram for this circuit is shown below. Note that transitions happen on the rising edge of the input (clock) pulses, and essentially the output is just toggling its output on each clock cycle. Hence the term “divide-by-2,” since the output pulse train oscillates at half the frequency of the input clock. More specifically, the flip-flop loads  $D = \bar{Q}$  to  $Q$  on each rising pulse-edge.

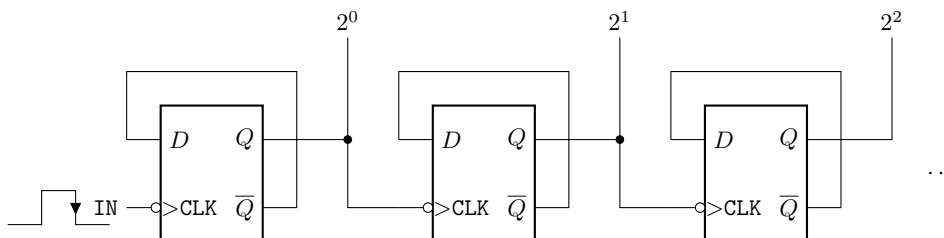


Since the output is just toggling, recall that we can also make a JK flip-flop do this by tying both  $J$  and  $K$  inputs **HIGH**, as shown below.

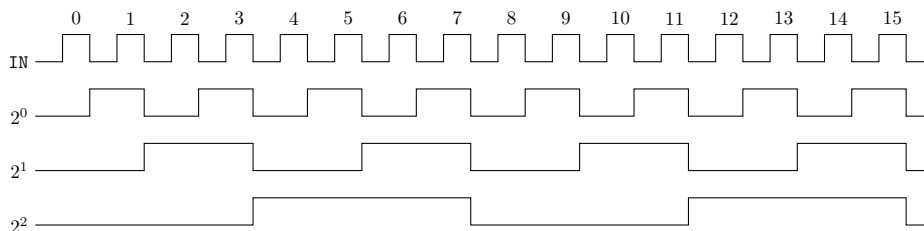


### 13.3.1 Asynchronous (Ripple) Counter

Generalizing the divide-by-2 is relatively easy. For example, we can make a divide-by-4 counter by cascading 2 divide-by-2 counters, and by chaining 3 of them, we make a divide-by-8 counter. Chaining  $n$  counters realizes a divide-by- $2^n$  counter, as shown schematically for D-type flip-flops below (first three bits are shown).



The timing diagram is shown below. Note that we changed the convention for the flip-flops, which now trigger on the *falling* edge of the clock pulse. (Why do we need to trigger on the falling edge? How would you modify the circuit if the flip-flops triggered on a *positive* edge?)



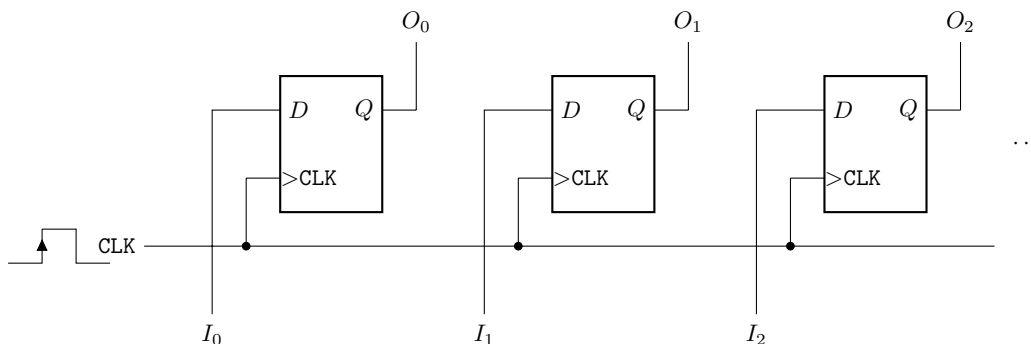
The main advantage of this circuit is that it is easy to build: it's easy to chain together flip-flops. The main disadvantage is the **asynchronous** or “**ripple**” operation of the circuit: since there is a finite propagation delay of the logic signal through each flip-flop, it takes some time for each clock pulse to “ripple” through a long chain of a many-bit counter, which can cause synchronization problems for fast input signals (i.e., spurious output states may be present for some or even all the time).

## 13.4 Memory and Registers

Flip-flops act as single-bit memory devices, as we have seen. Combining flip-flops, we can build up **registers**, which act as multi-bit memories.

### 13.4.1 Register

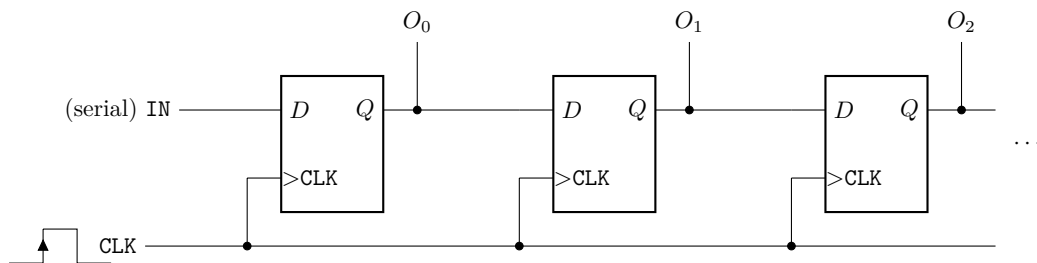
The basic register is an array of D-type flip-flops, which synchronized CLK inputs, the idea being to latch all the bits at once (to avoid timing problems, e.g., as in the ripple counter).



Once latched, the register holds the output state, independent of the inputs, until the next clock pulse. One application is where a shared set of data logic lines drives several devices; a register at the input of each device can hold the relevant logic data for each particular device while the data lines drive other devices.

### 13.4.2 Shift Register

The **shift register** shifts the data among the outputs, shifting all bits in one direction on each clock cycle.

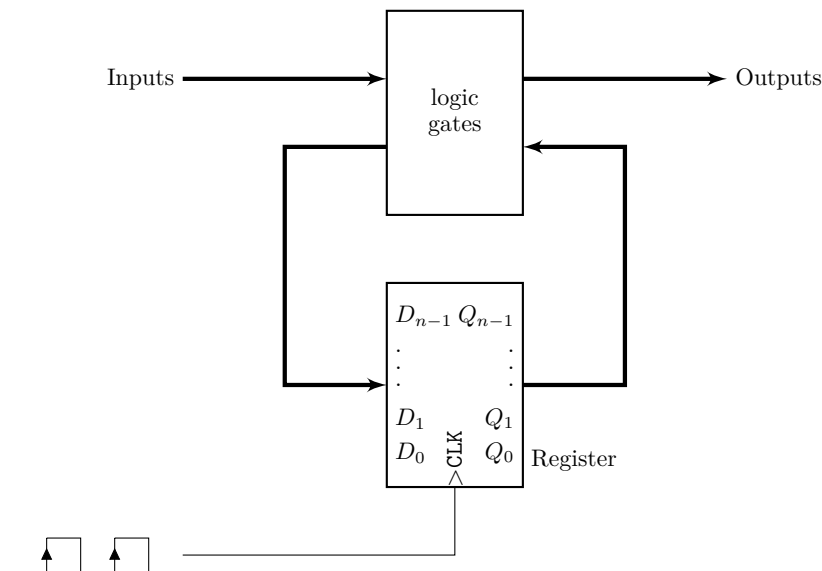


An example application is in converting serial data to parallel form (i.e., on the receiving end of a serial transmission channel). Also, note that a bit shift corresponds to a mathematical operation on binary data (divide/multiply by 2).

## 13.5 Sequential Logic and the State Machine

Recall in our discussion of asynchronous circuits (e.g., the ripple counter), we mentioned that there can be timing problems if the signals change rapidly, such that the gate delays are comparable to the time between transitions. The cure for this is to use **synchronous circuits**, where all logic transitions happen just after each clock pulse (or, more commonly and precisely, at an edge of each clock pulse). The transitions occur based on the logic levels present just *before* each clock pulse (edge). This is essentially what happens in a microprocessor, and this system of synchronous, clocked transitions is essential for high-speed and high-complexity logic systems.

The general scheme of sequential logic is shown in the diagram below.



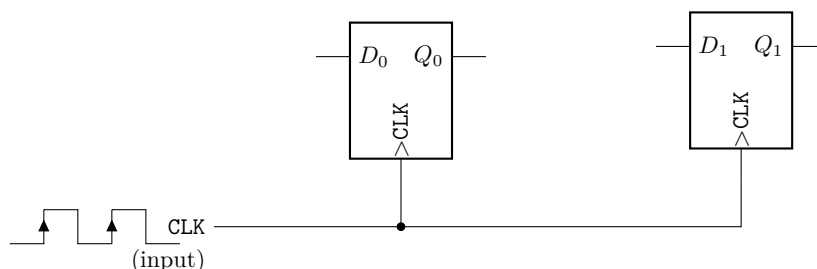
The first main ingredient is a register (Section 13.4.1), which is again an array of D-type flip-flops, with a common clock input. On the rising clock edge, the inputs  $D_j$  are all transferred to the outputs  $Q_j$ , and then held.

The next main ingredient is a set of logic gates in some combination, which work to transform the outputs  $Q_j$  into new inputs  $D_j$ . These can be implemented generically using **PAL (programmable array logic)** devices or **PLA (programmable logic array)** devices, which are basically configurable arrays of many logic gates. Also available are **registered PAL/PLA** chips, which contain flip-flops and gates all on one chip. These are usually called **PLD's (programmable logic devices)**.

Inputs and outputs to the logic gates also permit interaction with the outside world. The sequential logic scheme here is the most general form of digital logic, even though the idea is schematically simple.

### 13.5.1 Example: Synchronous, Divide-by-3 Counter

Here we will illustrate sequential logic by constructing a synchronous, divide-by-3 counter.<sup>2</sup> We will need two flip-flops (two bits to accommodate counting to 3), clocked from the counter input. We will call the register inputs  $D_0$  and  $D_1$  with corresponding outputs  $Q_0$  and  $Q_1$ .



To design the counter, first we will choose the sequence of states that we want. There are no external inputs here, so this is just a simple sequence, with no conditions (which we would represent as extra bits here). The counting sequence is thus as follows:

$Q_0$	$Q_1$
0	0
0	1
1	0
0	0

Here,  $Q_0$  functions (arbitrarily) as the MSB, and  $Q_1$  the LSB. We have also shown the first step in the repeat.

The next step is to find the appropriate  $D$ 's. Remember the  $D$ 's must be our desired  $Q$ 's on the next step, so we will explicitly write out the desired  $D$ 's as a function of the  $Q$ 's.

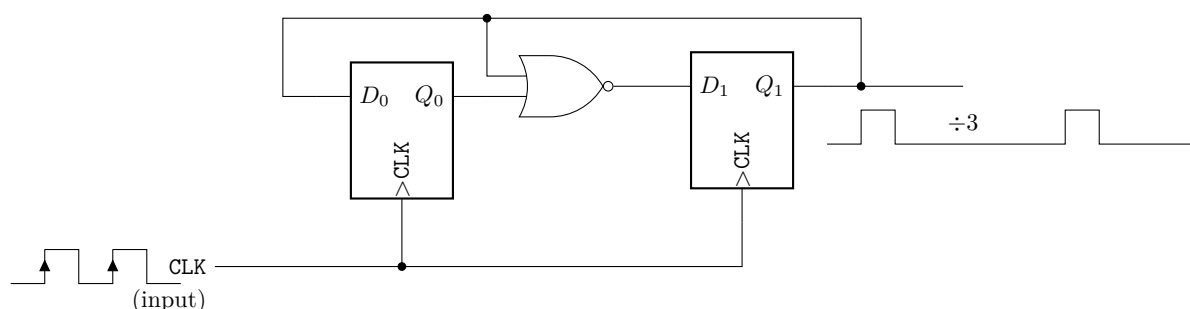
$Q_0$	$Q_1$	$D_0$	$D_1$
0	0	0	1
0	1	1	0
1	0	0	0
0	0	0	1

Finally we find a logic implementation of the functions  $D_j(Q_k)$ , using whatever means necessary (e.g., Karnaugh maps). Here, by inspection we can see that

$$D_0 = Q_1, \quad D_1 = \overline{Q_0} + Q_1. \quad (13.1)$$

Thus, the circuit implementing the counter is shown below.

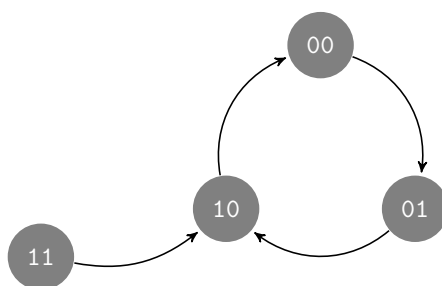
<sup>2</sup>This example from Paul Horowitz and Winfield Hill, *The Art of Electronics*, 2nd ed. (Cambridge, 1989), p. 513 (ISBN: 0521370957).



One more detail to worry about is the set of **excluded states**. For the divide-by-3 counter, we didn't use the state  $Q_0Q_1 = 11$ , but what if the flip-flops end up in this state (e.g., when the circuit is turned on)? Given our logic realization, we will then have  $D_0 = 1$  and  $D_1 = 0$ , so the counter will resume the normal counting cycle on the next cycle, with  $Q_0Q_1 = 10$ . But it's important to consider these, since the register could end up "frozen" in an excluded state or a cycle of excluded states if the logic gates don't handle them properly.

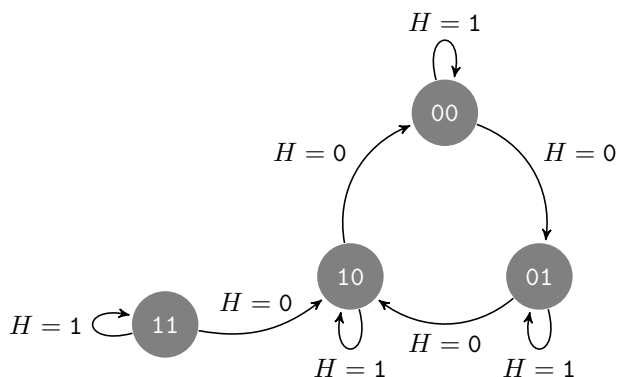
### 13.5.2 State Diagrams

A convenient way to represent the operation of a state machine is a **state diagram**. For example, the diagram (including the excluded state) for the divide-by-3 counter is shown below.



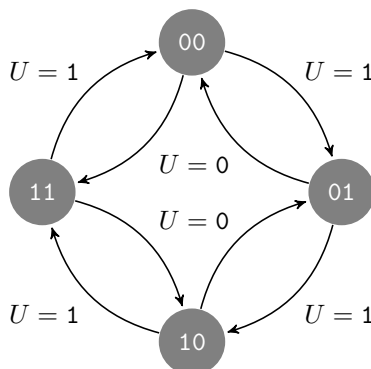
This is a directed graph of all the register states; the edges (arrows) show transitions between the states. You can use a diagram to start the design: you just begin with  $n$  flip-flops, where  $2^n$  equals or exceeds the number of distinct states. Then you use the procedure that we used for the divide-by-3 counter to generate the appropriate logic.

If there are inputs, then there can be multiple possible transitions depending on the inputs. In this case, you can modify the diagram, labeling the transitions according to the input state. For example, below is a divide-by-3 counter with hold. That is, a "hold" bit  $H$  causes the counter to hold its state when  $H = 1$  and counts as before when  $H = 0$ .



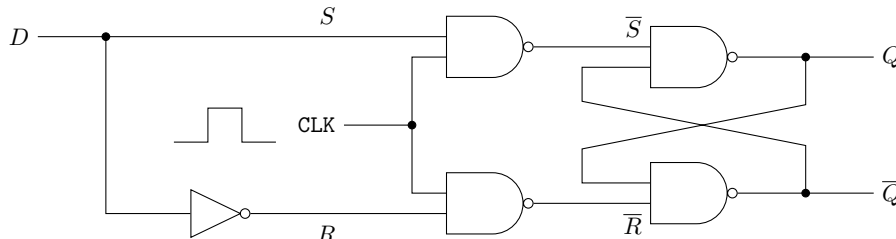
This example also illustrates “transitions” where the state remains the same. We will cover the implementation of this counter in Section 13.7.1.

Another example is the 2-bit up/down counter, where an input bit  $U$  controls whether the (divide-by-4) counter counts up or down.

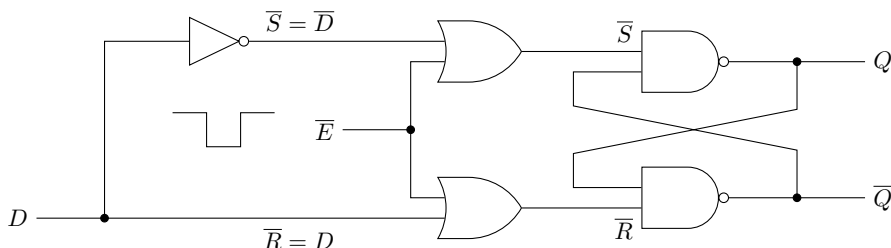


## 13.6 Memory

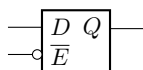
When we introduced the flip-flop, we noted that it is the basic building block of memory, since it is a simple 1-bit memory device. Now we will talk about how to build up many flip flops into a memory device. First, however, we will refer to the (non-edge-triggered) D-type flip-flop from Section 13.2.1. Recall that this is as shown below.



An RS flip-flop (with momentary-low inputs) has two extra gates on the inputs and a clock signal, so that the flip-flop only accepts input from the data line  $D$  when  $\text{CLK} = \text{HIGH}$ . The clocked version of the flip-flop is essential in scalable memories. We will actually refer to a modified version of this circuit, shown below.



Here we changed notation so that the clock signal is now an “enable-low” ( $\overline{E}$ ), and the flip-flop now accepts input only when  $\overline{E} = \text{LOW}$ . To keep things compact, we will refer to this circuit with the small block diagram below.



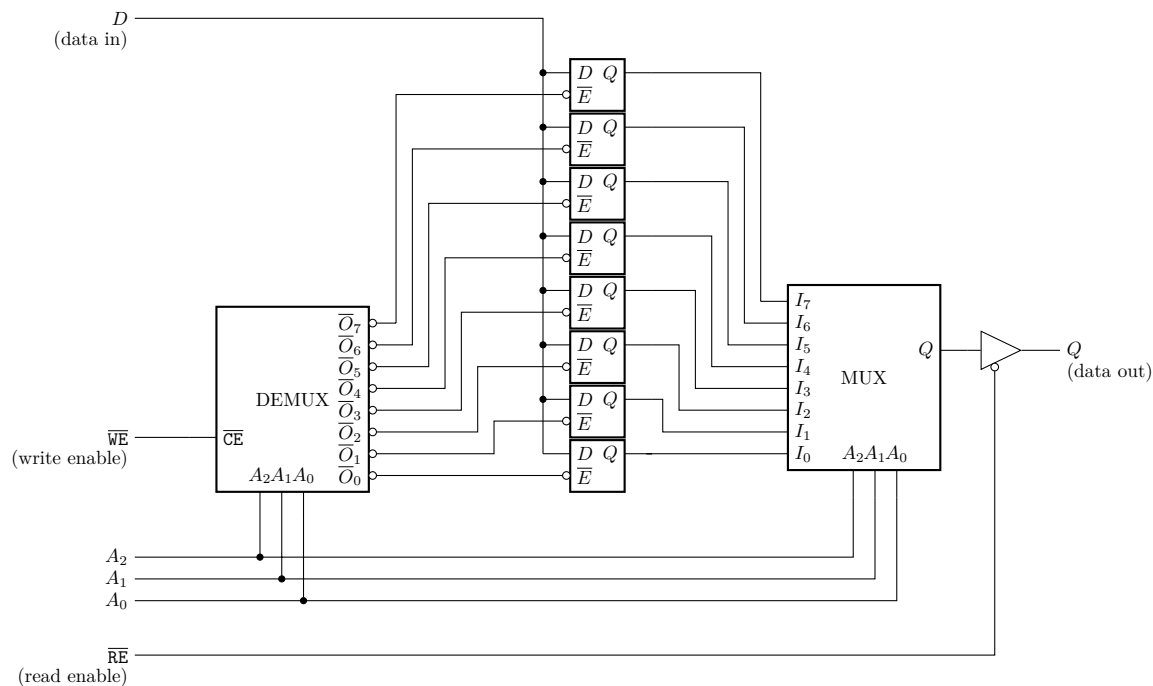
Now, modern computers can have GB of memory on the low end, easily exceeding  $10^{10}$  flip flops. Obviously it would be awkward to have  $10^{10}$  connections to a processor or other device from the memory, so how do we handle this? The answer is to go back to multiplexers and demultiplexers.

### 13.6.1 Example: $8 \times 1$ -bit RAM

Below is an example arrangement for a memory circuit with “ $8 \times 1$ -bit RAM.” This means

- There are 8 “slots” for 1 bit each of memory data (i.e., 1 flip-flop per slot).
- “**RAM**” means **random access memory**, meaning we can easily write and read data to and from any location in memory in any order (as opposed to sequential memory, as on a magnetic tape, or shingled magnetic storage on some modern hard disks).
- Flip-flop based memories like this are called **static RAM (SRAM)**.

The circuit is shown below, with external connections (larger versions of this circuit would come packaged in a single integrated circuit (IC) with similar external connections).



There are a few elements worth noting.

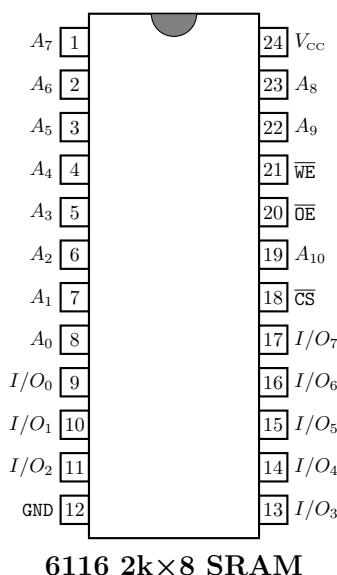
- 8 flip-flops here do the actual storage.
- The address lines  $A_0$ – $A_2$  select which of the 8 flip-flops (slots) is active for either reading or writing.
- For larger memories, the address lines would select a *register* of flip-flops so more data can be transferred in parallel. For example, if the flip-flops were each 4-bit registers, we would have  $8 \times 4$ -bit RAM.
- The data input  $D$  is wired to all of the flip-flop inputs; the DEMUX only enables one flip-flop to accept data according to the address.
- The DEMUX must also be enabled by an enable input here  $\overline{WE}$ , or write-enable-low. If this input is high, then *no* flip-flop is enabled for data input.



- On the readout side, all flip-flop outputs are fed into a single output via a single MUX, which is addressed by the same address lines as the DEMUX (which is not necessary, but saves address lines if we only want to read *or* write, but not both simultaneously).
- The output  $Q$  is buffered by a three-state buffer gate, which is enabled by the  $\overline{RE}$  (read-enable-low) input. In this way, several memory chips can share the same output line(s), with only the selected chip attempting to assert a logical value on the shared data output line.

### 13.6.2 Example: 6116 SRAM

In an example that is more typical than the toy example above, the data input and output lines are the *same*. (Another reason to have a three-state output, so it does not conflict with incoming data.) A old classic, the 6116 SRAM chip, is still available;<sup>3</sup> this is a 2-kB memory ( $2\text{ kB} = 2048 \times 8\text{ bits}$ ), and the connections are shown below.



Most elements are similar to the toy model above.

- The  $A_0$ – $A_{10}$  lines form the **address bus**: 11 bits are necessary to address from 0–2047.
- The  $I/O_0$ – $I/O_7$  lines form the **data bus**, which can serve as inputs or outputs (i.e., reading and writing) for the stored data, 1 byte at a time.
- The **chip-select-low** input  $\overline{CS}$  (i.e., chip enable) enables the action of the chip when **LOW**. Again, this is useful when several chips share the data bus, so only the enabled chip can write to the bus.
- The **write-enable-low** input  $\overline{WE}$  enables the latching of the input data (the  $I/O$  lines act as inputs).
- The **output-enable-low** input  $\overline{OE}$  enables the data lines to act as outputs.
- If  $\overline{CS} = \text{HIGH}$  or  $\overline{WE} = \text{LOW}$  or  $\overline{OE} = \text{HIGH}$ , then the output buffers are in the high- $Z$  state, again so the data busses of many chips can be connected together. (The address busses are also connected, but these need only act as inputs here.)

<sup>3</sup>For data sheet, see <http://www.idt.com/products/memory-logic/synchronous-and-asynchronous-sram-memory-devices/asynchronous-sram-async-sram/6116-50v-2k-x-8-asynchronous-static-ram>; availability for small quantities starts around \$6 each as of May 2015.

### 13.6.3 Other Memory Types

The SRAM above is the basic type of digital memory, but there are many other types. We will briefly review them here.

- **DRAM (dynamic RAM)** uses a small capacitor to store a bit of information as a charge state. The disadvantage is that because the capacitor leaks, it must continually be “refreshed” (on ms time scales), which greatly complicates the overhead circuitry. The advantage is that DRAM is cheap and highly scalable; DRAM is standard for large memory modules in modern computers.
- **SRAM (static RAM)** we have already talked about. Why do we want it? It is complicated to fabricate relative to DRAM, and is hard to scale to very large memory. However, it can be relatively power efficient (no refreshing is necessary), and due to the lack of refreshing overhead, it can be much easier to use in small projects. Note that while no refresh is needed, SRAM is volatile (the stored information vanishes when the circuit is not powered).
- **ROM (read-only memory)** is not intended to be written, just read. The data are written during manufacturing.
- **PROM (programmable ROM)** is ROM that can be written (once!) using special programming hardware, which burns fused connections inside the chip by applying relatively high currents.
- **EPROM (erasable PROM)** is PROM that is programmed electronically. The programming can be erased by exposing the integrated circuit (IC) to ultraviolet light, through a transparent window in the IC’s package. The photo below shows a nice example. It is a “PIC” microcontroller<sup>4</sup> with integrated EPROM, and you can see the circuit itself through the window, as well as the wire bonds that connect to the pins if you look closely.

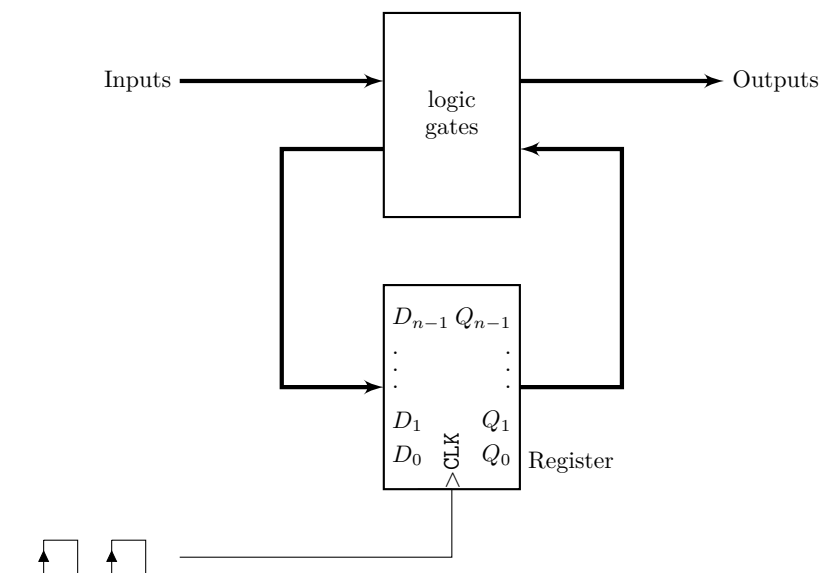


- **EEPROM (electronically erasable PROM)** is EPROM that can be erased electronically by the programmer (by applying high electric fields/voltages).

## 13.7 State Machines with Memory

Before, in Section 13.5, we covered the basic scheme of sequential logic, reproduced below.

<sup>4</sup>PIC16C57-/JW-S1, data sheet: <http://ww1.microchip.com/downloads/en/devicedoc/30453d.pdf>



Again, the basic idea is to use a register to hold logic values, the outputs of which are transformed and fed back to the register inputs via a logic-array block. Here we will discuss implementing the logic-array block in a very general way by replacing it with memory, either RAM or ROM. In the ROM case, the state machine is suited for fixed operation (e.g., as a counter), whereas with RAM we have the possibility that the state machine can adapt to input and even reprogram itself.

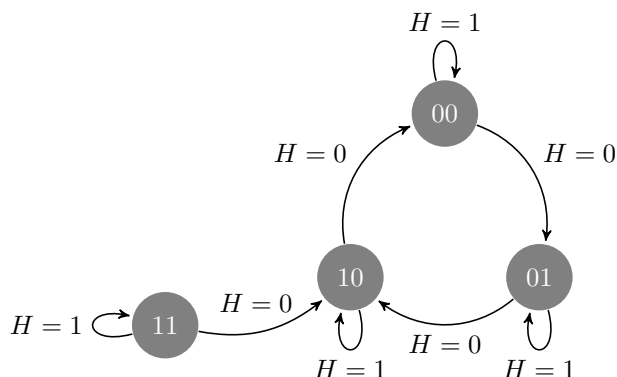
The general idea for implementing state machines with memory is:

- We will connect the register outputs ( $Q$ 's), representing the *present* state of the machine, to the memory address lines (inputs).
- We will connect the register inputs ( $D$ 's), representing the *future* state of the machine, to the memory data lines (outputs).
- Any external inputs (needed for the state machine to react to anything external), correspond to extra memory address lines.
- Any external outputs correspond to extra memory data lines.

It's easiest to see how this works in an example.

### 13.7.1 Example: Divide-by-3-With-Hold Counter

As an example of a memory-driven state machine, consider the divide-by-3 counter with a hold input  $H$ , whose state diagram we considered before in Section 13.5.2.

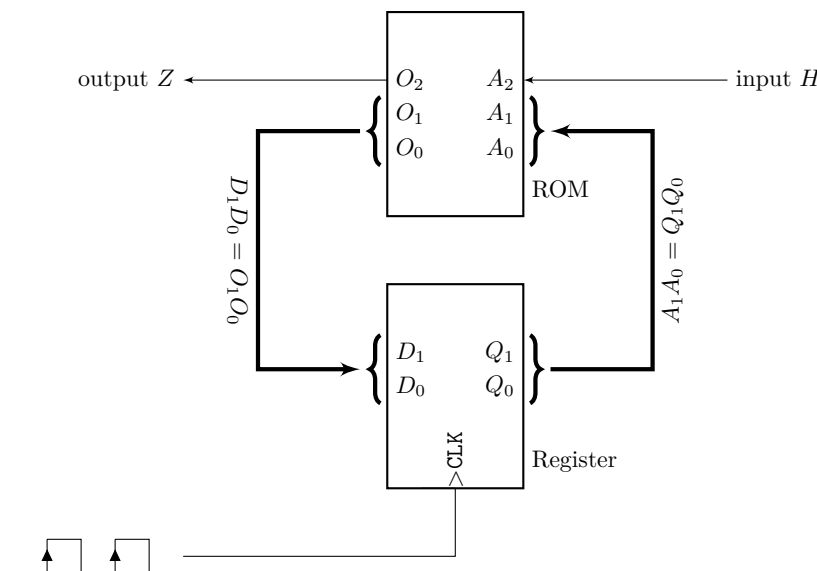


In addition, as an example output bit, we will define a “zero” output bit  $Z$  to be 1 if the counter’s output bit is 00, and is 0 otherwise.

The truth table is as follows.

(control)	(present state)		(future state)		(output)
$H$	$Q_1$	$Q_0$	$D_1$	$D_0$	$Z$
0	0	0	0	1	1
0	0	1	1	0	0
0	1	0	0	0	0
0	1	1	0	0	0
1	0	0	0	0	1
1	0	1	0	1	0
1	1	0	1	0	0
1	1	1	1	1	0

In this system, there will be 3 address bits ( $HQ_1Q_2$ ), for 8 memory slots, and 3 data bits ( $D_1D_0Z$ ), so the size of the memory is  $8 \times 3 = 24$  bits. The circuit to implement this is shown below.



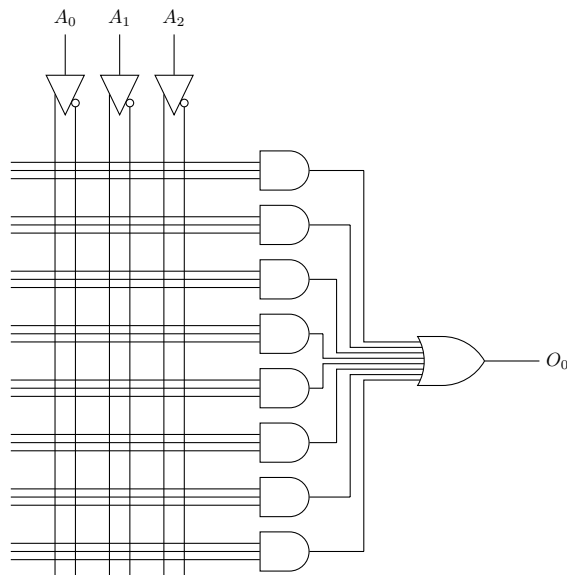
All the logic in the truth table, of course, must be programmed into the ROM. Thus, for example, in address  $A_2A_1A_0 = 000$ , corresponding to  $HQ_1Q_0 = 000$  in the truth table, we simply program in the value  $D_1D_0Z = 011$ , and so on for the rest of the 8 total memory locations.

### 13.7.2 General Considerations: Towards a Microprocessor

A few general remarks are in order. First, if there are no input bits, then basically we have some kind of counter (i.e., something that cycles through finitely many states, possibly with some outputs that are a boolean function of the state bits). If there is a single input bit, then it chooses between 2 possible actions (like the hold/count counter). If there are  $N$  input bits, then there are  $2^N$  possible operations. This grows quickly with the number of bits: for 8 inputs there are already 256 operations. We can also store sequences of input bits, e.g., in RAM, for effectively many more different possible operations for the same number of input bits (i.e., there could be a 1-bit serial input to a state machine that controls many different possible actions by using different sequences of input bits). These stored bit sequences in RAM correspond to a “program,” with 256 “instructions” in this 8-bit example, for a simple realization of a microprocessor. The input/output lines are connected to the data/address busses, which are also connected to input/output devices or interfaces. Real microprocessors often have more specific functionality, sophisticated instruction sets (with instructions that can take multiple clock cycles to complete), and have registers organized in more sophisticated ways than we have indicated here, but we still have the essence of a microprocessor.

### 13.7.3 Programmable ROM as Logic

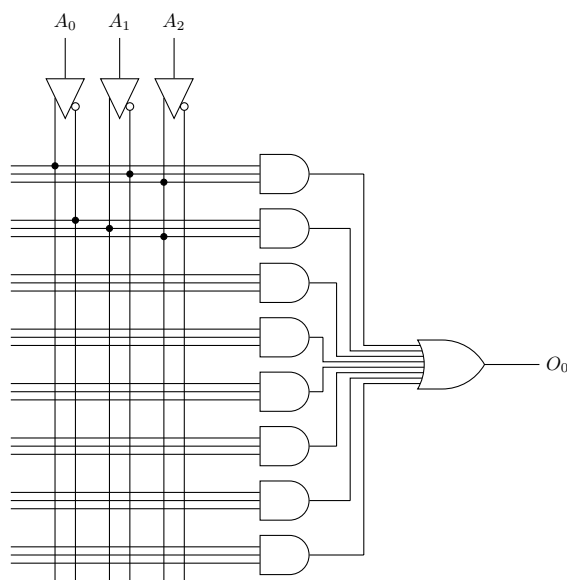
We have talked about replacing logic blocks with ROM, but we can think of ROM itself as a general array of gates, particularly the various forms of PROM. For example, consider the  $8 \times 1$ -bit memory shown below.



The main grid of lines is intended as a configurable grid, where we can make whatever connections we like. Then there are two ways to think of this. First, because the inputs go first into **AND** gates and then into an **OR** gate, we can realize Boolean-logic expressions if they are **sum-of-product** expressions. For example, we can realize

$$O_0 = A_0 \bar{A}_1 A_2 + \bar{A}_0 A_1 A_2 \quad (13.2)$$

by making the connections shown below.



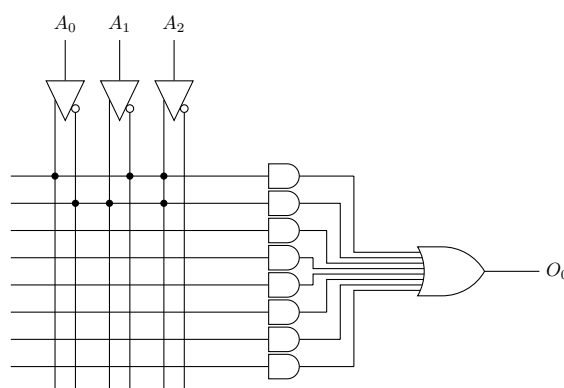
The default of the unconnected inputs is 0.

The other way to think of this particular example, is that this is a memory that stores the value 1 in the addresses  $A_0A_1A_2 = 101$  and  $A_0A_1A_2 = 011$ , and 0 in all the others. Thus we can store *any* value in any location by making proper connections for all locations that store a 1. This is why we need 8 **AND** gates

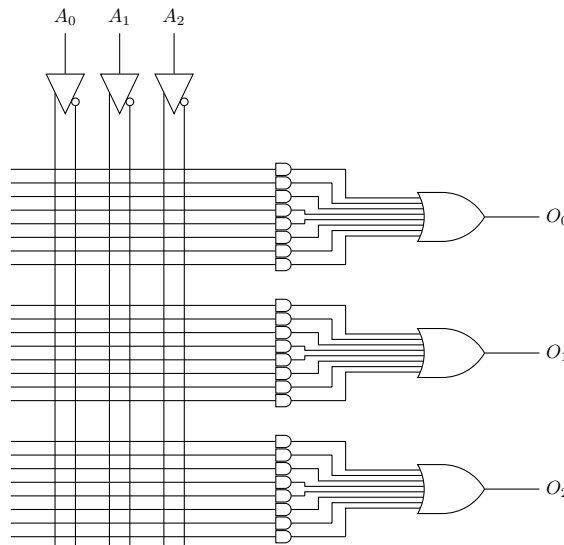
in this example—one to “recognize” each possible set of inputs, if needed. Of course, they are not all used unless all the stored bits are 1. In fact you don’t really need all of these AND gates. For example, suppose you only use 4 instead of 8 AND gates. This means that you can’t program any more than 4 1’s into memory. However, this is easy to handle: just program the *complement* of the desired logic expression, and invert the result with a programmable NOT gate at the output.

In (“write-once”) PROM, these connections are all made at the factory, and the undesired connections are “burned” away by flowing high current through the programmable fuses at each connection points.

Diagrams like this get to be complicated for larger memories, so they are often abbreviated by “collapsing” all the input lines for a given AND gate, as illustrated below. This circuit realizes the same example expression as the previous one.



In this way, we can draw out more complicated memories, like this  $3 \times 3$  memory.



### 13.7.4 Programmable Logic Devices

**Programmable logic devices (PLDs)** are chips that contain a register *and* PROM-type logic arrays like the ones we have shown above. Generally speaking, they do not have sufficient gates to implement *arbitrary* logic combinations (i.e., they generally have fewer than necessary AND gates), but they have enough to program a wide range of logic possibilities. Example of simple, but currently available, PLDs are the 22V10 **SPLD**, or **simple PLD** (e.g., the ATF22V10C from Atmel), with 12 dedicated input pins, 10 pins configurable as inputs or outputs, 10 D-type flip-flops, and a gate array (10 OR gates, with 10-16 AND gates feeding each OR gate). The flip-flop outputs can be connected to their corresponding output pins, or the

flip-flops can be bypassed altogether for non-registered outputs. A more powerful example is the ATF750C **CPLD**, or **complex PLD** from Atmel, which has the same pin configuration, but provides more gates, and 10 extra “internal” flip-flops that can be used as internal register variables that are not connected directly to outputs.

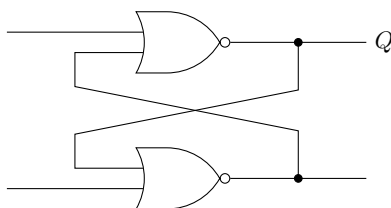
For a good example of programming a PLD as a state machine, see Appendix B, where two successive-approximation registers (Section 16.2.2) are implemented in two different PLDs.

## 13.8 Circuit Practice

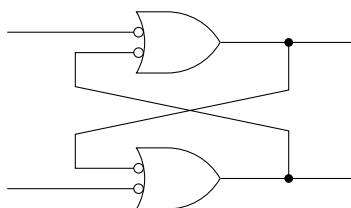
### 13.8.1 Basic Flip-Flops

For circuit practice, go through these three flip-flops.

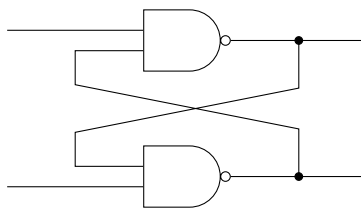
(a) First, label the inputs and remaining output. This is the first one we did, but do this without peeking! Think through the whole truth table.



(b) Work out the equivalence of the (a) circuit to this one (i.e., label the inputs and outputs). Try **not** to use a truth table to do this, use a logic theorem to connect these.

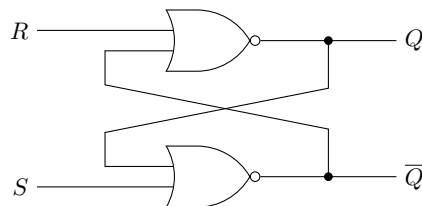


(c) Work out the equivalence of the (b) circuit to this one (i.e., label the inputs and outputs). Try **not** to use a truth table to do this, use a logic theorem to connect these.

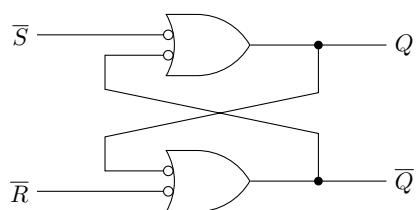


**Solution.**

(a) The labeled version is:

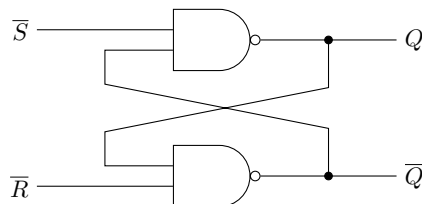


(b) This is basically the same circuit, but with all inputs and outputs of gates negated. Thus, the flip-flop inputs and outputs are similarly negated. The labeled version is:



Note that since we kept  $Q$  in the same spot, we had to swap the inputs as well.

(c) Using  $\overline{A + B} = \overline{A}\overline{B}$ , we simply change the negated-input OR gates to AND gates. The labeled version is:



### 13.8.2 Pulse-Area Stabilizer

We have seen the circuit on the next page before in Section 7.12.5. To review is designed to work as follows. To take photographs with a laser pulse, it is desirable to have the same exposure from each pulse. But the intensity of the laser drifts. Rather than try to stabilize the intensity of the laser, we can compensate for the drift by changing the *duration* of each laser pulse to compensate. By making the **pulse area** or integrated energy of each pulse the same, the photographs have exactly the same exposure, independent of the laser intensity.

Try to trace through the following features in the *digital* part of the circuit.

1. At the beginning of the pulse, the “pulse start trigger” input drives the CLK input of the D-type flip-flop (IC4a). Since the  $D$  input is tied HIGH, the  $Q$  output goes HIGH on the rising edge of this pulse, defining the beginning of the laser pulse.
2. The  $Q$  output is OR'd with the “pulse start trigger,” mainly to allow this input to override the flip-flop state, in case we want the laser to be on continuously for diagnostic purposes.
3. The pulse is finished when the integrated pulse area triggers the comparator IC3, when the inverting output goes low, triggering the  $\overline{\text{CLK}}$  input of the flip-flop (which triggers on the falling pulse-edge).
4. The  $\overline{Q}$  output drives the MOSFET to reset the integrator after the pulse is finished, until the next pulse starts.
5. The propagation delay of IC5a is matched by the propagation delay of the OR gate on the  $\overline{Q}$  output (IC5b). This is not critical, but we have extra OR gates around anyway, and it ensures an accurate  $t = 0$  for the integration.





### 13.8.3 Circuit Practice: Backwards Divide-by-3 Counter

As a simple variation on the state-machine example of a divide-by-3 counter from Section section:div3counter, design a *backwards* divide-by-3 counter ( $10 \rightarrow 01 \rightarrow 00$ , wash, rinse, repeat). Show the circuit and state diagram (and handle the excluded state).

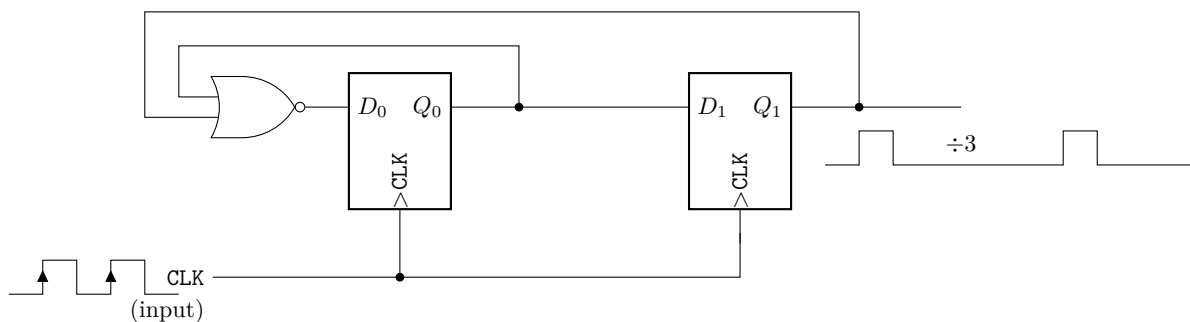
**Solution.** The counting sequence is as shown below, again with  $Q_0$  as the MSB.

$Q_0$	$Q_1$	$D_0$	$D_1$
1	0	0	1
0	1	0	0
0	0	1	0

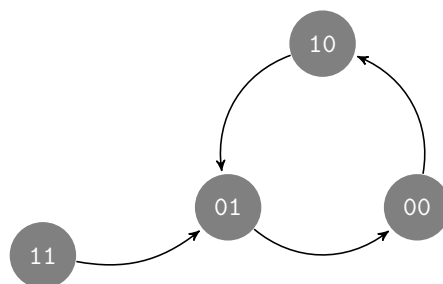
This sequence is realized by the logic

$$D_0 = \overline{Q_0 + Q_1}, \quad D_1 = Q_1, \quad (13.3)$$

which is just a reversed version of the forward counter. Thus, the circuit implementing the counter is shown below.



The excluded state is 11, which gets mapped to 01 by the logic here, which is no problem. Thus, the state diagram is as shown below.



### 13.8.4 Memory: RAM vs. ROM

Explain the following statement, about connecting memory to a CPU:

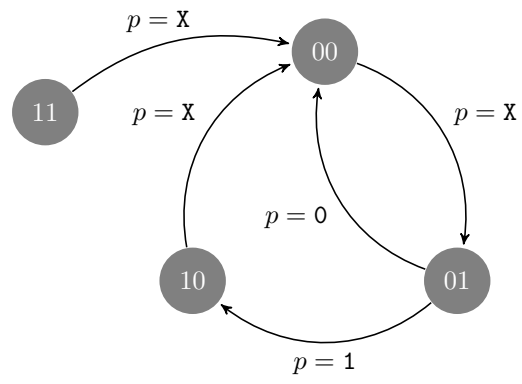
With RAM you can scramble the address lines in any order; the same is true of the address lines.

With ROM, you can't!

### 13.8.5 Circuit Practice: Divide-by-2-or-3 Counter

As practice with state diagrams, draw the state diagram for a divide-by-2-or-3 counter. That is, the counter counts differently based on an input bit  $p$ , and the counter counts  $00 \rightarrow 01 \rightarrow$  repeat if  $p = 0$ , and  $00 \rightarrow 01 \rightarrow 10 \rightarrow$  repeat if  $p = 1$ . Make sure to handle the excluded state.

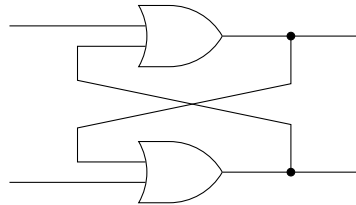
**Solution.** The diagram is sketched below.



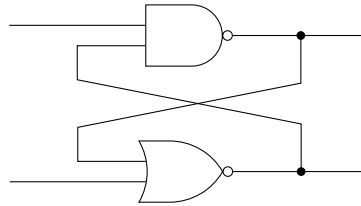
## 13.9 Exercises

### Problem 13.1

- (a) What is the essential property of a flip-flop?
- (b) Does the circuit below behave as a flip-flop? If so, label the inputs and outputs in flip-flop notation. If not, explain why not.



- (c) Does the circuit below behave as a flip-flop? If so, label the inputs and outputs in flip-flop notation. If not, explain why not.



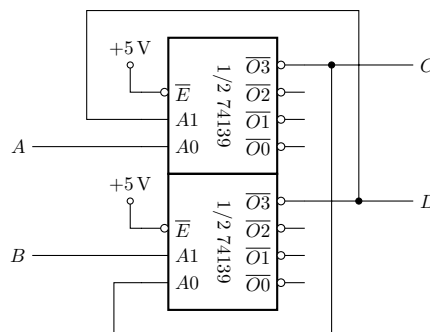
### Problem 13.2

Show how to make an RS flip-flop (with both normal and inverted outputs) using an AND gate, an OR gate, and a NOT gate. Analyze your circuit to show that this behaves as a flip-flop. What is the “bad” input state?

Note that the NOT gate isn’t strictly necessary, but you can use it to make the flip-flop better emulate the “usual” ones. (First try hooking up the AND and OR as in the usual flip-flop configuration, and then it should be more obvious where the INV should go.)

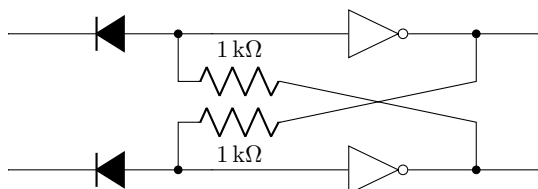
### Problem 13.3

Recall that the 74139 is a 2-bit, 4-output decoder/DEMUX, 2 per package, with 1 inverting enable per decoder. Does the circuit below behave as an SR-type flip-flop (with respect to the labeled inputs/outputs  $A$ ,  $B$ ,  $C$ , and  $D$ )? If so, label the inputs and outputs in flip-flop notation ( $S$ ,  $R$ ,  $Q$ ,  $\overline{Q}$ ). If not, explain how to change the circuit wiring to make it operate as a flip-flop, and label the inputs and outputs in flip-flop notation.



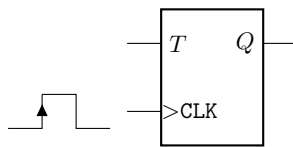
**Problem 13.4**

Does the circuit below behave as a flip-flop? If so, label the inputs and outputs in flip-flop notation, and discuss the memory and “bad” states. If not, explain why not. (You should assume **LOW** and **HIGH** logic levels of 0 and +5 V, respectively. You may also assume no current flows into the gate inputs; the resistors are intended to work only in conjunction with the diodes. You may want to begin by analyzing the circuit with nothing connected to the inputs.)

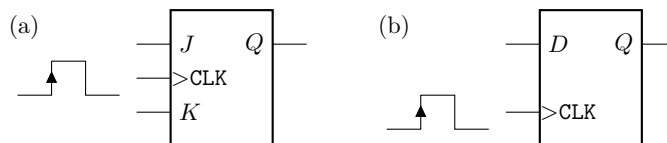
**Problem 13.5**

A **T flip-flop** has a single (“ $T$ ”) input, which causes the output  $Q$  to toggle if the  $T = \text{HIGH}$ , and to hold if  $T = \text{LOW}$ . The truth table and schematic diagram are shown below ( $Q_n$  is the output state  $Q$  after the  $n$ th clock pulse).

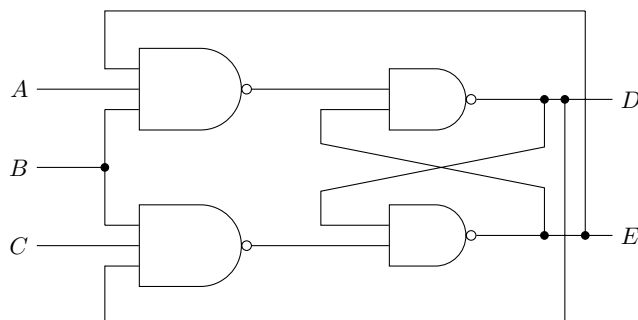
$T$	$Q_n$	$Q_{n+1}$
0	0	0
0	1	1
1	0	1
1	1	0



- (a) Show how to connect the JK flip-flop below as a T flip-flop.  
 (b) Show how to connect the D flip-flop below as a T flip-flop. You will need to use one additional gate.

**Problem 13.6**

The circuit below realizes a circuit from this chapter. What is it? (Be specific.) Label the inputs and outputs to indicate the function and operation of the circuit, being as specific as possible.

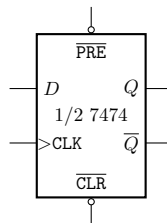


**Problem 13.7**

- (a) Show how to connect 3 JK flip-flops to make a 3-bit ripple counter, and then show how to add a single 2-input NAND gate to change this to a **modulo-5** counter—that is, the counter resets to zero when the output reaches 5, and continues counting.
- (b) Show how to design a (ripple) counter that counts from 0 to 5 and then stops. Your circuit should include a **RESET** input that resets the counter back to zero (and then continue counting) after a negative pulse.
- (c) Show how to design a circuit that passes only 5 (positive) pulses and blocks subsequent other pulses, after a negative reset pulse that “arms” the circuit.

**Problem 13.8**

- (a) Show how to connect 3 flip-flops to make an asynchronous (ripple), 3-bit **down** counter. To be specific, use 7474 flip-flops (dual D-type, positive-edge-triggered, with complementary outputs and jam preset and clear), whose connections are shown below.

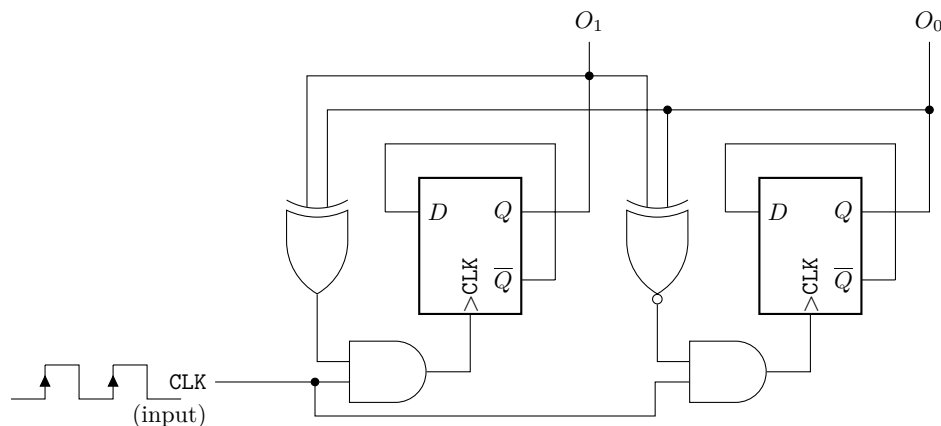


Recall that for the jam inputs, “preset” is the same thing as “set,” and “clear” is the same thing as “reset.”

- (b) Show how to realize an **asynchronous, divide-by-5 down counter**, made from the same D-type flip flops. That is, your counter should count 4, 3, 2, 1, 0, 4, 3, 2, 1, 0, ...

**Problem 13.9**

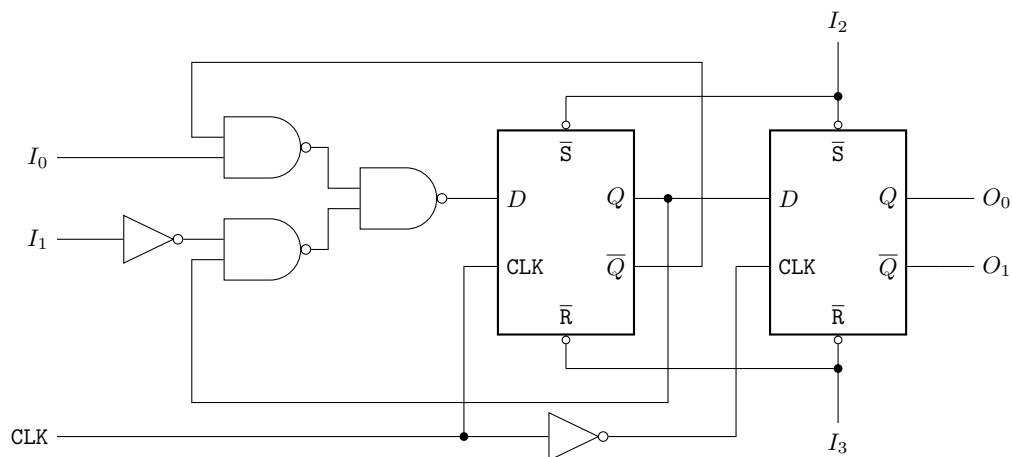
Consider the asynchronous counter circuit shown below.



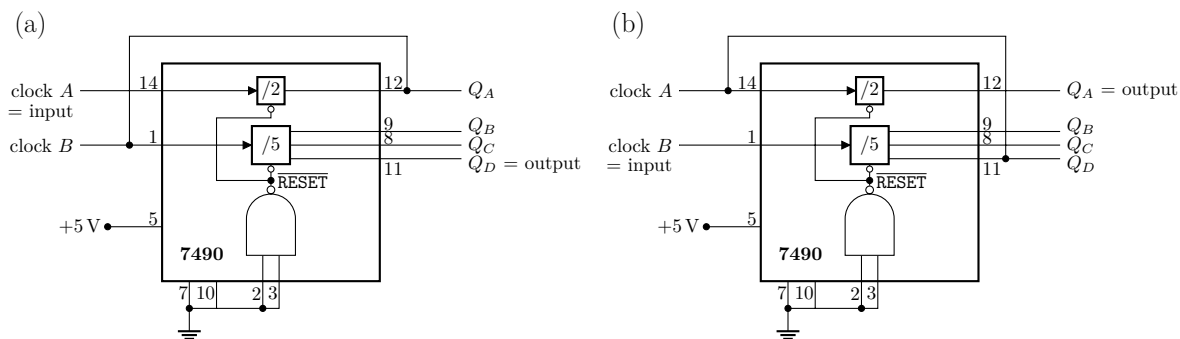
- (a) Why is this circuit “asynchronous” and not “synchronous”?
- (b) Assuming the counter starts in the state  $O_1O_0 = 00$ , give the truth table for the counter output on subsequent clock cycles.
- (c) Can this counter get “stuck” in any particular state? Why or why not?
- (d) What kind of counter is this?

**Problem 13.10**

Describe *precisely and completely but concisely* the operation of the circuit below, including all inputs  $I_0$ – $I_3$  and outputs  $O_0$ – $O_1$ .

**Problem 13.11**

The 7490 is a decade counter, with clock inputs that trigger on falling edges. Below are two ways to connect the 7490 as a divide-by-10 counter. However, they are not equivalent. What, specifically, is the difference in the output waveforms? (Be *quantitative*.) Also, remember  $Q_B$  is the least significant bit of the divide-by-5 subcounter.

**Problem 13.12**

Design a synchronous 2-bit UP/DOWN counter: It has a clock input, and a control input ( $U/\overline{D}$ ); the outputs are the two flip-flop outputs  $Q_1$  and  $Q_2$ . If  $U/\overline{D}$  is HIGH, it goes through a normal binary counting sequence; if LOW, it counts backward— $Q_2Q_1 = 00, 11, 10, 01, 00, \dots$ <sup>5</sup>

**Problem 13.13**

Design a **synchronous**, 3-bit Fibonacci counter (i.e., count through the Fibonacci numbers 0, 1, 2, 3, 5, and then repeat). Use three flip-flops (of the 7474 type, as shown above) and whatever gates you

<sup>5</sup>Paul Horowitz and Winfield Hill, *The Art of Electronics*, 2nd ed. (Cambridge, 1989), Exercise 8.25 (ISBN: 0521370957).

like. Be sure to show a state diagram, *including* all excluded states. Can your counter get “stuck” in any excluded state?

### Problem 13.14

Design a synchronous divide-by-3 circuit using two JK flip-flops, but without any gates or inverters. One hint: When you construct the table of required  $J_1$ ,  $K_1$  and  $J_2$ ,  $K_2$  inputs, keep in mind that there are two possibilities for  $J$ ,  $K$  at each point. For instance, if a flip-flop output needs to go from 0 to 1,  $J, K = 1, X$  ( $X$  = doesn’t matter). Finally, check to see if your circuit will get stuck in the excluded state.<sup>6</sup>

### Problem 13.15

(a) (15 points) Design a **synchronous** circuit (state machine) using flip-flops and logic gates that makes a 4-bit “Knight-Rider” pattern. That is, the output counts:

0001

0010

0100

1000

0100

0010

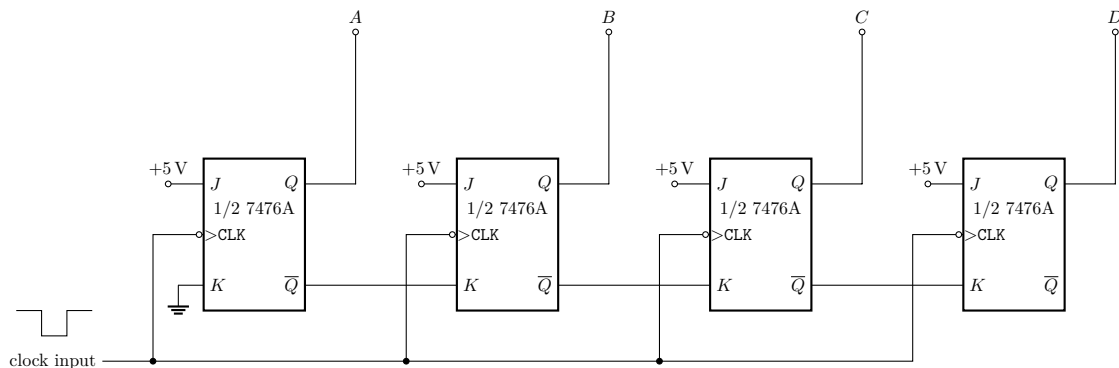
(repeat).

Don’t draw the circuit, just come up with the required logic expressions. *Hint:* there are 4 bits shown here, but you will need an extra bit to keep track of the direction.

(b) (10 points) Of course, if the Knight Rider car is driving down the street with the lights stuck in an excluded state, it would be less than awesome. Either show that your circuit doesn’t get stuck in excluded states, or show how to modify your circuit to make sure that your circuit settles into the above pattern for any starting state. Use whatever logic you like (again, no circuit, just the required logic).

### Problem 13.16

Analyze the circuit below by drawing a timing diagram for the outputs. Assume that the initial state is  $ABCD = 0111$ , and analyze the circuit output for 5 clock pulses.

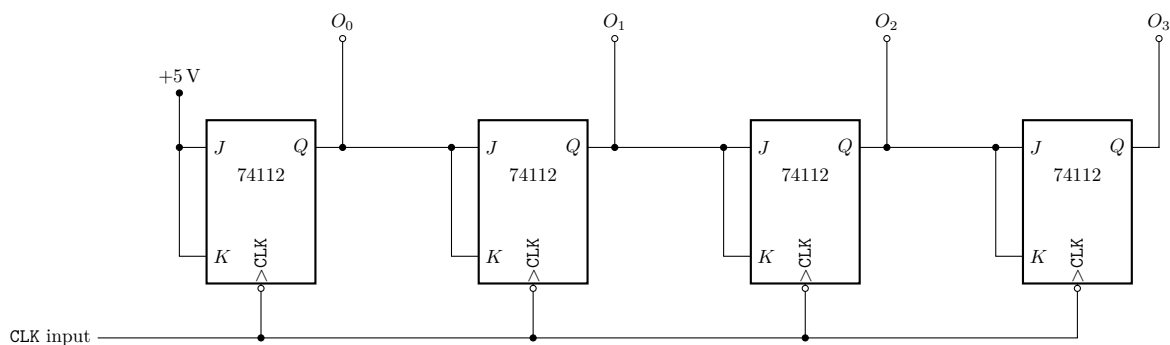


### Problem 13.17

Work out the counting sequence (starting at  $O_3O_2O_1O_0 = 0000$ ) for the synchronous counter shown below.

<sup>6</sup>Paul Horowitz and Winfield Hill, *op. cit.*, Exercise 8.24.

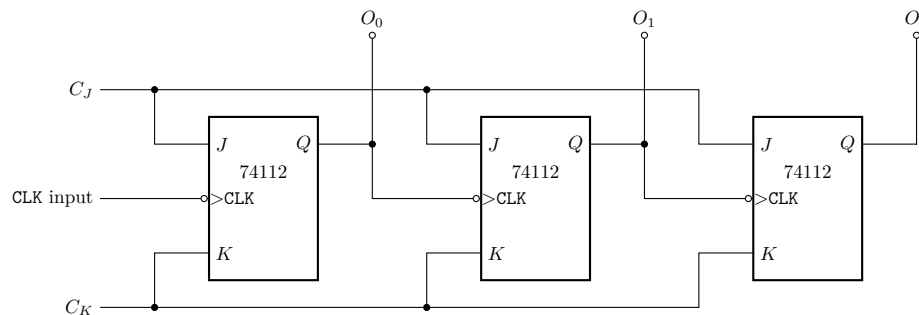


**Problem 13.18**

Design a **synchronous**, divide-by-16 counter. Just work out Boolean expressions for the required logic using the design procedure for the state machine; no need to draw a circuit schematic. (Write your output bits as  $Q_3Q_2Q_1Q_0$  from MSB to LSB.)

**Problem 13.19**

Consider the ripple-counter circuit shown below, with control inputs  $C_J$  and  $C_K$ .



What control-input combination  $C_JC_K$  is necessary for the circuit to count, and what kind of counter is it? Describe the operation of this circuit for all other control-input combinations  $C_JC_K$ .

**Problem 13.20**

- Design a **synchronous** 4-bit, binary counter, using JK flip-flops (use the same 74112's as shown in Problem 13.19), by considering the following: Think about how binary counting works, and in particular, what is the condition for a particular bit to toggle its state during the counting sequence? How do you switch between toggling or holding in a JK flip-flop?
- Design a **synchronous** 4-bit, binary counter, as in part (a), but this time using D flip-flops. Use the same hint from part (a). What kind of gate acts as a conditional inverter?

**Problem 13.21**

- Design an **asynchronous/ripple** BCD counter (i.e., counts up in binary from 0–9 and then resets). Use JK flip-flops as in Problem 1, but include (and use) jam-reset inputs  $\overline{R}$ .
- To drive a two-digit decimal display that counts from 00–99 (then resets), you would need to provide an *output* to drive the counter for the next digit. Indicate how you would do this in your circuit, and how it connects to the counter for the next digit.

**Problem 13.22**

(a) Design a **synchronous** version of the BCD counter in Problem 13.21. Again, use JK flip-flops as in Problem 1 (and no jam inputs). That is, treat this circuit using the normal design process for a state machine/synchronous counter.

*Hint.* Simplify your analysis by tying the  $J$  and  $K$  inputs together for each flip-flop (like in Problem 1).

(b) Repeat part (b) of Problem 13.21 for this circuit. But now your counter circuit will need an input as well as an output, to interface to the digits of lesser and greater significance, respectively.

**Problem 13.23**

Design a 3-bit, **synchronous** counter circuit with one control input bit  $E$  (for “even”), that works as follows. If  $E = 1$ , then the circuit counts up through the *even* numbers, and the circuit counts up through the *odd* numbers if  $E = 0$ . On clock cycles where  $E$  has just changed state, the counter should always start on 0 or 1 for  $E = 1$  or 0, respectively.

**Problem 13.24**

Design a 3-bit **synchronous** counter, with the following specifications.

- Use JK flip-flops for the three output bits; call them  $O_0$ – $O_2$ , with  $O_0$  the LSB.
- Include a 2-bit “increment” input  $I_1I_0$ , which gives the increment of the counter on each clock cycle. That is,  $I_1I_0 = 00$  means that the counter just holds;  $I_1I_0 = 01$  means the counter counts like a regular divide-by-8 counter;  $I_1I_0 = 10$  means the counter counts only odds or evens (depending on the initial state); and  $I_1I_0 = 11$  counts 0, 3, 6, 1, 4, 7, 2, 5, ... (in binary).
- Use whatever logic gates you need.
- You need only write out the logic expressions for the flip-flop inputs to make the counter work; no need to sketch a circuit diagram.

*Hint:* start by working *separately* with each possible input state.

**Problem 13.25**

The goal of this problem is to design a circuit that adds two binary numbers, one bit at a time, starting with the LSB. That is, on the first clock cycle, the circuit should take the LSBs on the input lines  $A$  and  $B$ , and produce the LSB of the result as output  $Q$ . On the second cycle, the circuit should take the second bits of the two numbers on the (same) input lines, generating the second bit of the result on the same output. The circuit should continue until the addition is complete (i.e., indefinitely). You may assume that you already have a circuit that presents the proper bits serially to  $A$  and  $B$ , one each per clock cycle, as required.

(a) Draw the state diagram for this circuit, showing explicitly the operation of the output bit and a “carry” bit, as well as how the inputs  $A$  and  $B$  drive state transitions. (You may assume that the carry bit starts in the appropriate state at the beginning of the addition.)

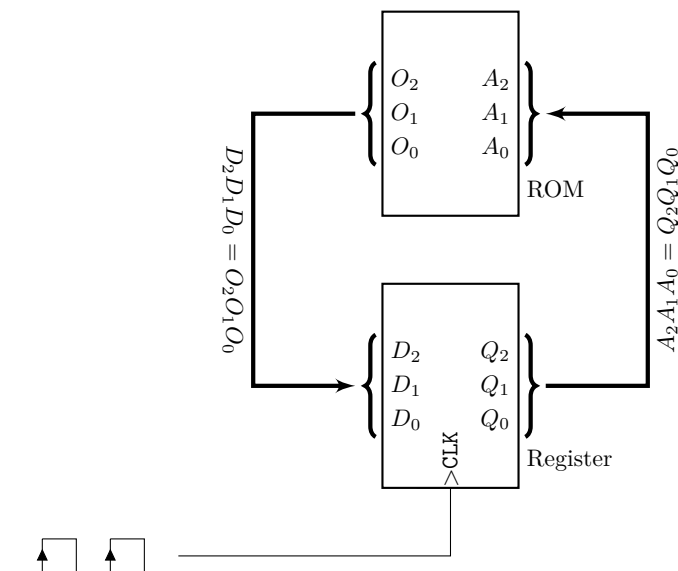
(b) Design a circuit that realizes this state machine. (Just work out the required logic, no need to draw the circuit.)

**Problem 13.26**

Design a circuit that combines 4 of the 6116 SRAM IC’s (2k×8-bit) to make a single, 4k×16-bit memory. The resulting circuit should have the same behavior as one of the original 6116’s, just with more address/data bits.

**Problem 13.27**

A 3-bit register and an  $8 \times 3$  ROM are connected as shown below.



The ROM is programmed as follows:

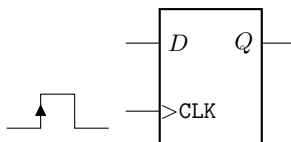
$A_2$	$A_1$	$A_0$	$O_2$	$O_1$	$O_0$
0	0	0	1	0	0
0	0	1	0	1	0
0	1	0	0	0	0
0	1	1	0	1	1
1	0	0	1	1	0
1	0	1	0	0	0
1	1	0	1	1	1
1	1	1	0	0	1

Assume that the initial state of the register is  $Q_2Q_1Q_0 = 010$ .

- Give the state of the register outputs  $Q_2Q_1Q_0$  after *each* of the next three clock pulses.
- Even after arbitrarily many clock pulses, (at least) one possible value of  $Q_2Q_1Q_0$  will never occur. Which?

**Problem 13.28**

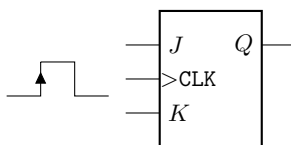
Show how to connect a T flip-flop as a D flip-flop (shown below). (The operation of the T flip-flop is defined in Problem 13.5.) Use whatever extra logic gates you need (simple two-input gates only).



*Hint:* if you don't see how to get started, treat the problem as a state machine with current state  $Q_n$  and external input  $D$ .

**Problem 13.29**

Show how to connect a T flip-flop as a JK flip-flop (shown below). Recall that the operation of the T flip-flop is defined in Problem 13.5. You can use whatever extra logic gates you need (simple two-input gates only).

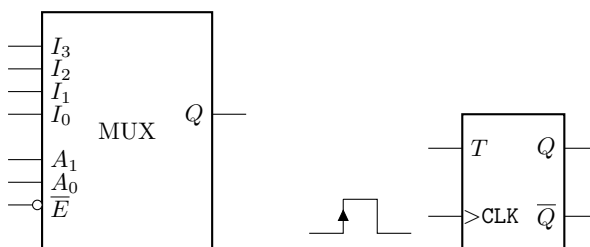


For this problem, follow the outline below.

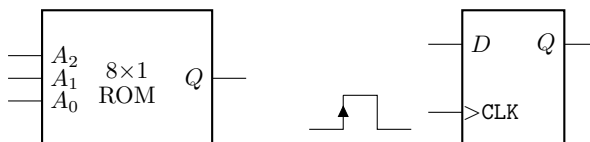
- Begin the design as a state machine by making a truth table, with columns:  $J$ ,  $K$ , and  $Q_n$ ,  $Q_{n+1}$ , and  $T$ . (The  $J$ ,  $K$ , and  $Q_n$  variables are the inputs and current state,  $Q_{n+1}$  the future state, and  $T$  the flip-flop input where you need to supply the appropriate logic function as input.)
- Make a Karnaugh map for the  $T$  variable in terms of the input variables and current state, and use it to write down a simple logic expression for  $T$ .
- Complete the design by making a schematic diagram that realizes the JK flip-flop in terms of the T.

**Problem 13.30**

Show how to connect a T flip-flop as a JK flip-flop, but using only an extra 4-input MUX, and *no extra logic gates*. The T flip-flop and MUX you should use are shown below. (The operation of the T flip-flop is defined in Problem 13.5.)

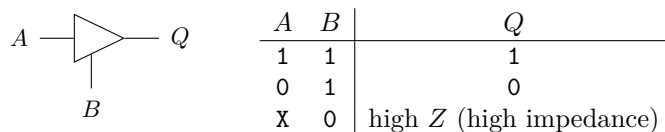
**Problem 13.31**

Show how to connect an  $8 \times 1$  ROM (shown below) as a JK flip-flop. You will need an extra D flip-flop (also shown below), but *no other logic gates*. Also specify how the ROM should be programmed to make your circuit work properly.

**Problem 13.32**

Design a  $1 \times 1$  RAM circuit according to the following specifications:

- The circuit should have  $\overline{CS}$ ,  $\overline{WE}$ , and  $\overline{RE}$  inputs that work as usual (the circuit behavior is *undefined* if both  $\overline{WE}$  and  $\overline{RE}$  are simultaneously LOW, meaning it is allowed to do anything except self-destruct).
- The circuit should also have a single I/O line for both writing data to and reading data from the memory.
- You can use as many logic gates as you like, but do *not* use any “prepackaged” flip-flops. You should also include a single 3-state buffer, as shown below (with the truth table).



### Problem 13.33

Design a push-button interface circuit, which makes use of a 1-kHz clock signal, with the following specifications.

- The circuit output is normally HIGH, and goes LOW in response to a button press.
- The switch is a simple, normally open, SPST switch. (That is, it’s a “make or break” switch which is either conducting during a button push or otherwise nonconducting.)
- The LOW output should last only one clock cycle, even if the duration of the button press is longer.
- Once the LOW output has been triggered, the button should be “locked out,” so that the circuit cannot produce any more LOW pulses for  $\sim 50$  ms. Note that since switch bounces typically occur on ms time scales, this lockout effectively debounces the switch.
- The circuit should only produce one output pulse per button press; that is, a long press (greater than 50 ms) will only produce one pulse.
- Use whatever flip-flops, logic gates, resistors, and switches you like (but not any more-complicated chips that happen to contain these); assume the existence of a 5-V power supply.
- It would be a good idea to include a state machine as part of your circuit, in which case you can skip drawing it and just specify the number and type of flip-flops, as well as the logic expressions needed to make the state machine work. But explicitly handle any excluded states.

### Problem 13.34

A typical cruise-control interface in a modern car might have four buttons (plus a pedal switch), with the following functions:

- $\overline{ON}$ : Activates the cruise-control system if it is not currently active (but does not hold the car’s speed until the  $\overline{SET}$  command). An active system is prerequisite for the system to maintain speed.
- $\overline{OFF}/\overline{CANCEL}$  (multipurpose button): If the cruise-control system is currently active, but not holding speed, this button deactivates it; all other input (except  $\overline{ON}$ ) is ignored in the deactivated state. If the cruise-control system is maintaining speed, this button stops maintaining speed (i.e.,  $\overline{CANCEL}$ ), but the system remains active until a second press.
- $\overline{PEDAL CANCEL}$ : a switch (basically equivalent to another button) in the brake pedal performs the same  $\overline{CANCEL}$  operation if the pedal is depressed (but it cannot turn the system off).
- $\overline{SET}/\overline{COAST}$  (multipurpose button): If the cruise-control system is active but not maintaining speed, this button activates the servo so the car maintains the *present* speed. If the system is already maintaining speed, a button press reduces the target speed by 1 mph. (Multiple button presses are needed to reduce the speed by more than 1 mph.) This button has no effect if the system is in the deactivated state.

- $\overline{\text{RESUME}}/\overline{\text{ACCEL}}$  (multipurpose button): If the cruise-control system is active but not maintaining speed, this button activates the servo so the car maintains the *previously held* target speed (assuming some speed was in fact previously held; otherwise the button should do nothing). If the system is already maintaining speed, a button press increases the target speed by 1 mph. Multiple button presses are needed to increase the speed by more than 1 mph. This button has no effect if the system is in the deactivated state.

The bars over the button names emphasize that the button output is normally HIGH, and active LOW (momentarily, during a button press).

Your task as a newly employed engineer of the Schmoyota Motor Corporation is to design the interface logic for the cruise-control system for next year's Schmoyota Schmius. That is, **design a state machine and its logic implementation** to take the five buttons as input, and to provide outputs to an existing cruise-control servo module (which does all the work of sensing and holding the current speed; this means you don't have to track or deal directly with the speed in your circuit). The five inputs to the servo module (all to be connected to your circuit, but to reiterate, all of the below servo functionality has been implemented already by another engineer) are as follows:

- $\overline{\text{ACTIVE}}$ : When LOW, this input activates a light on the instrument cluster to show the driver that the cruise-control system is active and ready to start maintaining speed.
- $\overline{\text{MAINTAIN}}$ : This input activates the speed servo, causing the car to hold constant speed. Activating this input also causes the target speed (in mph) to be displayed on the instrument cluster.
- $\text{SET}/\overline{\text{RESUME}}$ : Right after the falling edge of  $\overline{\text{MAINTAIN}}$ , when the car begins to hold speed, the system needs to know if it should hold the *current* speed or the *last held* speed. If HIGH, this input tells the servo to hold the current speed, while if LOW, the servo should set the target speed to the previously held target. **Do not** tell the servo to hold a previous target unless the system has in fact held some previous target speed; otherwise the system will attempt to hold some undefined (random) speed, which may not be very nice for the driver!
- $\overline{\text{SPEEDUP}}$ : Pulse this input LOW to increment the target speed by 1 mph (pulse only while  $\overline{\text{MAINTAIN}} = \text{LOW}$ ).
- $\overline{\text{SLOWDOWN}}$ : Pulse this input LOW to decrement the target speed by 1 mph (pulse only while  $\overline{\text{MAINTAIN}} = \text{LOW}$ ).

For your state machine, you should use 4 D-type flip-flops: three will drive the servo inputs  $\overline{\text{ACTIVE}}$ ,  $\overline{\text{MAINTAIN}}$ , and  $\text{SET}/\overline{\text{RESUME}}$ , plus you should use another flip-flop with output  $\overline{\text{ALREADY}}$ , which is LOW if (and only if) a speed has already been held since the system was activated. You can take  $\overline{\text{SPEEDUP}}$  and  $\overline{\text{SLOWDOWN}}$  to be extra outputs of your state machine.

Assume that a clock source has already been provided (at some reasonably high frequency so the interface seems to respond instantaneously, as far as the driver is concerned). Also assume that the buttons are connected via the circuit of Problem 13.33, so that a button press coincides with exactly one clock cycle, and bounce is not an issue.

(a) Write down a truth table defining the operation of your state machine according to the above specifications. Your table should include the current and future states of the four flip-flops, plus the five button inputs. Make use of the  $X = \text{"either"}$  notation to avoid writing out all  $2^9$  individual truth-table entries!

(b) Derive logic expressions for the inputs of the four flip-flops, as well as the  $\overline{\text{SPEEDUP}}$  and  $\overline{\text{SLOWDOWN}}$  inputs to the servo. (No need to draw the circuit.)

### Problem 13.35

In the circuit of Problem 13.34, which of the buttons would produce undesired behavior if the button were **not** debounced? Briefly *explain*.

(Note that you don't need to solve that problem to solve this one.)



# Chapter 14

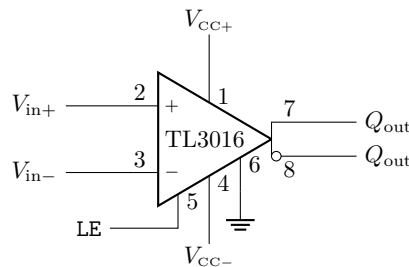
## Comparators

### 14.1 Overview and Review

We talked before about comparators in the context of analog electronics in Section 7.9. There, comparators were variations on the basic op-amp, which compares two analog input voltages, and outputs something like a 1 or 0, depending on the comparison. Here, we will review and extend our discussion of comparators, because in the context of digital electronics, comparators are the fundamental way to interface analog to digital circuits. For example, they form the fundamental building block of the **analog-to-digital converter (ADC)**. They also represent a basic method for generating logic pulses (generally, the pulse timing is based on an analog signal, such as an RC decay). Also, comparators allow you to do “level translation” between different logic types (e.g., interfacing high-voltage logic to TTL).

#### 14.1.1 Example: TL3016

As an example of a comparator that is really optimized for driving logic circuits, consider the TL3016 (or LT3016) fast comparator (“fast” here means a 7.6-ns propagation delay). This comparator and its pin connections are shown schematically below.



This comparator has both regular and inverted outputs, and the outputs are TTL-compatible. To summarize the regular operation:

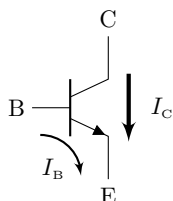
- if  $V_{in+} > V_{in-}$ :  $Q = \text{HIGH}$ , and  $\overline{Q} = \text{LOW}$
- if  $V_{in+} < V_{in-}$ :  $Q = \text{LOW}$ , and  $\overline{Q} = \text{HIGH}$
- here, “HIGH” means nominally around +5 V (actually, about +3.8 V)
- and “LOW” means nominally around 0 V (actually, about +0.6 V)
- for proper operation, the inputs should be in the range of the supply voltages; the positive supply  $V_{CC+}$  is a +5 V, while the negative supply  $V_{CC-}$  is either 0 V or -5 V
- the “LE” pin is a “latch enable”, which latches the output when held HIGH



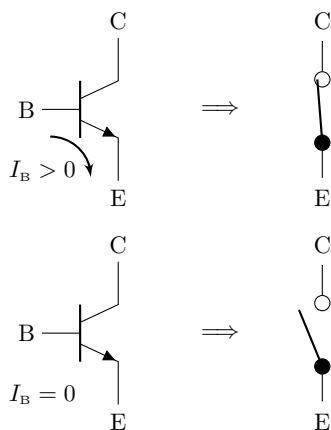
## 14.2 Open-Collector Output

A common configuration for comparator outputs (and logic-gate outputs, too) is the **open-collector output**. We discussed this before in Section 7.9, but we will review the basic idea here.

First, we should review the basic switch-type operation of the bipolar transistor. The transistor acts as a switch for current, based on another current. The two important currents are  $I_B$  from the base to the emitter, and  $I_C$  from the collector to the emitter.

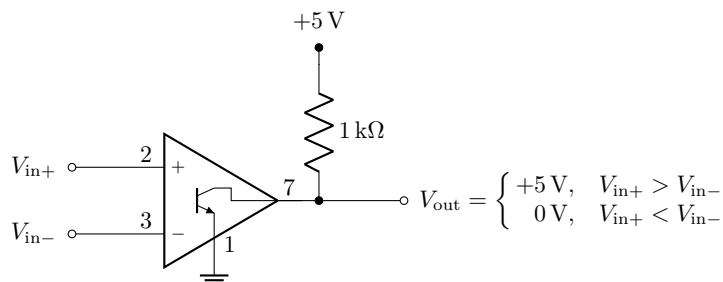


The base current  $I_B$  acts as the control current, and  $I_C$  is the current to be switched. Simplistically, if there is some current  $I_B$ , then  $I_C$  can flow, so the C–E path acts as a closed switch.



However, if  $I_B = 0$ , then the C–E path acts as an open switch. There are some extra voltage drops to consider here, but as in RTL logic (Section 11.3), this simple model is sufficient to understand open-collector outputs.

An inexpensive and popular comparator with open-collector output is the LM311 (also LF311, henceforth just the “311”). Schematically, this comparator is shown below, with a typical “pull-up resistor” to +5 V on the output.

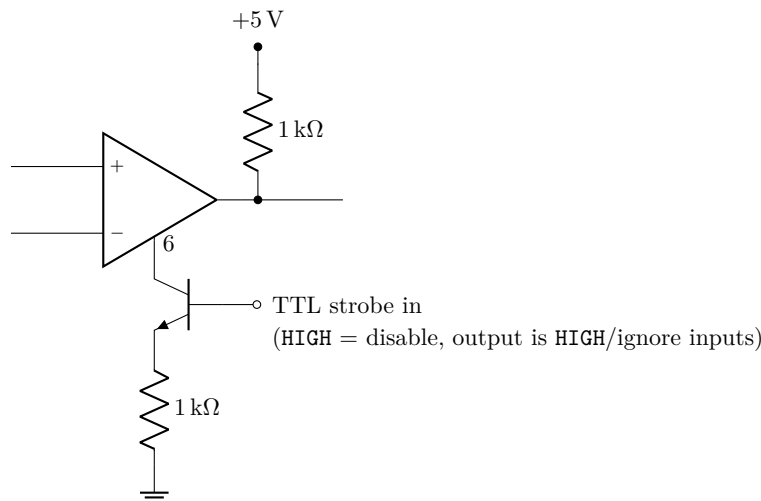


Not shown here are power-supply connections (to supplies of up to  $\pm 18$  V). The inputs must stay within the supply-voltage range. To summarize the operation here:

- if  $V_{in+} > V_{in-}$ : the transistor is **OFF**, and  $V_{out} = +5$  V
- if  $V_{in+} < V_{in-}$ : the transistor is **ON**, and  $V_{out} = 0$  V (actually, about 0.2 V or higher)

The point of the open-collector output is its flexibility: it's not restricted to particular voltages like the outputs of normal TTL gates. For example, the 311 can drive loads up to 40 V and 50 mA, so it can directly control LEDs, lamps, relays, and so on. By contrast, a TTL-gate outputs would need a buffer transistor.

Another feature of the 311 is a TTL strobe input. A typical connection for this pin is shown below.

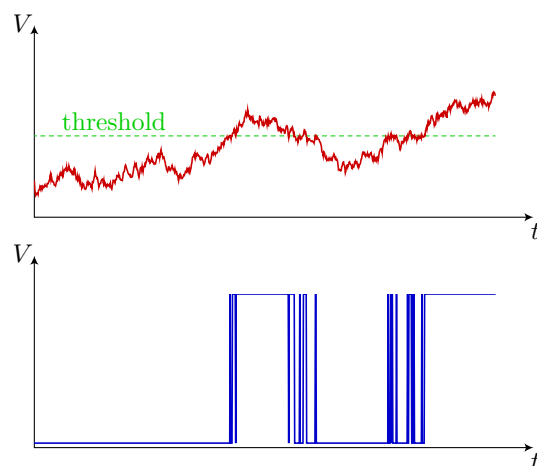


If the strobe input is held HIGH, then the output is disabled: the output is HIGH (i.e., the output transistor opens), and the comparator ignores the inputs. This is useful for “gated” operation, if the comparator should only trigger during some time interval or some condition determined by other logic.

### 14.3 Schmitt Trigger

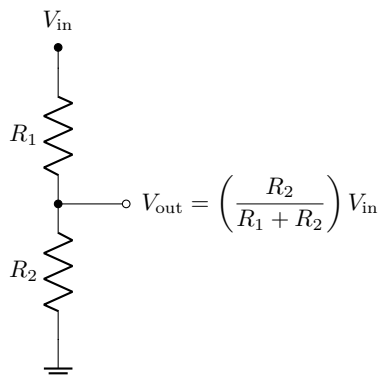
We also discussed the **Schmitt trigger** in the context of analog circuits before in Section 7.9.1. We will review the basic idea again here, since these are so important in digital applications that many logic devices have integrated Schmitt triggers.

The motivation for the Schmitt trigger comes from noisy inputs signals. Typically, these are slowly changing analog signals, but even “digital” signals are fundamentally analog, and are thus similarly susceptible to noise. As illustrated below, a noisy signal that is rising and falling can cause many spurious transitions while crossing the threshold.

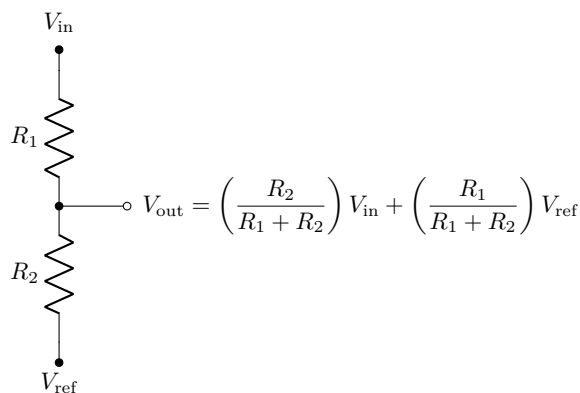


Ideally, with a noise-free signal, the output (shown in the bottom plot) would have one up transition, then later a down transition, and then an up transition again, instead of the many clustered around the three “ideal” trigger times.

To analyze the fix for this, we will briefly review the voltage divider (see Section 1.3.3 and Problem 1.2). Given a series pair of resistors supplied by  $V_{\text{in}}$ , the output voltage at the tap point is given in terms of the “fractional resistance” at the tap point.



If there are two voltages at either end of the resistor pair, the output voltage is a linear combination of the two voltages, given by the fractional resistances:

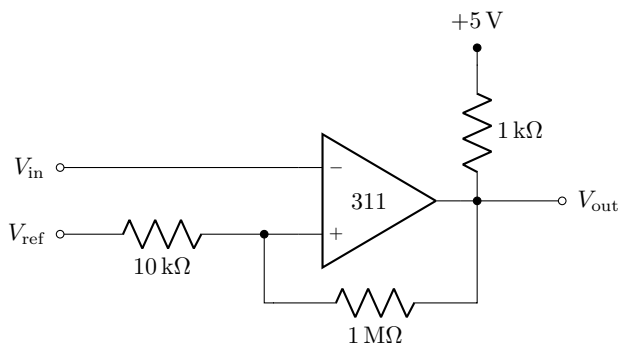


We can deduce the second result from the first by subtracting  $V_{\text{ref}}$  from all voltages, and applying the first (grounded) result, to give

$$V_{\text{out}} - V_{\text{ref}} = \left( \frac{R_2}{R_1 + R_2} \right) (V_{\text{in}} - V_{\text{ref}}). \quad (14.1)$$

Rearranging gives the result in the figure. *Note:* you should *memorize* both of these voltage-divider formulas, and be able to quickly come up with resistor combinations that divide a voltage by 2, 3, etc.

Now consider the following circuit, which is a comparator with two added resistors. One resistor is in series with the “trigger voltage”  $V_{\text{ref}}$ , and another, large resistor ties the output to the noninverting input.



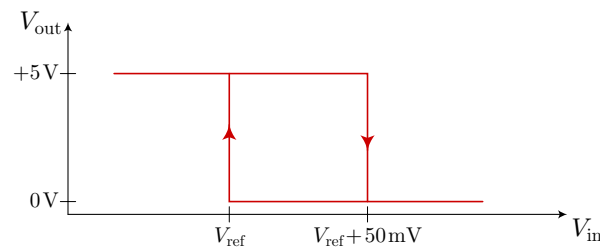
The extra resistors introduce hysteresis into the operation of the comparator as follows.

- if  $V_{\text{out}} = 0$ ,  $V_+$  (the voltage at the noninverting input) is  $0.99 V_{\text{ref}}$ , using the first voltage-divider formula.
- if  $V_{\text{out}} = +5 \text{ V}$ ,  $V_+ = 0.99 V_{\text{ref}} + 0.01 \cdot 5 \text{ V}$ , using the second voltage-divider formula.

Thus, the reference voltage changes by about 50 mV, depending on the output. Notice that

- if  $V_{\text{in}} = \text{HIGH}$ , the trigger point is *lower*
- if  $V_{\text{in}} = \text{LOW}$ , the trigger point is *higher*

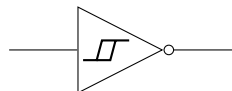
So  $V_{\text{in}}$  “repels” the trigger point, which gives immunity to noise. The hysteresis of the Schmitt trigger is sketched in the output-response plot below.



Note that the Schmitt trigger is bistable in the 50-mV-wide region near the nominal trigger voltage  $V_{\text{ref}}$ . Thus, once the Schmitt trigger makes a transition, noise of less than around 50 mV will not cause another spurious transition.

Note that the configuration above is inverting, since  $V_{\text{in}}$  goes into the inverting input (i.e., large  $V_{\text{in}}$  means the output is LOW). For a noninverting comparator, you can swap the  $V_{\text{in}}$  and  $V_{\text{ref}}$  labels, but note that the input is no longer well-isolated from the input, which may be a problem if the input source has high impedance.

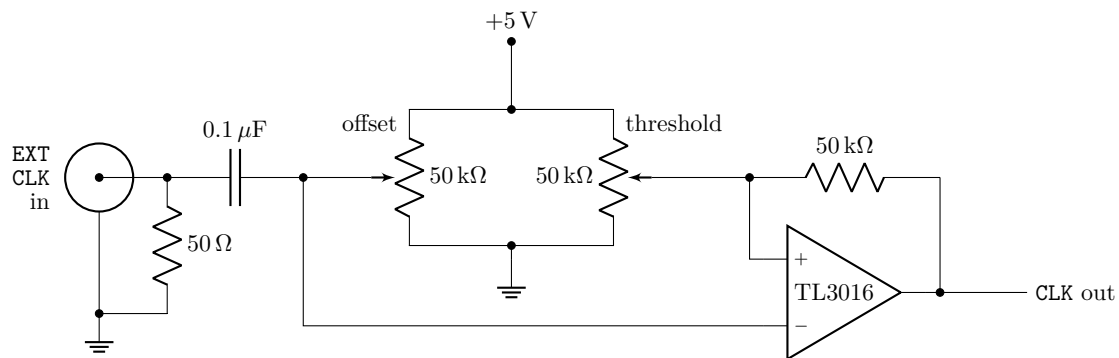
In digital circuits, many gates and logic devices are available with Schmitt-trigger inputs. An example is the 7414, a hex Schmitt-trigger inverter. This is shown schematically below.



Note the Schmitt-trigger symbol on the gate, which suggests the hysteresis curve. Schmitt-trigger-input gates are good for signals without well-defined edges (like a sinusoidal input), or signals from an external source with a long cable run, which could be susceptible to noise pickup.

### 14.3.1 Example: Analog-to-Digital Clock-Signal Conversion

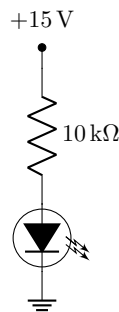
As an example application of a Schmitt trigger, consider the conversion of an analog clock signal to a digital clock. An example of an analog clock source is a rubidium atomic clock, a relatively inexpensive (few \$K) instrument that typically provides a 10-MHz sine wave, with a stability/accuracy of better than a part in  $10^{10}$ . To drive digital circuits, however, this should be converted into an appropriate square wave with logic-level inputs. The circuit below accomplishes this, using a Schmitt trigger to avoid spurious extra clock pulses due to noise on the clock signal.



Note the ac-coupled input, a 50- $\Omega$ -terminated input (as appropriate for a 50- $\Omega$  cable connecting the clock to this circuit), and the Schmitt trigger, with TTL-compatible output. The clock offset and threshold voltages are adjustable to account for different input clock amplitudes, and to adjust the duty cycle of the resulting square wave (typically, these can both be set to 2.5 V).

## 14.4 Circuit Practice

For comparator practice, first note that you can light an LED with a power supply and a current-limiting resistor as follows.

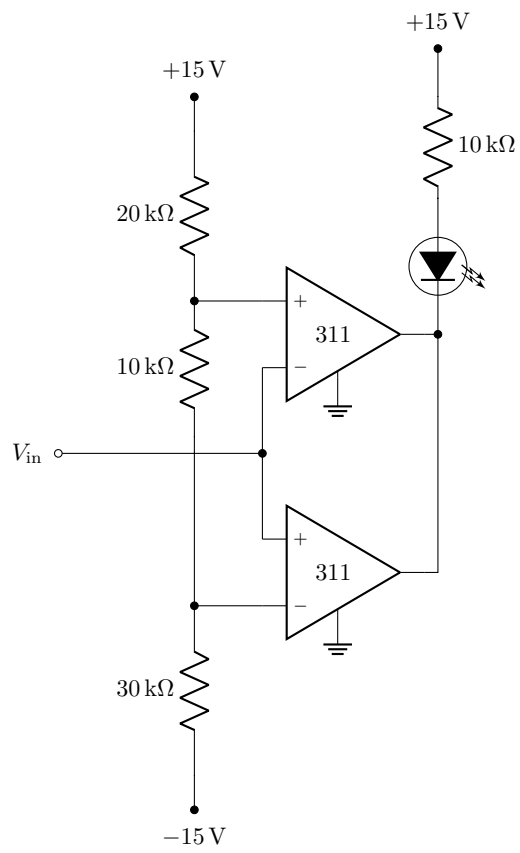


Thus, an open-collector output can control the LED, since the output is either open or shorted to ground (corresponding to an OFF or ON LED, respectively).

Now design a “TTL out-of-range alarm,” given the following components and requirements:

- two 311’s
- any resistors you like
- an LED
- $\pm 15$ -V power supplies (and ground)
- the circuit operates as follows: if  $V_{in} > +5$  V or  $V_{in} < 0$  V, the LED “alarm” lights; otherwise, the LED is off

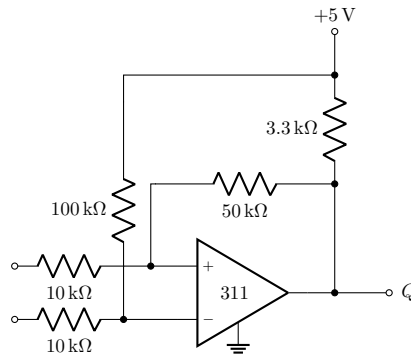
**Solution.**



## 14.5 Exercises

### Problem 14.1

(a) Show that the comparator circuit below acts as a flip-flop (“**bistable multivibrator**”). Label the set and reset inputs.



(b) Show how to add another (311) comparator to also produce the  $\overline{Q}$  output. (Use whatever resistors you like.)

### Problem 14.2

Briefly describe how a Schmitt trigger in an input of a logic gate can improve the performance of a circuit. Under what conditions on the input signal do you expect an improvement?

## Chapter 15

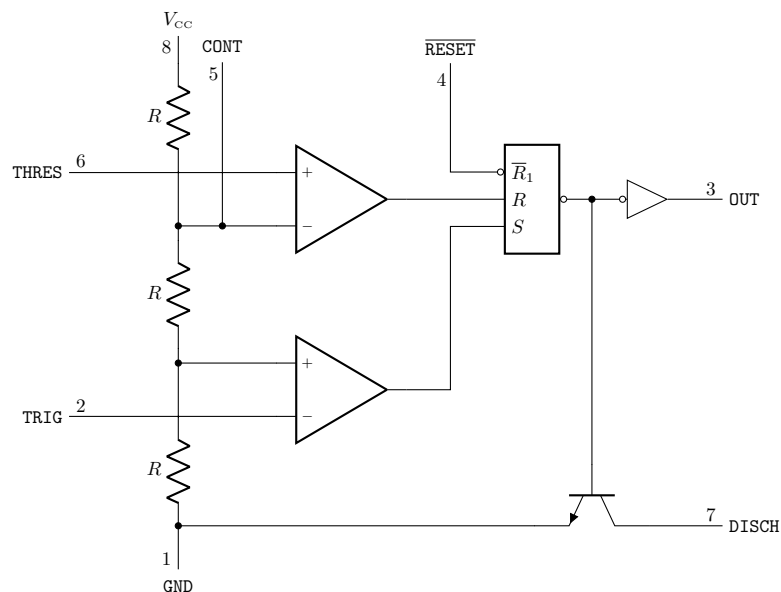
# Pulse and Waveform Generation

### 15.1 The Classic 555 Timer

The 555 timer is an old, classic workhorse for timing and pulse-signal generation, in applications that require moderate accuracy (1%) and relatively slow signals (good performance up to a few hundred kHz, and can be pushed up to  $\sim 1$  MHz). It's a versatile chip, and can produce square waves, arbitrary-length pulses, and can perform more complicated tasks such as pulse-width modulation—there is a lot of functionality packed into this 8-pin chip.

#### 15.1.1 Equivalent Circuit

The functional equivalent circuit for the 555 chip is shown below, with pin assignments for an 8-pin DIP package.<sup>1</sup>



A few things to notice here:

1.  $V_{CC}$  powers all the components in the chip. The output is TTL compatible when  $V_{CC} = +5$  V, but  $V_{CC}$  can go up to +18 V. The minimum for the standard chip is +4.5 V, but some CMOS variants can make use of lower voltages (e.g., +2 V for the ICL7555, +1 V for the TLC551). The GND (ground) input sets the ground reference for the circuit.

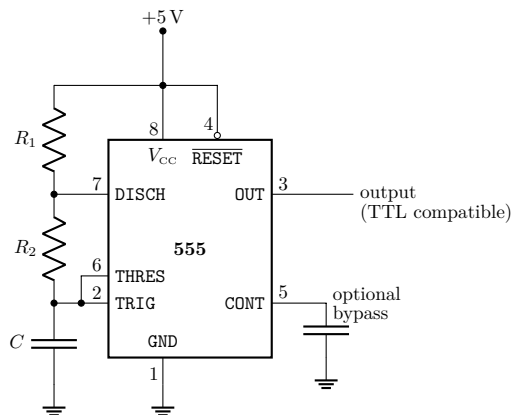
<sup>1</sup>For a nice view inside the package and discussion of the “guts” of the 555, see <http://www.righto.com/2016/02/555-timer-teardown-inside-worlds-most.html>.



2.  $V_{CC}$  also drives a resistor divider chain, which sets voltage-reference points at  $(2/3)V_{CC}$  and  $(1/3)V_{CC}$ . The **CONT** (control) input taps into the  $(2/3)V_{CC}$  point, and can be used as an “override” for the reference voltages.
3. Two comparators compare two input voltages (**THRES** or threshold, and **TRIG** or trigger) to these reference voltages.
4. The comparator outputs drive the inputs of an SR flip-flop. The flip-flop also has an externally connected, direct-reset input ( $\overline{\text{RESET}}$ ).
5. The flip-flop deals with “bad” inputs in way that could depend on the exact flavor of the 555—it’s best to check things with the data sheet. The original NE555, for example, specifies that the  $\overline{\text{RESET}}$  input can override the other inputs. The TLC551, as another example, specifies that  $\overline{\text{RESET}}$  can override **TRIG**, which can in turn override **THRES**.
6. The output of the flip-flop is buffered and set to the **OUT** port. The standard 555 (NE555) has a lot of “oomph,” and can handle  $\pm 200$  mA of current. Note that “setting” the flip-flop takes the output **HIGH**.
7. The flip-flop output also connects to the base of an open-collector, NPN transistor, whose collector is the **DISCH** (discharge) output. This output is useful, for example, for dumping the charge of a timing capacitor. In this case, “resetting” the flip-flop will turn the transistor on (i.e., capacitor charge gets dumped), setting the flip-flop will turn the transistor off.

### 15.1.2 Astable Multivibrator

A typical 555 circuit is shown below. This is called an **astable multivibrator**, which just means that the output is a square wave, whose timing is set by the external components.



Let’s analyze how this works:

1. The bypass capacitor at pin 5 (**CONT**) helps to stabilize the reference voltages in the resistor chain. One problem with the 555 is that it can cause large power-supply transients while switching, which can feed into the reference chain and cause multiple (unintended) transitions and glitching, and a capacitor here helps to fight these problems. It’s a good idea to bypass pin 8 ( $V_{CC}$ ) with a large capacitor as well.
2. The main timing of the square wave is controlled by charging and discharging the capacitor  $C$ . The capacitor charges from  $V_{CC}$  via  $R_1 + R_2$ , while the capacitor discharges only via  $R_2$  to pin 7, when the transistor is on and this pin is connected to ground.

3. We will assume that the capacitor is initially uncharged. Note that as long as THRES (pin 6) is below  $(2/3)V_{CC}$ , the corresponding comparator output is LOW, which does not reset the flip-flop. The other comparator sets the flip-flop since the capacitor voltage at pin 2 (TRIG) is below  $(1/3)V_{CC}$ . This sets the output HIGH and turns the transistor off, so the capacitor charges via  $R_1 + R_2$ .
4. When the capacitor voltage reaches  $(2/3)V_{CC}$ , the flip-flop gets reset by the THRES comparator, setting the output LOW and turning the transistor on, shorting pin 7 (DISCH) to ground.
5. The capacitor discharges through  $R_2$  until the voltage drops to  $(1/3)V_{CC}$ , when the flip-flop sets again, and the process repeats.
6. For either phase, we can use the exponential relaxation of the RC circuit to figure out the times in each state. For the output LOW-output phase, the capacitor voltage relaxes exponentially (with time  $R_2C$ ) from  $(2/3)V_{CC}$  towards 0 (neglecting the small collector-emitter voltage across the discharge transistor). That is, the capacitor voltage is

$$V(t) = \frac{2}{3}V_{CC} e^{-t/R_2C}, \quad (15.1)$$

and if we set  $V(\tau_{\text{low}}) = (1/3)V_{CC}$ , where  $\tau_{\text{low}}$  is the LOW-cycle time, we get

$$\tau_{\text{low}} = (\log 2)R_2C \approx 0.693 R_2C. \quad (15.2)$$

Note that “log” here is the natural logarithm.

7. For the output HIGH-output phase, the capacitor voltage charges via an exponential “decay” [with time constant  $(R_1 + R_2)C$ ] from  $(1/3)V_{CC}$  towards  $V_{CC}$ . The shortcut to evaluating the HIGH time is to note that, except for the time constant, this is the same as the LOW phase, except for an inversion and a shift in voltage. Thus, a similar result applies, with

$$\tau_{\text{high}} = (\log 2)(R_1 + R_2)C \approx 0.693 (R_1 + R_2)C \quad (15.3)$$

as the high-time dwell state. To do this a little more slowly, the upwards exponential decay is

$$V(t) = V_{CC} - \frac{2}{3}V_{CC} e^{-t/(R_1+R_2)C} \quad (15.4)$$

(check the limits at  $t = 0$  and  $t \rightarrow \infty$  to convince yourself this expression does the right thing). The result (15.3) follows from setting  $V(\tau_{\text{high}}) = (2/3)V_{CC}$  and solving for  $\tau_{\text{high}}$ .

Observe that  $\tau_{\text{high}} > \tau_{\text{low}}$  in this circuit; we can't have  $R_1 = 0$ , otherwise the discharge transistor would attempt to short  $V_{CC}$  to ground—not a good idea.

8. The oscillation period is then

$$T = \tau_{\text{low}} + \tau_{\text{high}} = (\log 2)(R_1 + 2R_2)C \approx 0.693 (R_1 + 2R_2)C. \quad (15.5)$$

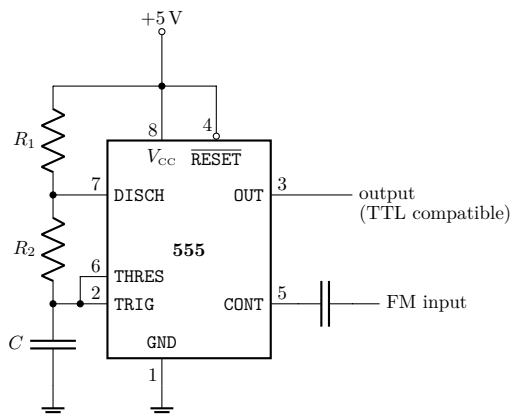
Note that  $V_{CC}$  dropped out of this expression, which is convenient, since it means the timing is insensitive to the power-supply voltage. Of course, we assumed that the power supply is *constant*; if  $V_{CC}$  varies on the time scale of the period, such that the voltage is differs between successive transitions, then  $V_{CC}$  *doesn't* drop out of the period. Hence, again, the importance of a bypass capacitor for the power supply.

9. Of course, there are some limitations to what is possible with the above timing. Lancaster<sup>2</sup> recommends:  $R_1 + R_2 \leq 3.3\text{M}\Omega$ ,  $R_1, R_2 \geq 1\text{k}\Omega$ ,  $C \geq 500\text{pF}$ . The capacitor can be large—the  $RC$  times may ultimately be limited by capacitor leakage, but hours-long periods are possible.
10. A nearly symmetric square wave results if  $R_1 \ll R_2$ , but if the asymmetry is a problem, the oscillator frequency can be doubled, and the output fed through a divide-by-two circuit, resulting in a symmetric wave, independent of the initial asymmetry (why?).

<sup>2</sup>Don Lancaster, *The TTL Cookbook* (SAMS/Prentice Hall, 1974), p. 173.

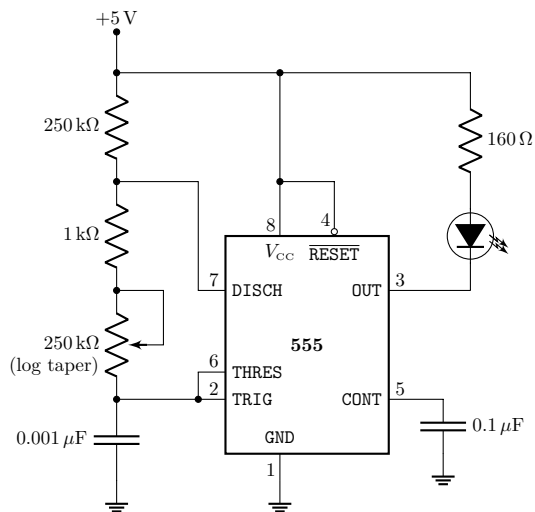
### 15.1.2.1 Frequency Modulation

A variation on this circuit is to use the **CONT** input to modulate the  $(2/3)V_{CC}$  reference point, which modulates the frequency (raising the reference voltage = longer time to trigger the comparator = lower frequency). A coupling capacitor makes the frequency modulation more convenient, as a zero voltage corresponds to no modulation.



### 15.1.2.2 Pulse-Width Modulation: LED Dimmer

The asymmetry of the 555 output is sometimes useful, as in **pulse-width modulation**. You can control the brightness of an LED by controlling the supply current; however, if the LED is driven by a fixed-voltage, digital output, this isn't feasible. The solution is to blink it rapidly on and off, and vary the fraction of time it is on (i.e., vary the **duty cycle**—the fraction of the period where the signal is **HIGH**—of the pulse). The  $160\,\Omega$  resistor here sets the LED current to 20 mA, assuming a 1.8-V drop across the LED (as appropriate for a standard red LED), and assuming the output goes down to 0 V (almost true). The component values give a minimum frequency of about 1.3 kHz (plenty fast to make the LED appear continuous—it should be over about 50 Hz), and the duty cycle varies from a maximum of about 50% down to a minimum of about 0.4%. (Note that the LED is on when the output is **LOW**.) This is also an *efficient* dimming scheme, which is why it is so common: restricting the current generally requires some wasted power, due to power dissipated in a current-limiting resistance. In this circuit, the efficiency is roughly the same at any value (neglecting power dissipated due to output transitions).

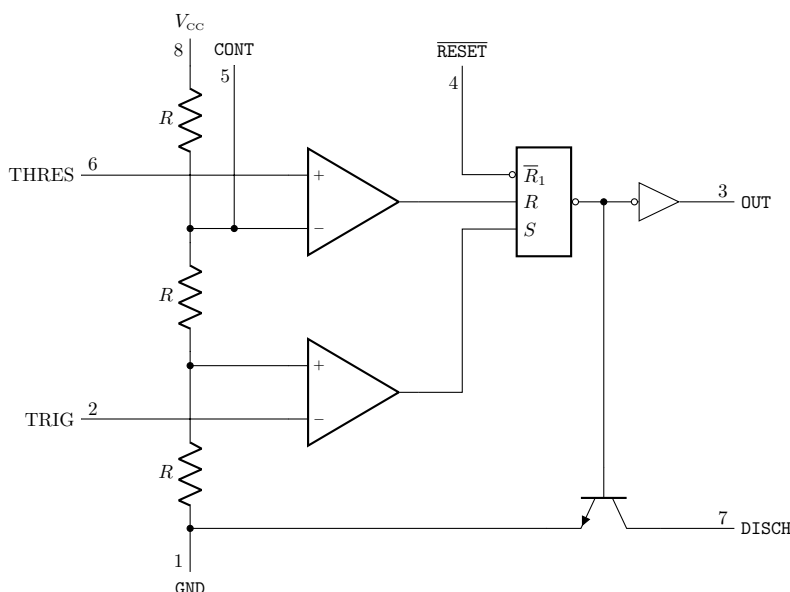


## 15.2 Monostable Multivibrators

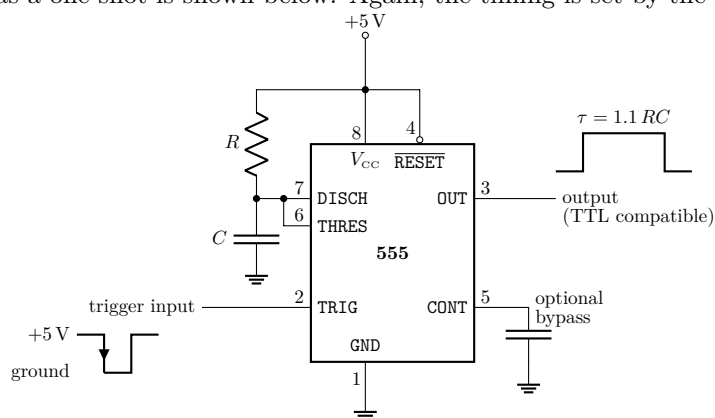
A **monostable multivibrator**, or **one-shot**, is a circuit that simply supplies a single digital pulse of some defined duration. This is the fundamental building block of timing circuits—circuits that control the timing of sequences of events (laser pulses, data acquisition, camera triggers, etc.).

### 15.2.1 555 as a One-Shot

The 555 can be conveniently hooked up as a one-shot. Recall the internals of the 555:



The connection as a one-shot is shown below. Again, the timing is set by the external components.



Let's see how this works:

1. The TRIG input is normally high, which is consistent with the internal flip-flop being reset, and thus the output being LOW. In this case the transistor is on, so the capacitor is shorted to ground. We will assume all this to be the case; if the flip-flop is indeed set, it is as if the circuit has been triggered, and the circuit will go through the rest of the cycle and reset, and then we can proceed with the assumption of a reset flip-flop.
2. The negative edge of a trigger pulse starts the output pulse. This happens when the TRIG input crosses below  $(1/3)V_{CC}$ , and the lower comparator sets the flip-flop. This takes the output HIGH, and turns the transistor off.

3. The timing of the output pulse is controlled by charging the capacitor  $C$ . The capacitor charges from  $V_{CC}$  via  $R$ .
4. The capacitor is uncharged at the beginning of the pulse, and it will charge until the capacitor voltage (and thus **THRES** input) rises to  $(2/3)V_{CC}$ . At this time, the upper comparator resets the flip-flop, ending the pulse and dumping the capacitor charge.
5. Again, we can use the exponential relaxation of the RC circuit to figure out the pulse duration. The capacitor voltage rises exponentially (with time  $RC$ ) from 0 to  $(2/3)V_{CC}$  towards  $V_{CC}$ . This is equivalent to an exponential decay from  $V_{CC}$  towards 0 to  $(1/3)V_{CC}$ . Thus we can consider the decay

$$V(t) = V_{CC} e^{-t/RC}, \quad (15.6)$$

and if we set  $V(\tau) = (1/3)V_{CC}$ , where  $\tau$  is the pulse duration, we get

$$\tau = (\log 3)RC \approx 1.1 RC. \quad (15.7)$$

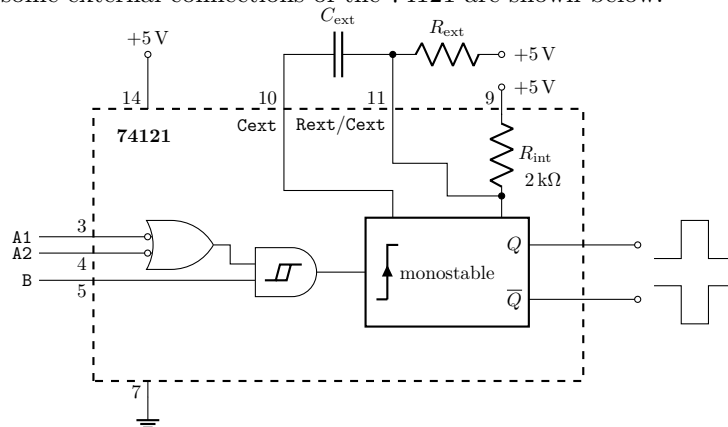
Again, “log” here is the natural logarithm, and note that this is independent of the supply voltage.

Note that we have assumed that by the time the flip-flop is reset, the trigger pulse at the **TRIG** input is back to **HIGH**; otherwise the **S** input to the flip-flop will also be high, putting the flip-flop into the “bad” state. (In reality, typical 555’s will allow the **TRIG** to override **THRES** to avoid any such bad states.) The net effect is that the output pulse will be “stretched” beyond the  $RC$  duration until the trigger input goes high. One workaround for this (if you’re stuck with long pulses but want to trigger short pulses) is to run the trigger input through a differentiator (with a short  $RC$  time) and then into the 555.

Also, note that the usual intent of this circuit is to make fairly short (ms-scale) timing pulses, so that the timing capacitor is typically small ( $\sim \mu\text{F}$  or smaller). In this case, there is no problem with the 555’s transistor effectively shorting the capacitor to dump the charge. However, in the case of a very large capacitance, especially with a large  $V_{CC}$ , the large current that may flow when dumping the capacitor charge could cause problems, in which case a small current-limiting resistance in the capacitor-DISCH path would be prudent.

### 15.2.1.1 The 74121

A number of other monostable multivibrators are available, more-or-less prepackaged. A good example is the 74121, which is faster than the 555 circuit—it can be programmed for pulses down to 35 ns, and up to 28 s. The “guts” and some external connections of the 74121 are shown below.



Some operation notes:

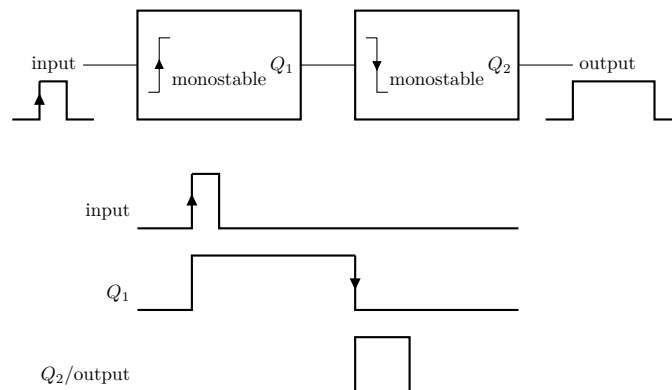
1. As in the 555, the internal monostable is controlled by the time for the external capacitor to charge via a resistor. The resistor can be internal ( $R_{int}$ ) if pin 9 is connected to  $+5V$ , or external, if  $R_{ext}$  is connected between pin 11 and  $+5V$ . (Of course, if both are connected, then we have the parallel resistance of the two.) The pulse duration is given by

$$\tau = (\log 2)RC \approx 0.693 RC. \quad (15.8)$$

2. The internal monostable is triggered by a rising pulse edge, but there is some extra input logic to make this more flexible. For example, suppose **A1** and **B** are held **HIGH**. Then a falling edge on **A2** triggers the pulse. Of course, **A1** and **A2** can be exchanged here.
3. Similarly, if either **A1** or **A2** are held **LOW**, then a rising edge into **B** will trigger the flip-flop. In this case, the **B** input goes into a Schmitt trigger (0.2-V hysteresis), and thus can handle slow/noisy inputs.
4. The 74121 is **nonretriggerable**; this means the device will ignore any input edges while it is generating an output pulse. Other one-shots, like the 74123, is **retriggerable**, which means that a new triggering edge will always start a new timing cycle, even if already in the middle of a timing cycle.

### 15.2.1.2 Combining One-Shots: Pulse Delay

In complex timing systems, where many things must happen at the proper times, many one-shots can be chained together to generate the proper timing sequence. As simple example, two one-shots can be chained together to generate a *delayed* pulse from an initial trigger pulse. This times two events with a fixed delay, such as launching a projectile, and then a short time later triggering a photographic flash.

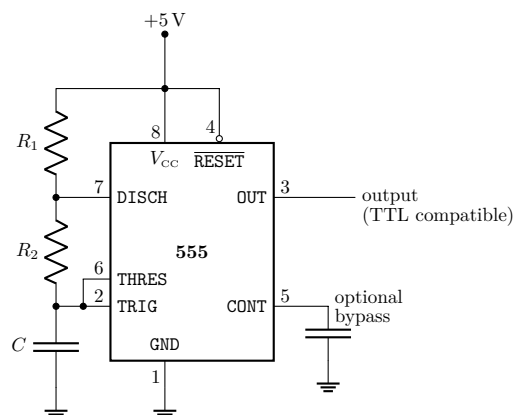


Note the different edge triggers of the two monostables.

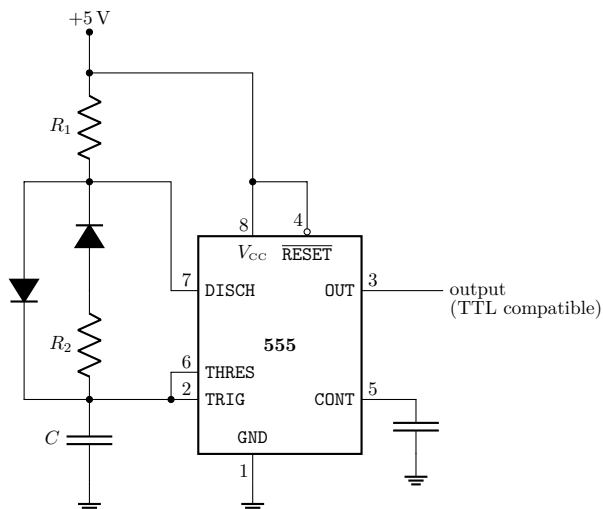
## 15.3 Circuit Practice

### 15.3.1 Duty-Cycle Control

What is the duty cycle of the basic astable circuit (reproduced below)?



What is the duty cycle of the modified astable shown below? (Ignore any voltage drops across the diodes.) How does it allow better control over the duty cycle?



**Solution.** With this arrangement, the capacitor charges through  $R_1$  and the left-hand diode, and discharges through  $R_2$  and the right-hand diode. In either case, we can work out the timing as follows. Suppose the capacitor discharges from voltage  $(2/3)V$  towards zero, with time constant  $RC$ . We need to solve for the time  $\tau$  when the capacitor voltage is  $(1/3)V$ :

$$\frac{2}{3}V e^{-\tau/RC} = \frac{V}{3}. \quad (15.9)$$

The solution is

$$\tau = (\log 2)RC. \quad (15.10)$$

The HIGH output cycle is the charging cycle, so the high time is

$$\tau_{\text{HIGH}} = (\log 2)R_1C. \quad (15.11)$$

The LOW output cycle is the discharge cycle, so the high time is

$$\tau_{\text{LOW}} = (\log 2)R_2C. \quad (15.12)$$

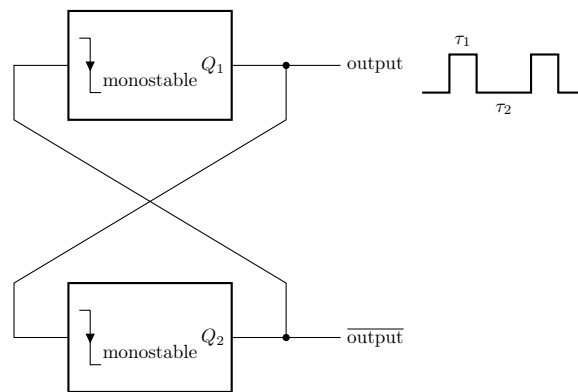
The period is the sum of these, or

$$T = (\log 2)(R_1 + R_2)C. \quad (15.13)$$

### 15.3.2 Astable Multivibrator

Another example of combining one-shots is to combine two, in order to make an astable multivibrator. How can you do this?

**Solution.** The circuit for this is shown below.



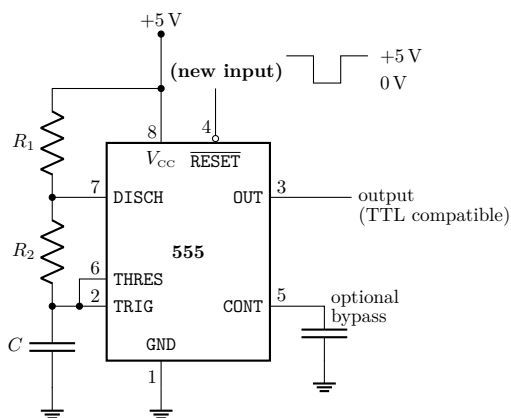
Note that the HIGH and LOW times are controlled separately by the durations of the two one-shots, and each one-shot triggers on the falling edge of the other.



## 15.4 Exercises

### Problem 15.1

Shown below is the classic configuration of the 555 timer as an astable multivibrator, **except** that we're treating the  $\overline{\text{RESET}}$  line as an extra input (recall that normally it's just tied to +5 V for normal astable operation).

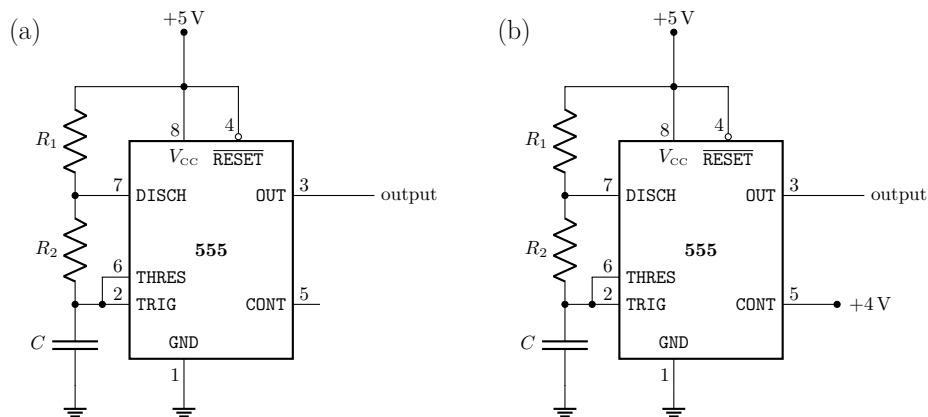


Describe what happens when the  $\overline{\text{RESET}}$  input goes through one HIGH–LOW–HIGH pulse cycle as shown. More specifically, you should say what happens to the OUT line and capacitor voltage on each transition (HIGH–LOW and LOW–HIGH), and what happens after the input pulse is finished.

*Note:* on the internal flip-flop of the 555, you should assume the  $\overline{R}_1$  input overrides the others.

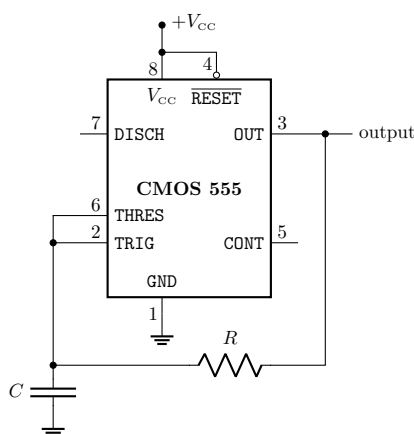
### Problem 15.2

Shown below is (a) the basic 555 astable multivibrator and (b) a modified version. For each circuit compute the period and duty cycle in terms of  $R_1$ ,  $R_2$ ,  $C$ , and any relevant voltages.



### Problem 15.3

Shown below is a 555-based astable multivibrator with 50% duty cycle. Show that it works as advertised, and compute the oscillation period.



Note: this circuit only works with CMOS variants of the 555 (like the 7555), where the output will swing all the way from ground to  $+V_{CC}$ .

#### Problem 15.4

Suppose you have a standard TTL logic gate, and you want to put a Schmitt trigger on one of the inputs to address a problematic input signal. However, suppose that all you have is a 555; show how to wire a 555 to act as a Schmitt trigger (actually, an inverting Schmitt trigger). Make sure to explain *why* your circuit acts as a Schmitt trigger, and what are the input logic levels.

#### Problem 15.5

Work out the period of the basic 555 astable multivibrator circuit, but this time account for the voltage drop  $V_d$  across the discharge transistor when it is turned on. That is, the supply voltage should no longer drop out of the result. Assuming  $V_d = 0.2$  V, how much does this affect the period compared to the idealized value if  $V_{CC} = +5$  V?

#### Problem 15.6

In one incarnation, the **bicolor LED** looks like an ordinary LED (plastic package with two leads), but is really 2 LEDs in parallel, with one reversed. That is, if current flows “forwards,” the LED lights green, and if it flows “backwards,” the LED lights red.

Design a circuit *two* 555’s to light a bicolor LED, alternating between red and green. For concreteness, design for a  $\sim 50\%$  duty cycle (approximately equal time in each color), with a period of 1 s. Use whatever passive components you like, but be specific about their values. Also, show all pin connections on the 555’s.

*Hint:* note that it would be a bad idea to configure two 555 as independent oscillators. Why? Instead try using one of the 555’s as a NOT gate.

#### Problem 15.7

Design a 10 kHz square-wave oscillator (50% duty cycle) using only 74121’s and capacitors.

#### Problem 15.8

Show how to build a 4-bit ripple counter (mis)using **only** 555 chips. (Plus sufficient wire and a +5-V power supply; no other components needed/allowed. You may assume the clock signal already exists.) Some notes:

- Your solution could (should) end up involving a considerable number of 555's. Stay *organized* and present your solution without drawing out each and every 555 explicitly. Otherwise, use as many 555's as you like—they're free, at least on paper. Start by figuring out what logic elements you'll need (gates, flip-flops), and how to repurpose 555's to make them happen.
- For concreteness of operation, assume the actual 555 involved to be a TLC555. Recall that the internal flip-flop of this chip deals with potential "bad" inputs by allowing  $\overline{\text{RESET}}$  to override TRIG, which can in turn override THRES.
- Be careful with the clock-input operation of the counter; remember that a ripple counter typically employs edge-triggered clock inputs. You may want to seek out (and study!) an example of how clock-edge triggering is implemented in another circuit.

## Chapter 16

# Digital–Analog Interfaces

One of the most important concepts in digital electronics is *interfacing* digital circuits to analog circuits. If an analog signal serves as the input to a digital circuit, then we need **analog-to-digital conversion (ADC)**, while a digital circuit generating an analog signal requires **digital-to-analog conversion (DAC)**. We will consider the latter first, which is simpler, and ADC often relies on DAC.

### 16.1 Digital-to-Analog Conversion

Digital-to-analog conversion is very common in everyday circuits. This is required to generate the audio signals in cell phones and CD/DVD/MP3 players, and to generate the output intensity (or color) of displays in LCD projectors or in CRT/plasma/LCD displays. Essentially, any analog signal coming out of a computer must have gone through the DAC process.

#### 16.1.1 Resolution

Before understanding how DAC circuitry works, let’s review some of the resolution requirements for representing analog signals. Analog signals must be **sampled**—that is instead of a continuous function  $y(t)$ , we must represent it via samples  $y_j := y(t_j)$  at sample times  $t_j$  (typically regularly spaced), and the values of  $y_j$  must be represented with some finite precision (i.e., it must be represented with a finite number of bits). In terms of amplitude resolution, if there are  $N$  bits of data, then there are  $2^N$  different signal levels available within a defined range (e.g., within some voltage range). The signal levels could be positive only, represented by unsigned integers, or positive/negative using signed integers (or unsigned integers after adding an offset that ensures the signal is always positive). In this case, since the “real” signal must always be rounded to the nearest available level, the *fractional* sampling resolution is  $2^{-N}$ , and the absolute resolution is  $2^{-N}V_{\text{range}}$  for a voltage signal if  $V_{\text{range}}$  is the total voltage range available for sampling the signal. Since the rounded value should be to the nearest sampling value, the maximum error is  $1/2$  of the resolution, or  $2^{-(N+1)}$  maximum fractional error for signals within  $V_{\text{range}}$ . So, for 16-bit sampling, the error is at worst about 8 ppm.

In terms of timing resolution, the requirement is set via the **sampling theorem**. Suppose we sample a signal every  $\Delta t$  in time. Then the **sampling rate** is given by

$$\text{sampling rate} = \frac{1}{\Delta t}. \quad (16.1)$$

Then we can also define the **Nyquist frequency** by

$$\text{Nyquist frequency} = \frac{1}{2\Delta t} = \frac{\text{sampling rate}}{2}. \quad (16.2)$$

It turns out that, according to the sampling theorem, the Nyquist frequency is the largest frequency that is accurately reproduced by the sampled signal.<sup>1</sup> For example, in compact-disc (CD) audio, the sampling rate

---

<sup>1</sup>For details, see Daniel Adam Steck, *Quantum and Atom Optics*, available online at <http://steck.us/teaching>.

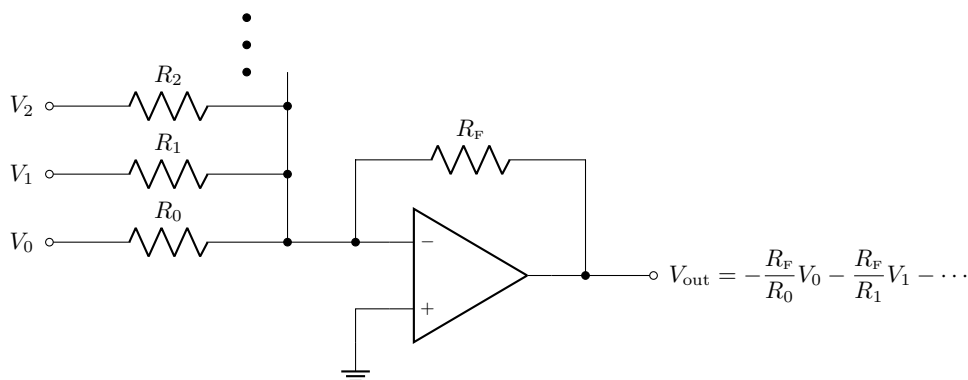
is 44.1 kHz, since the goal is to reproduce audio frequencies up to about 20 kHz. Note that 44.1 kHz is then a bit above the Nyquist frequency, which allows for extra tricks, like an **anti-aliasing filter**, to improve the quality of the reconstructed audio. The idea here is to guard against **aliasing**, which is the error suffered by frequencies *above* the Nyquist frequency—they are spuriously represented as lower (sub-Nyquist) frequencies in the sampled signal. Generally, a low-pass filter is used to remove these high frequencies, but since the filter does not have a perfectly sharp cutoff, the extra sampling rate above the Nyquist frequency accommodates the desired audio range, while giving some bandwidth for the low-pass filter to have a significant effect before aliasing errors occur.

### 16.1.2 DAC Circuitry

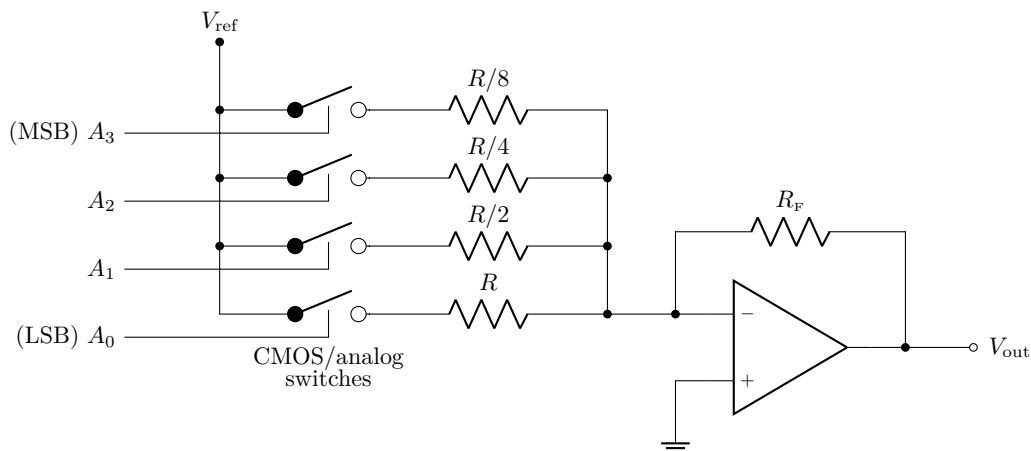
Now, how do we make a DAC? The basic ingredient is a summing (inverting) amplifier, which you may recall from Section 7.3.4. Recall that this takes a number of voltage inputs  $V_0$ ,  $V_1$ , ..., and has as output

$$V_{\text{out}} = -\frac{R_F}{R_0}V_0 - \frac{R_F}{R_1}V_1 - \frac{R_F}{R_2}V_2 - \dots, \quad (16.3)$$

which is an inverted, weighted sum, where the relative weights are controlled by the input resistors, and the weights have the feedback resistance  $R_F$  in common.



Then, for example, we can build a 4-bit DAC as in the diagram below.



Here,  $A_0$ – $A_3$  are digital inputs, with  $A_0$  the LSB and  $A_3$  the MSB of an unsigned integer. The inputs drive analog switches, which conduct when  $A_j = 1$  and are open when  $A_j = 0$ . The voltage  $V_{\text{ref}}$  sets the (absolute)

voltage resolution and the range of the conversion. Then using the summing-amp formula,

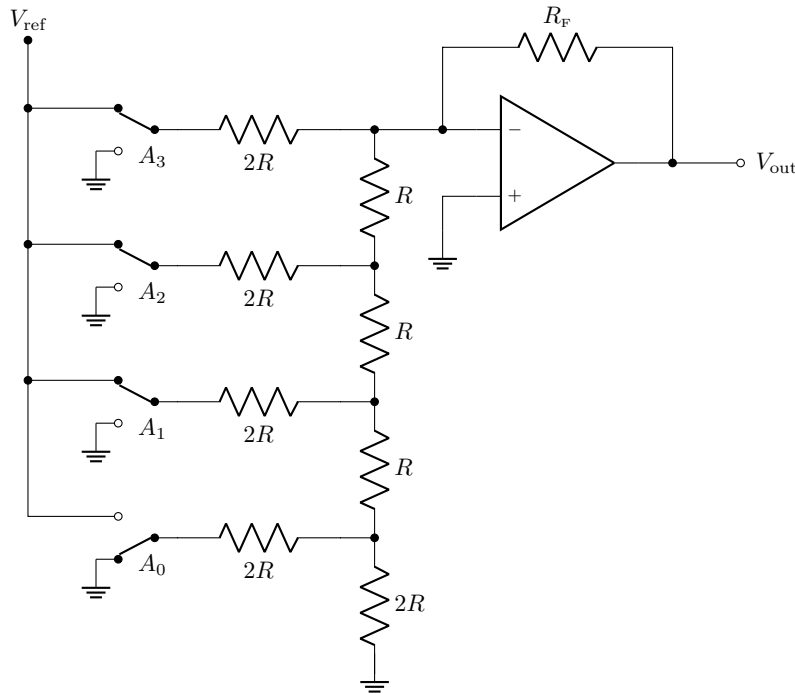
$$\begin{aligned} V_{\text{out}} &= -A_0 \frac{R_F}{R} V_{\text{ref}} - 2A_1 \frac{R_F}{R} V_{\text{ref}} - 4A_2 \frac{R_F}{R} V_{\text{ref}} - 8A_3 \frac{R_F}{R} V_{\text{ref}} \\ &= -\frac{R_F}{R} V_{\text{ref}} (A_0 2^0 + A_1 2^1 + A_2 2^2 + A_3 2^3). \end{aligned} \quad (16.4)$$

(DAC output)

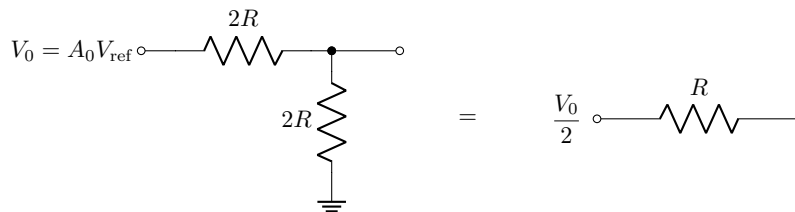
So, for example, an input of 0011 corresponds to  $V_{\text{out}} = -3(R_F/R)V_{\text{ref}}$ . Of course, usually we want a *positive* output, which requires another inverting amplifier.

### 16.1.3 R–2R Ladder

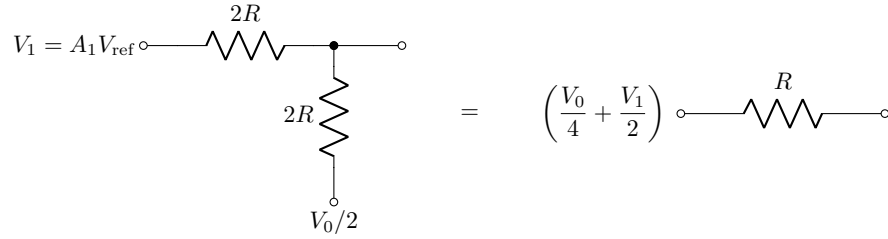
One problem with the above circuit is that it requires a series of resistors of different values to be made to high accuracy, which is difficult, especially for a high-resolution DAC. It is much easier to make sets of *matching* resistors, and there is a circuit that takes advantage of this, called the **R–2R ladder**, which uses only resistors of size  $R$  and  $2R$ . The idea is below.



Again, the inputs are  $A_0$ – $A_3$  here, controlling analog/CMOS switches (here SPDT). To see how this works, consider the Thévenin-equivalent circuit for the  $A_0$  input. This is a simple 50% voltage divider, so the equivalent resistance is  $R$ , and the voltage is half the input  $A_0 V_{\text{ref}}$ . (Note that we consider  $A_0 = 0$  if the switch is down,  $A_0 = 1$  if the switch is up.)



Now lumping this equivalent circuit into the next “stage” with the  $A_1$  input, we have a similar voltage-divider situation.



Note that the Thévenin-equivalent voltage is just the average of the two input voltages. Continuing this process, we find that at each stage we add in half of the next input voltage and divide the remaining ones by two. The result is

$$\begin{aligned}
 V_{\text{out}} &= -V_{\text{ref}} \frac{R_F}{R} \left( \frac{A_0}{16} + \frac{A_1}{4} + \frac{A_2}{8} + \frac{A_3}{2} \right) \\
 &= -\frac{V_{\text{ref}} R_F}{2^4 R} (A_0 2^0 + A_1 2^1 + A_2 2^2 + A_3 2^3).
 \end{aligned}
 \tag{16.5}$$

(output of  $R$ – $2R$  ladder)

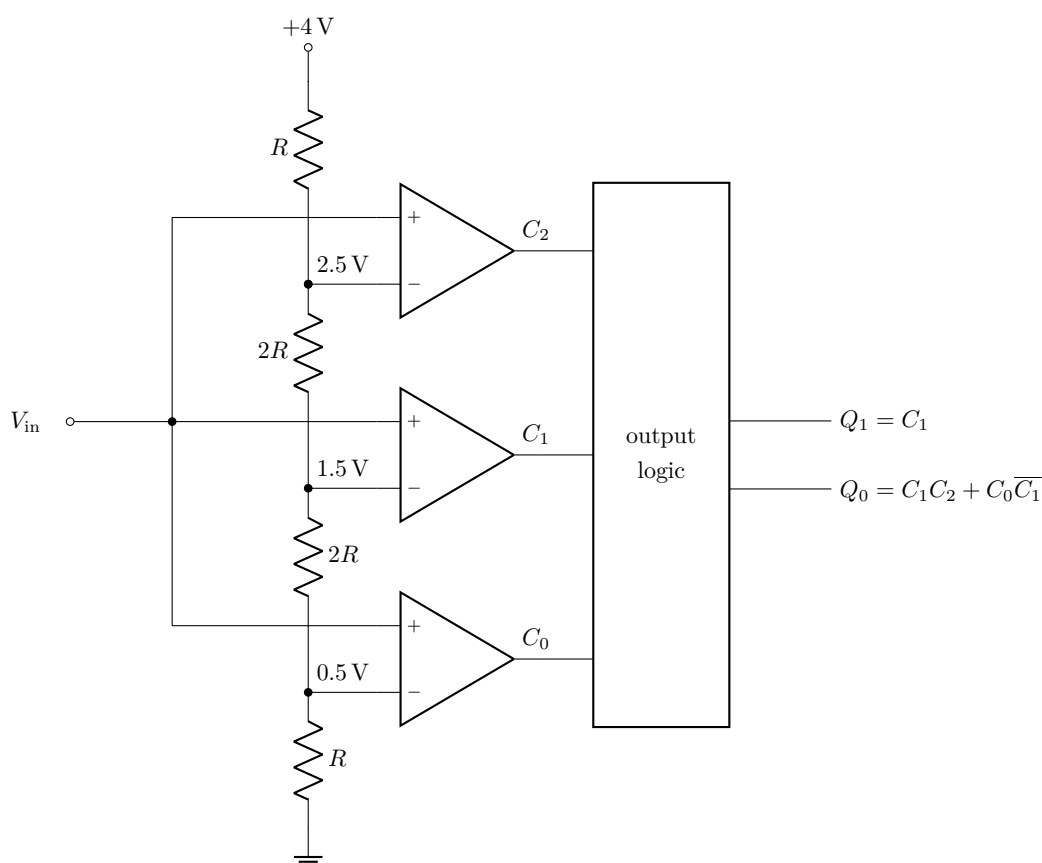
Thus, up to an overall factor, we obtain the same output as in Eq. (16.4).

## 16.2 Analog-to-Digital Conversion

The complementary process to digital-to-analog conversion is **analog-to-digital conversion (ADC)**. We will go through several ADC methods.

### 16.2.1 Flash ADC

A conceptually simple method for ADC is **flash ADC** or **parallel-encoding ADC**. The idea is to use a voltage-divider chain to create many reference voltage, and a separate comparator is used to compare the input voltage to each reference. Then output logic is needed to properly encode the digital output as a binary number. As a simple example, consider the 2-bit flash ADC below.



The idea is to choose digital voltage levels (“quantization level”) of 0, 1, 2, and 3 V. Then the maximum input range is  $-0.5$ – $3.5$  V with a maximum error of  $0.5$  V. Then the conversion ranges with comparator outputs are enumerated below.

voltage range	digital output	comparator output $C_2C_1C_0$
$-0.5$ – $0.5$ V	00	000
$0.5$ – $1.5$ V	01	001
$1.5$ – $2.5$ V	10	011
$2.5$ – $3.5$ V	11	111

Note the logical expressions included to encode the three comparator outputs into two-bit binary.

The main advantage of a flash ADC is that it is *fast*: the signal just needs to propagate through the comparators and gates, and the ADC can sample rapidly changing signals. The main disadvantage is that for  $N$  bits, there must be  $2^{N-1}$  comparators, which is difficult for more than about 10 bits of resolution.

### 16.2.2 Successive Approximation

A slower, but more generally useful ADC method is **successive approximation**. This method is analogous to the root-finding problem: Suppose  $f(x)$  is a continuous function with a single root in  $(a, b)$ . That is,  $f(a)f(b) < 0$ . Then how do we find the root; i.e., how do we find  $x_0$  such that  $f(x_0) = 0$ ? The **bisection method** for root-finding works as follows.

1. We know  $(a, b)$  brackets  $x_0$ . So let  $\bar{x}_0 = (a + b)/2$  be the initial best estimate for  $x_0$ .
2. If  $f(a)f(\bar{x}_0) < 0$ , then  $(a, \bar{x}_0)$  brackets the root.
3. Otherwise,  $(\bar{x}_0, b)$  brackets the root.

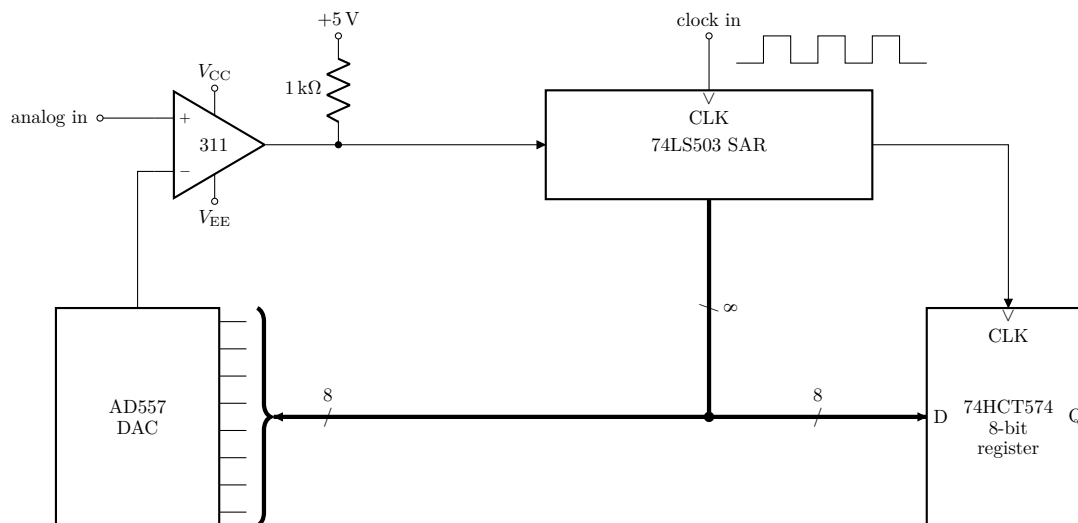


4. Redefine  $(a, b)$  to be the new, tighter bracketing interval for  $x_0$ , and repeat.

This process converges exponentially, because the width of the bracketing interval is halved on each iteration. If  $\Delta x = b - a$ , then the initial worst-case error in the estimate  $\bar{x}_0$  is  $\Delta x/2$ . After  $N$  iterations, the error is  $\Delta x/2^{N+1}$ .

The successive-approximation approach to ADC is the same problem, but to find what digital voltage  $\bar{V}$  best corresponds to  $V_{in}$ . This is an iterative process to a predetermined accuracy, given by the number of digital bits. A circuit to implement this procedure is shown below. This uses the 74LS503 **successive-approximation register (SAR)**,<sup>2</sup>

which controls the bisection process. Most ADCs nowadays using successive approximation have all these components integrated into a single chip, with serial data output—handy for keeping pin counts low, but it’s harder to understand what is going on inside.



The SAR, the “mastermind” of the conversion process, works to find 1 bit of the digital result on each clock cycle, starting with the MSB. It does this by writing the “midpoint value” of the DAC’s range to the DAC (as the first approximation), and reads the comparison result from the comparator, which tells the SAR which half of the DAC range brackets  $V_{in}$ . On the next cycle, the SAR writes out the new midpoint of the smaller bracketing range, and records the comparison result as the next converted bit, and so on. The 8-bit latch (‘574) holds the completed conversion, while the SAR is performing the next conversion (and thus its outputs are changing).

As a 2-bit example, consider the same analog range (digital levels of 0, 1, 2, and 3 V). But now, we will consider the conversion ranges to be

$$\begin{aligned}
 < 0 \text{ V} = 00 \\
 0\text{--}1 \text{ V} &= 01 \\
 1\text{--}2 \text{ V} &= 10 \\
 > 2 \text{ V} &= 11,
 \end{aligned}
 \tag{16.6}$$

as we will see. Note the different offset compared to the flash-ADC example. Let’s assume a 1.3-V input. Then the process is as follows:

- During clock cycle #1, the SAR tries the midpoint of the whole range. We can take this to be 01 (more generally, 0111111... for  $N$  bits). The DAC voltage is then 1 V, so the comparator is HIGH, and so the MSB is 1.

<sup>2</sup>The 74LS503 is now obsolete, but see Appendix B for an implementation of this chip in a PLD. In modern chips, all the various parts are typically integrated into one chip. See, for example, the ADC0804 SA-ADC: <http://www.ti.com/lit/ds/symlink/ad0804-n.pdf>.

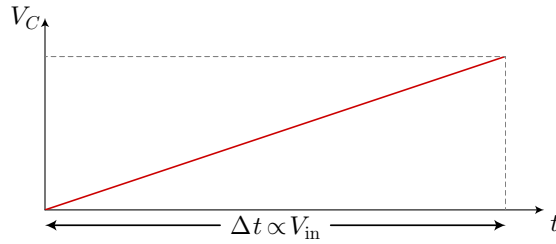
- During clock cycle #2, the SAR tries the midpoint of the remaining range (10–11). The “midpoint” is 10. The DAC voltage is then 2 V, so the comparator is now LOW, and so the LSB is 0.

Thus, the converted result is 10, in agreement with the table above.

The advantages of SA-ADC is that the timing is guaranteed (measured in clock cycles), the result can be very accurate (if a good DAC is used), and the circuit is not too complicated. The main disadvantage are that SA-ADC is slower than flash conversion, and thus may need a sample/hold circuit to deal with rapidly changing input signals.

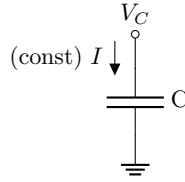
### 16.2.3 Single/Dual-Slope ADC

Another pair of important ADC methods goes under the name(s) of **single/dual-slope ADC**. The basic idea in **single-slope ADC** is to use a constant-current source to charge a capacitor, and then use a counter/clock combination to measure the time for the capacitor voltage to reach the input  $V_{in}$ .



The charge time is then proportional to the voltage. The idea is that time is very easy to measure accurately and precisely, so the method can be very accurate.

To show this mathematically, consider the capacitor-charging situation diagramed below, where a constant-current source  $I$  charges a capacitor to voltage  $V_C(t)$ .



Then using  $Q = CV_C$  and differentiating,

$$I = \frac{dQ}{dt} = C \frac{dV_C}{dt}, \quad (16.7)$$

so that for constant  $I$ ,

$$V_C(t) = \frac{It}{C}. \quad (16.8)$$

Since  $I/C$  is a constant, this can be calibrated precisely to yield  $V_C$  (and thus  $V_{in}$ ) in terms of  $t$ .

In **dual-slope ADC**, the conversion is done in two steps.

1.  $C$  is charged for a fixed time  $\tau$  by a constant current  $I \propto V_{in}$ . If we let  $\alpha$  be constant, then we can write  $I = \alpha V_{in}$ , so that

$$V_C(\tau) = \frac{\alpha V_{in} \tau}{C}. \quad (16.9)$$

2. Then,  $C$  is discharged at a constant current  $I'$ , and the discharge time  $\delta t$  is measured. Then the discharge time is fixed as in single-slope ADC by

$$V_C(\tau) = \frac{I' \delta t}{C}. \quad (16.10)$$

Thus,

$$\delta t = \frac{CV_C(\tau)}{I'} = \frac{C(\alpha V_{in} \tau / C)}{I'} = \left( \frac{\alpha \tau}{I'} \right) V_{in}. \quad (16.11)$$

In this way, we are left to calibrate  $\alpha\tau/I'$ , which is a combination of the output of a current source ( $\alpha$ ), a time  $\tau$ , and a current  $I'$ , all of which can be well-calibrated. Notably, the capacitance  $C$  dropped out; capacitances are difficult to fabricate in a way that is accurate and stable. Also, note that the first stage takes time  $\tau$ , which has an averaging effect over noise in  $V_{\text{in}}$ , whereas single-slope conversion is more apt to trigger early on a downward noise fluctuation of  $V_{\text{in}}$ .

## 16.3 Circuit Practice

### 16.3.1 Computer-Interface DAC Controller

For circuit practice, see the DAC controller board design by Todd Meyrath and Florian Schreck.<sup>3</sup> Trace through the circuit and note the following.

- The DAC7744 chips have 4 analog outputs, for 8 total output channels per board.
- An 8-bit address bus selects which DAC and output to use. The two LSBs (bits 0 and 1) select which output on a particular chip, bit 2 selects which of the two DACs on the board to activate, and the other bits select which (of possibly many) boards to address. Trace through the logic leading up to and including the NAND gates to verify that it works as advertised.
- This allows only one 16-bit data bus to feed all the outputs. The desired output is selected, the desired data is presented to the data bus, and then a strobe signal causes the addressed DAC to latch the desired output value.
- Note that the strobe pulse, which just amounts to matching the proper address, must be delayed behind the data and address signals so the inputs are settled before “load DAC” is triggered. The pulse is delayed by a buffered, RC circuit.
- Note that the NOT gates have Schmitt-trigger inputs. What part of the circuit justifies having Schmitt-input NOTs?

### 16.3.2 3-Bit ADC

Suppose you have a 3-bit DAC, with voltage levels 0 V, 0.1 V, 0.2 V, ..., in a successive-approximation ADC. If  $V_{\text{in}} = 0.35$  V,

- make a plot of the DAC output vs. time
- what is the final, converted digital value?

**Solution.** For the converted digital value: 100. The comparison voltages will be: 0.3 V, 0.5 V, 0.4 V.

<sup>3</sup><http://strontiumbec.com/Control/DAC.pdf>, in particular the schematic on p. 13.

## 16.4 Exercises

### Problem 16.1

(a) Derive an expression for the dynamic range (the largest vs. the smallest nonzero-amplitude signal) of an  $N$ -bit sampled signal. Recall that when you compare two amplitudes  $A$  and  $A_0$  in dB, the expression is

$$(\text{ratio in dB}) = 20 \log_{10} \left( \frac{A}{A_0} \right) \quad (16.12)$$

(b) What is the dynamic range for CD audio (16 bits) in dB? Bluray audio (24 bits)? (For comparison, the dynamic range of human hearing is usually quoted as 120 dB.)

### Problem 16.2

A 3-bit flash ADC uses 7 comparators to compare an input voltage to a set of voltages from a voltage-divider chain (say, 0.5 V, 1.5 V, ..., 6.5 V).

(a) Describe the possible values for the comparator outputs  $C_6 \cdots C_1 C_0$ . ( $C_0$  does the comparison to 0.5 V, and  $C_6$  handles the comparison to 6.5 V.)

(b) Give logic expressions required to transform the comparator output  $C_6 \cdots C_1 C_0$  into a 3-bit binary number (ranging from 000 to represent 0 V and 111 for 7 V). Try to end up with reasonably simple final expressions.

### Problem 16.3

Use a counter to design a simple DAC based on pulse-width modulation as follows: Your circuit should take an 8-bit digital input, representing an 8-bit unsigned integer, and then control the brightness of an LED to be proportional to this integer. Assume the clock signal to be given and to be as fast as you need it to be. You may use whatever support logic you like, but you may find it useful to use a flip-flop, and look into the binary magnitude comparator (read up on the 74688). Make sure to properly limit the LED current.

Be *specific* about any ICs you use (i.e., give the model number, like 74688, and if it matters, specify which logic family, e.g., 74HCT688). You should show all important connections, but don't bother with universal stuff like power supplies, grounds, chip enables, etc. Also, you don't need to explicitly show the clock source, just indicate the existence of the clock signal and show any connections where it enters your circuit.

### Problem 16.4

(a) Use a counter (specifically, an 8-bit up counter), a clock source (astable multivibrator), a comparator, and a DAC to design a simple ADC. That is, the counter should count upwards starting from zero, and freeze at the appropriate conversion value (i.e., the frozen counter is the output). Use whatever support logic you like, and don't worry about latching the output. You should include a start/reset input that resets the counter and allows the next conversion to start.

As in Problem 2, be *specific* about any ICs you use; show all important connections, but don't bother with universal stuff; and you don't need to explicitly show the clock source, just indicate the existence of the clock signal and show any connections where it enters your circuit.

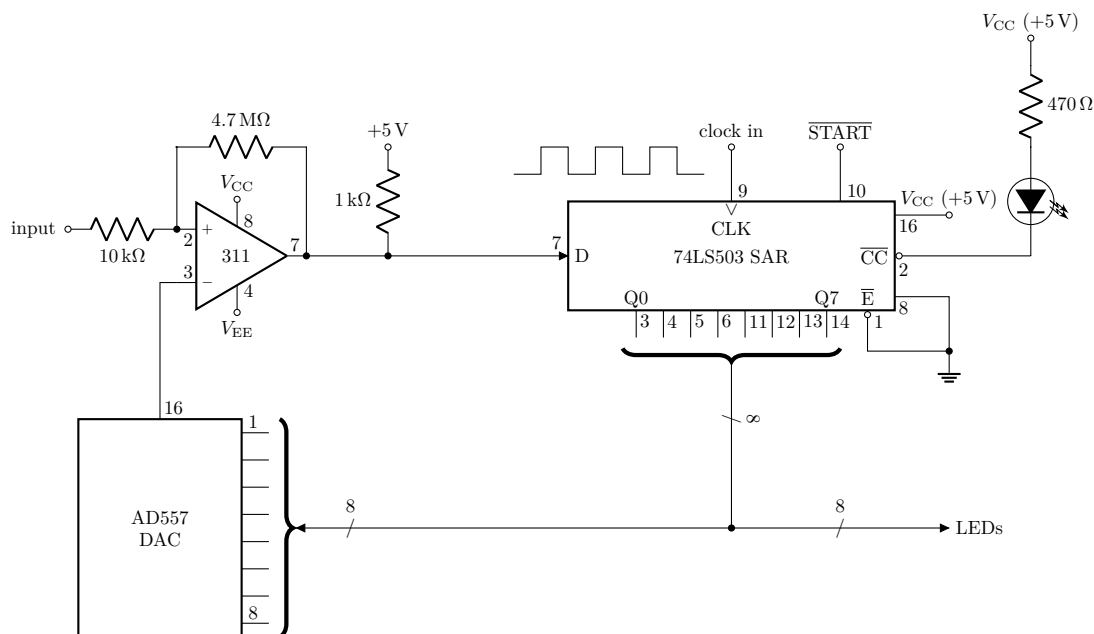
(b) Why is this ADC slower than a successive-approximation ADC? (And by how much is this ADC slower on average?)

(c) What kind of "bias" does this ADC have in terms of converting a noisy input signal?

**Problem 16.5**

(a) In the successive-approximation ADC examples from Section 16.2.2, the convention is that the SAR starts with the “midpoint” word 01111111 for an 8-bit ADC (as in the 74LS502 and 74LS503 SAR’s in Appendix B). Another possible (and reasonable) convention is to start off with the *alternate* midpoint word 10000000. Briefly describe the difference in the *end result* of the two schemes for an arbitrary input.

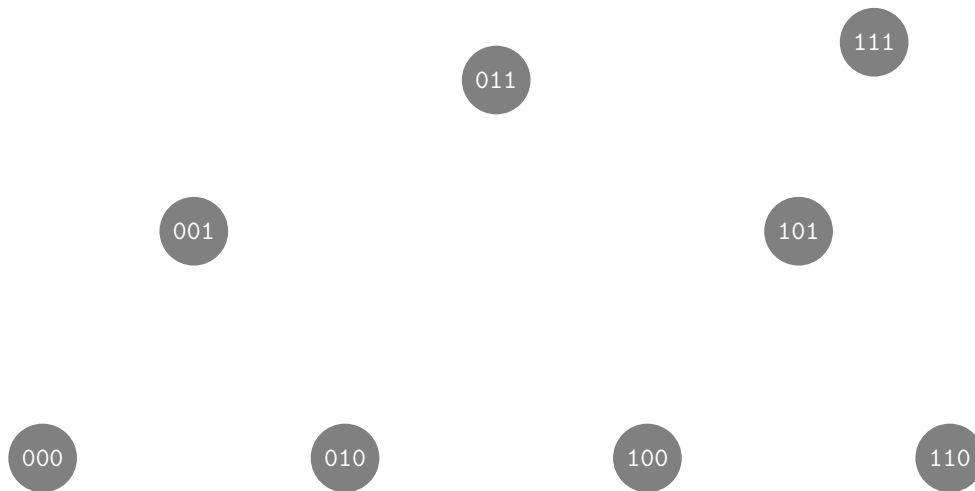
(b) Consider the ADC circuit shown below, where the DAC conversion levels are 0, 0.01, 0.02, ..., 2.55 V. Suppose the input analog voltage of 0.45 V, and as mentioned the SAR starts with 01111111 on the first clock cycle (the “start” cycle). What is the SAR output after 3 more clock cycles?

**Problem 16.6**

In this problem, you should design a 2-bit successive-approximation register (SAR), with data ( $D$ ) and start-LOW ( $\overline{S}$ ) inputs,  $Q_1$  (MSB) and  $Q_0$  (LSB) outputs, and conversion-complete-LOW ( $\overline{CC}$ ) output. On start, the outputs should initialize to  $Q_1Q_0\overline{CC} = 011$ .

(a) Draw a state diagram for the SAR, enumerating all possible output states  $Q_1Q_0\overline{CC}$  as nodes. Specify **all** possible transitions, and label transition arrows with the appropriate input states wherever multiple transitions are possible. Also, make sure to handle all possible states and eliminate the possibility that the state machine will get “stuck.” (Remember you can use “X” for “doesn’t matter” for logic states in diagrams and truth tables.)

A suggested template for your state diagram is shown below. (That is, these are the states, you should fill in the transitions.)



(b) Write down a truth table for register inputs  $D_1$ ,  $D_0$ , and  $D_{\overline{CC}}$  in terms of the other variables to implement this SAR in sequential logic.

Again, a suggested template for your solution is shown below.

$\overline{S}$	$D$	$Q_1$	$Q_0$	$\overline{CC}$	$D_1$	$D_0$	$D_{\overline{CC}}$
0	X	X	X	X			
X	X	1	1	1			
1	0	0	1	1			
1	1	0	1	1			
1	0	0	0	1			
1	1	0	0	1			
1	0	1	0	1			
1	1	1	0	1			
1	X	X	X	0			

(c) Finish the design: write down logical expressions for  $D_1$ ,  $D_0$ , and  $D_{\overline{CC}}$ .



## Chapter 17

# Phase-Locked Loops

Simply put, a **phase-locked loop (PLL)** is a feedback-loop circuit that compares two oscillating signals. It attempts to adjust the frequency of the second one so that it exactly matches the first in terms of phase (and thus also in terms of frequency).

Strictly speaking, a phase-locked loop can be implemented in an analog circuit, where, for example, the circuit makes one sine wave copy another one. However, it is common to implement phase-locked loops using digital gates, so we are covering these as digital circuits.

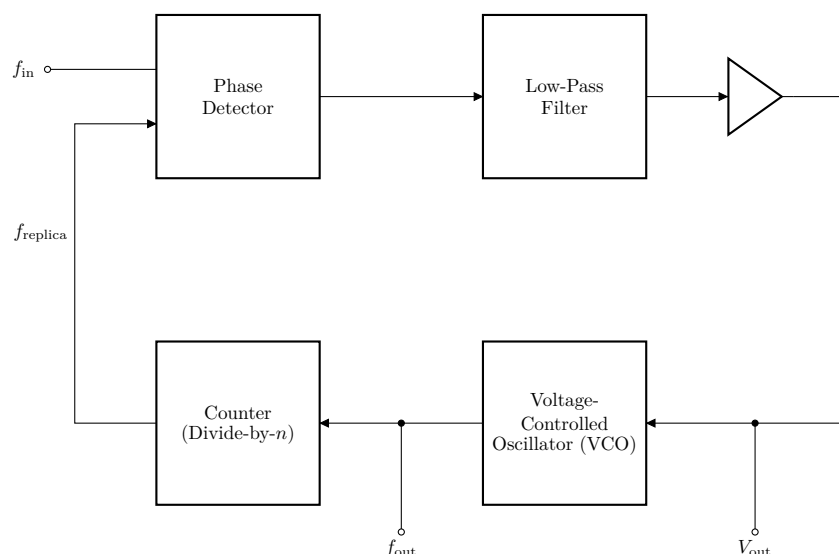
It may also sound a bit weird to use a feedback loop to make a copy of a signal, when you could just directly make a copy of a signal, e.g., with a buffer amplifier or gate. However, the magic comes in taking advantage of the feedback loop. As our first main example, recall that using a counter, it is relatively straightforward to *divide* the frequency of a square-wave clock signal. But how do we *multiply* the frequency of a signal? The answer: a phase-locked loop.

### 17.1 Frequency Multiplier

The idea behind a frequency multiplier is to start with the original clock signal. Suppose we want to multiply the frequency by  $N$ . Then generate a new signal, *divide it by  $N$* , and compare the *divided* signal to the original (i.e., phase synchronized). Adjust the frequency of the new signal until the divided version matches the original, and *voilà*, you have a new signal with a frequency  $N$  times the original, with matching phases of the two signals.

The block diagram of a circuit that accomplishes this is shown below.

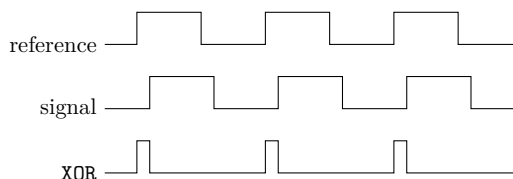




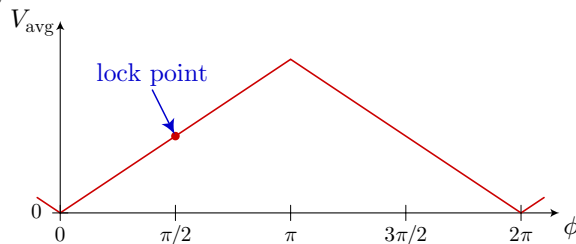
Let's look at each of the components here.

1. The **phase detector** compares two oscillating signals, and the output gives some measure of the relative phase. A general requirement is that the “in sync” state should give a “zero output”—this need not actually be zero (i.e., it could be offset to some other voltage), but the point is that the signal should go up if the phase is perturbed one way, and down if the phase is perturbed the other way. There are two basic classes of digital phase detectors in PLLs.

- **Type I phase detector.** This is simply an XOR gate, and it can be driven by digital signals, or also by analog signals, provided they have been converted to digital via a comparator or Schmitt trigger. The output of the XOR gate is illustrated below.



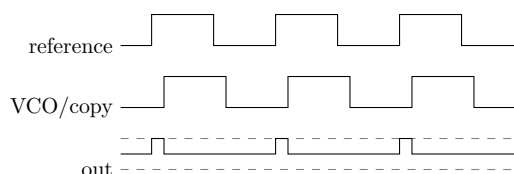
The output is a sequence of pulses; the output is HIGH whenever the two input signals mismatch. The XOR output is zero if the reference input matches the signal input, and the duty cycle of the output increases to 100% if the signals have a  $\pi$  phase difference. Only the *average* signal will matter, because the output is fed through a low-pass filter. So the average output is proportional to the phase difference, as shown below.



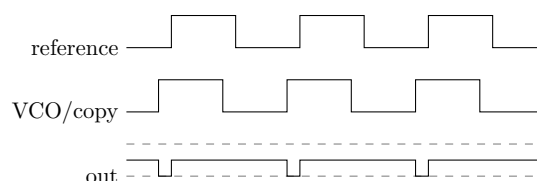
Then we should choose a  $\pi/2$  phase shift as the lock point (i.e., the PLL will force the signal to be a  $90^\circ$ -phase-shifted copy of the reference). Remember this is because we need the output signal to vary both up *and* down if the phase moves away from the lock point.

- **Type II phase detector.** This phase detector is sensitive to digital *edges*, so it is really suited to digital signals, although in principle if analog signals are converted to digital, this would amount to

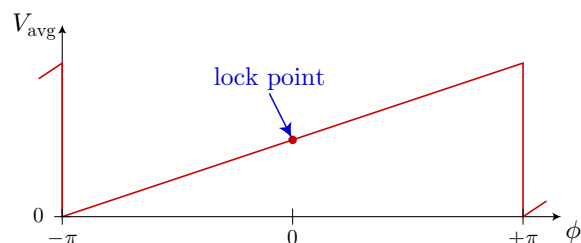
the detector responding to the zero-crossings of the analog signals. To see how it works, consider the first timing diagram below.



The “normal” state of the phase detector is a *middle* voltage midway between LOW and HIGH. If the detector finds a rising edge from the reference signal first, then the output changes HIGH. When the edge from the other signal (“VCO/copy” signal, which we will explain below), then the output changes back to MID. If the relative phase has the opposite sign, then we have the situation shown below.



Now the VCO/copy signal presents its rising edge first, so the output goes from MID to LOW. It goes back to MID when the reference edge arrives. Once this signal is time-averaged, the result is shown below.



Now the lock point is at zero phase, and because the output can move LOW or HIGH relative to the normal MID state, the output can vary in either direction. The advantage is that the lock point is in perfect sync: at the lock point, the error signal is identically zero, even before the time average. Contrast this to the type-I case, where the output was a 50%-duty-cycle square wave. Some of this will leak through the time average, and end up frequency-modulating the output signal. Note that the operation of this circuit is independent of the duty cycles of the two signals, unlike the type-I case where we assumed 50% duty cycles. (Otherwise the locked phase may differ from  $\pi/2$ , and it may be necessary to choose a different lock voltage.)

Note that for sine-wave analog signals, the phase detector can be as simple as a multiplying amplifier (for rf frequencies, you would use an **rf mixer**, which has just this function).

2. The **low-pass filter** (Section 2.3.5) “keeps” low frequencies, and “removes” high frequencies. It thus acts to time-average the phase-detector signal. It also limits the speed with which the PLL can respond to frequency changes in the reference signal. This may be a disadvantage if you want to perfectly track the frequency. However, this allows the PLL to act as a frequency “flywheel,” so it ignores some of the noise in the incoming signal to “clean it up.” It also induces a phase shift, which we will return to below.
3. The **voltage-controlled oscillator (VCO)** is an oscillator (clock), where the output frequency  $f$  depends on an input control voltage. The frequency may depend nonlinearly on the control voltage, but it should be at least monotonic.
4. The **counter** is here as we described for frequency-multiplier applications. This should be omitted in other applications.

### 17.1.1 Feedback Loop

In the feedback loop here, the time-averaged output of the phase detector is fed into the VCO control input. Recall that in PID control (Chapter 8), it is necessary to integrate the error signal in order to get zero steady-state error. Here, this is automatic, since frequency  $\omega$  and phase  $\phi$  are related by

$$\omega = \frac{d\phi}{dt}. \quad (17.1)$$

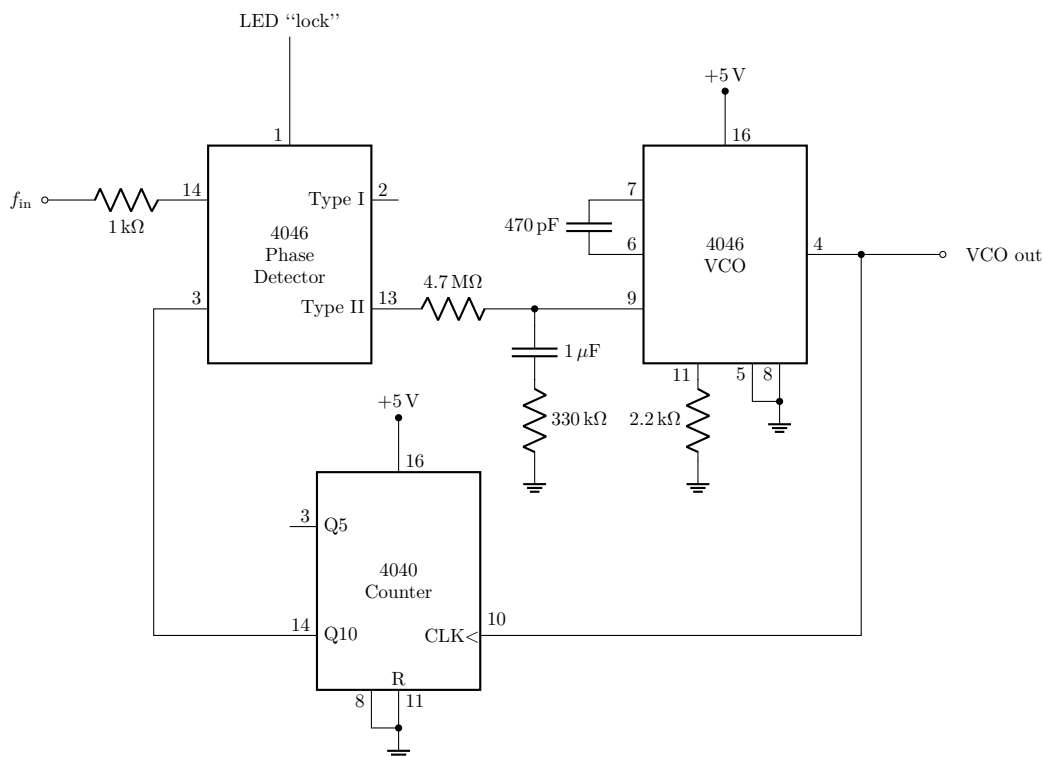
Since we detect phase and are feeding back to a frequency control, we are controlling  $\omega = \int \phi dt$ , and thus we effectively have the integral of the error signal.

If the phase is *below* the lock point (i.e., phase lag), then the output is positive (relative to the lock point), and so the frequency increases. If the phase is *above* the lock point (i.e., phase lag), then the output is negative (relative to the lock point), and so the frequency decreases. Then there are different options for loops.

1. In a **first-order loop**, there is no low-pass filter, so there is just the  $90^\circ$  phase shift associated with the phase-frequency integration. In this case, we don't have the time averaging as in the analysis above, but the idea is the same, because the integral VCO response makes the loop behave in essentially the same way.
2. In a **second-order loop**, there *is* a low-pass filter as in the diagram, so you need to be careful about any extra phase shifts to guard against instabilities. Again, the low-pass filter limits the rate of change (i.e., the bandwidth) of the control signal.

## 17.2 Example PLL

Below is a more detailed example of a PLL circuit, based on the 4046 PLL IC. This IC includes both type I and type II phase detectors, as well as a VCO. The circuit below is designed to multiply an input 60-Hz signal by  $2^5$  or  $2^{10}$ , depending on whether the  $Q_5$  or  $Q_{10}$  output of the 4040 counter is used. Note the passive low-pass filter and the use of the type II detector, so this is a second-order feedback loop.



## 17.3 Other Applications

### 17.3.1 FM Demodulation

Another important application of PLL circuits is frequency-modulation (FM) demodulation (i.e., the demodulator in an FM radio). In this case, the receiver receives an FM signal, and a PLL attempts to reproduce this signal. The FM signal was generated by changing the frequency according to some signal to be transmitted (e.g., audio). Then the control voltage to the VCO is a copy of the original signal, provided the VCO control voltage is related to frequency in the same way as in the original FM process. Typically, this just means that the VCO frequency should be linear in the control voltage.

### 17.3.2 Direct Digital Synthesis

Another application is in **direct digital synthesis (DDS)**. The idea here is to take a precision clock input (e.g., from an atomic clock or oven-stabilized crystal oscillator), multiply it to a high frequency, and then use a counter to divide it to some other frequency. This can produce frequencies with high resolution over a wide range if the frequency-multiplication factor is large. The divider is digitally programmable so the final frequency is dynamically programmable. Then the digital counter output drives an analog “look-up table” of voltages to get a high-quality (low-distortion), timing-accurate sine wave. An example is the AD985L DDS IC, which can take a 10-MHz clock in, and produce a sine-wave output in the range of 0–135 MHz.

## 17.4 Dynamical Model

To understand the behavior of a phase-locked loop in more depth, here we will develop a simple dynamical model. Consider a reference signal  $V_{\text{ref}}$ , given by

$$V_{\text{ref}}(t) = V_{\text{r0}} \cos \phi_{\text{ref}}(t), \quad (17.2)$$

where the reference frequency is

$$\omega_{\text{ref}} = \dot{\phi}_{\text{ref}}. \quad (17.3)$$

We will take this frequency to be constant, so that  $\phi_{\text{ref}} = \omega_{\text{ref}} t$ . The signal that we want to phase-lock to the reference is similarly

$$V_{\text{sig}}(t) = V_{\text{s0}} \sin \phi_{\text{sig}}(t), \quad (17.4)$$

where the signal frequency is

$$\omega_{\text{sig}}(t) = \dot{\phi}_{\text{sig}}(t), \quad (17.5)$$

or inverting this relation,

$$\phi_{\text{sig}}(t) = \int_0^t \omega_{\text{sig}}(t') dt'. \quad (17.6)$$

Note that we have already built in a relative phase of  $\pi/2$  between reference and signal, anticipating that the two signals will prefer to lock with this phase difference. Thus, the locking condition is that  $\phi_{\text{sig}}(t) = \phi_{\text{ref}}(t)$ , modulo  $2\pi$ .

The simplest phase detector for the analog signals here is a multiplier for two analog signals, which is called a **mixer** for radio-frequency (rf) signals. This is the analog equivalent of the Type-I phase detector (the XOR gate). Think of the XOR gate operating on logic state of 1 and  $-1$ . The XOR (or product) is 1 if two input signals are the same, and the XOR (or product) is  $-1$  if the two signals are opposite. We can write the output of the mixer as

$$V_{\text{mix}}(t) = \frac{V_{\text{ref}}(t) V_{\text{sig}}(t)}{V_0} = \frac{V_{\text{r0}} V_{\text{s0}}}{V_0} \cos \phi_{\text{ref}} \sin \phi_{\text{sig}}. \quad (17.7)$$

Then defining  $V_{\text{m0}} := V_{\text{r0}} V_{\text{s0}} / 2V_0$  and using the identity  $\sin \alpha \cos \beta = (1/2)[\sin(\alpha - \beta) + \sin(\alpha + \beta)]$ , the mixer signal becomes

$$V_{\text{mix}}(t) = V_{\text{m0}} [\sin(\phi_{\text{sig}} - \phi_{\text{ref}}) + \sin(\phi_{\text{sig}} + \phi_{\text{ref}})]. \quad (17.8)$$

In phase-locked-loop operation, only the first term will be important here. At or near the locking condition,  $\omega_{\text{sig}} \approx \omega_{\text{ref}}$ , so the first term varies slowly (close to dc), while the second term has a frequency of  $\omega_{\text{sig}} + \omega_{\text{ref}} \approx 2\omega_{\text{ref}}$ , which is much faster. The effect of this fast oscillation will tend to be averaged away to zero, especially as we will be feeding this mixer signal through a low-pass filter, which will greatly suppress the second term. Thus, we will write

$$V_{\text{mix}}(t) = V_{\text{m0}} \sin(\phi_{\text{sig}} - \phi_{\text{ref}}), \quad (17.9)$$

as far as the operation of the feedback loop is concerned.

To complete the loop, we must connect the output of the mixer to the input of the VCO, via a low-pass filter (to make a second-order loop). Then the VCO outputs a signal at frequency

$$\omega_{\text{sig}}(t) = \omega_{\text{s0}} - \frac{g_{\text{I}}}{V_{\text{m0}}} \int_0^t V_{\text{mix}}(t') dt', \quad (17.10)$$

where  $\omega_{\text{s0}}$  is the “natural” frequency of the VCO (i.e., the frequency with zero input voltage),  $g_{\text{I}}$  is an “integral” gain factor, and we are modeling the low-pass filter via an integral (recalling from Section 2.2.1 that a low-pass filter acts as an integrator provided the output signal is small compared to the input). Note that we have included a minus sign here to provide negative feedback. However, this is optional, but without it the relative phase at the lock point will differ by  $\pi$  from the analysis here.

### 17.4.1 Equation of Motion

Now we can write down a dynamical equation for the phase-locked loop. Consider the phase difference

$$\Delta\phi(t) := \phi_{\text{sig}}(t) - \phi_{\text{ref}}. \quad (17.11)$$

The first derivative gives the frequency difference

$$\Delta\dot{\phi}(t) = \omega_{\text{sig}}(t) - \omega_{\text{ref}}. \quad (17.12)$$

Putting in Eq. (17.10) for  $\omega_{\text{sig}}$ ,

$$\Delta\dot{\phi}(t) = \omega_{\text{s0}} - \omega_{\text{ref}} - \frac{g_{\text{I}}}{V_{\text{m0}}} \int_0^t V_{\text{mix}}(t') dt'. \quad (17.13)$$

Differentiating this equation,

$$\Delta\ddot{\phi}(t) = -\frac{g_{\text{I}}}{V_{\text{m0}}} V_{\text{mix}}(t), \quad (17.14)$$

and then using Eq. (17.9), the result is

$$\Delta\ddot{\phi}(t) = -g_{\text{I}} \sin(\Delta\phi). \quad (17.15)$$

This equation has the form of a mechanical pendulum,  $\ddot{\theta} = -(g/\ell) \sin \theta$ . This means that  $\Delta\phi = 0$  (modulo  $2\pi$ ) is a steady state, meaning that once locked, the circuit can stay locked. However, there is no means to *become* locked: if the relative phase is displaced from the lock point, it will oscillate back and forth about it without settling. Thus, we need to introduce some *damping*.

### 17.4.2 Damping

To introduce damping, we will introduce a “proportional” term in the feedback in Eq. (17.13):

$$\Delta\dot{\phi}(t) = \omega_{\text{s0}} - \omega_{\text{ref}} - \frac{g_{\text{I}}}{V_{\text{m0}}} \int_0^t V_{\text{mix}}(t') dt' - \frac{g_{\text{P}}}{V_{\text{m0}}} V_{\text{mix}}(t). \quad (17.16)$$

Here  $g_{\text{P}}$  is the “proportional” gain. This model corresponds to the two-resistor, one-capacitor filter in the example PLL circuit on p. 388. At low frequencies, the filter acts as an ordinary low-pass filter (integrator), while at high frequencies the capacitor acts as a short, so the two resistors form a divider that determine  $g_{\text{P}}$ .

Now putting this feedback into Eq. (17.12), the result is

$$\Delta\ddot{\phi}(t) = -g_i \sin(\Delta\phi) - g_p \Delta\dot{\phi} \cos(\Delta\phi), \quad (17.17)$$

in place of Eq. (17.15). The new equation has the form of a damping term. We should compare this equation of motion to that of the damped pendulum,

$$\ddot{\theta} = -\frac{g}{\ell} \sin \theta - \frac{\gamma}{\ell} \dot{\theta}, \quad (17.18)$$

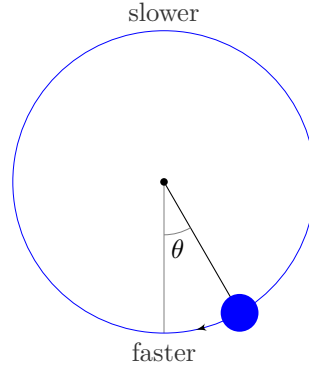
where  $\gamma$  is the damping rate (i.e., the last term here applies a torque that opposes the angular velocity, slowing the pendulum). The PLL equation (17.17) has a similar damping term, but modulated by the cosine of the relative phase. It isn't completely clear that this helps: the sign of the damping changes depending on the phase, and a negative damping is no damping at all (it tries to *speed* up the rate of phase change).

If we look close to lock, then this equation works out. That is, suppose  $\Delta\phi$  is small. Then  $\sin(\Delta\phi) \approx \Delta\phi$ , and  $\cos(\Delta\phi) \approx 1$ . Thus, Eq. (17.17) becomes

$$\Delta\ddot{\phi}(t) \approx -g_i \Delta\phi - g_p \Delta\dot{\phi}, \quad (17.19)$$

which is the equation for a damped harmonic oscillator. Thus,  $\Delta\phi$  will settle to zero (or a multiple of  $2\pi$ ), possibly exhibiting damped oscillations about the lock point along the way.

If the circuit is far from lock, then  $\Delta\phi$  varies rapidly, and it isn't clear that the damping helps, because  $\cos(\Delta\phi)$  seems like it should average to zero. The key to understanding how the damping works is to consider a mechanical pendulum, as shown below, rotating at high angular velocity.



Because the kinetic energy is lower when the pendulum is going “over the top,” the angular velocity is lower at the top than at the bottom. In the phase-locked loop, this means  $|\Delta\dot{\phi}|$  is larger when  $\Delta\phi$  is near  $0, \pm 2\pi, \pm 4\pi, \dots$ , while it is smaller when  $\Delta\phi$  is near  $\pm\pi, \pm 3\pi, \dots$ . We can model this by writing (assuming  $\Delta\dot{\phi} > 0$ )

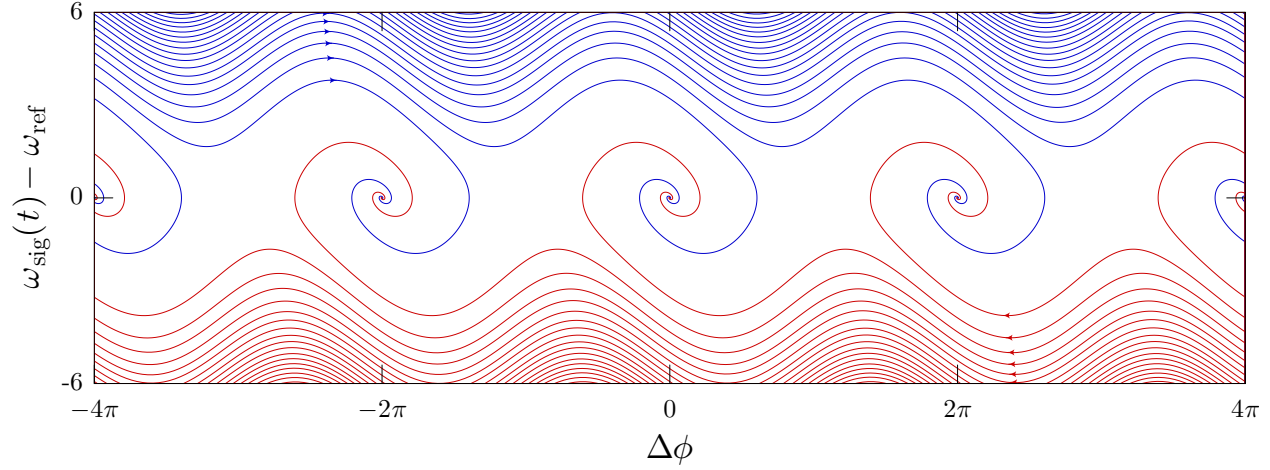
$$\Delta\dot{\phi} = (\omega_{\text{sig}} - \omega_{\text{ref}}) + \delta\omega \cos(\Delta\phi), \quad (17.20)$$

where we assume  $(\omega_{\text{sig}} - \omega_{\text{ref}})$  to be slowly varying (i.e., constant on the time scale that  $\Delta\phi$  rotates through  $2\pi$ ), and  $\delta\omega$  is small. Then the damping term from Eq. (17.17) becomes

$$-g_p \Delta\dot{\phi} \cos(\Delta\phi) = -g_p (\omega_{\text{sig}} - \omega_{\text{ref}}) \cos(\Delta\phi) - g_p \delta\omega \cos^2(\Delta\phi). \quad (17.21)$$

Since  $\Delta\phi$  changes rapidly, the first term on the right-hand side will average to zero (because the average value of  $\cos x$  is zero), while the second term on the right will average to  $-g_p \delta\omega/2$  (because the average value of  $\cos^2 x$  is  $1/2$ ). This leads to a net damping that opposes  $\Delta\dot{\phi}$ . Recalling that  $\Delta\dot{\phi}$  is the relative frequency  $\omega_{\text{sig}}(t) - \omega_{\text{ref}}$ , this means that the damping pushes  $\omega_{\text{sig}}(t)$  towards  $\omega_{\text{ref}}$ , until it is close enough that the phase-locked loop can “capture” the signal, and then the phase stabilizes. This damping effect becomes smaller the further  $\omega_{\text{sig}}$  is from  $\omega_{\text{ref}}$ , so if the signal frequency is initially highly mismatched, it may take a while for the loop to attain lock.

The process is illustrated in the plot below, which shows the relative frequency  $\Delta\dot{\phi}$  plotted against the relative phase itself, from numerical solutions of Eq. (17.17), with  $g_p = 1$  and  $g_i = 2$ . The trajectories starting far away from the correct frequency are pushed towards the correct frequency, with “bumps” in the frequency along the way, and the capturing process is evident, where the trajectories settle down to a point (the lock point).



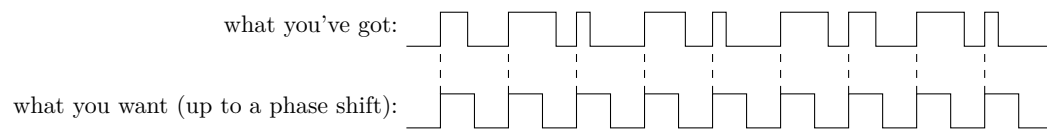
## 17.5 Exercises

### Problem 17.1

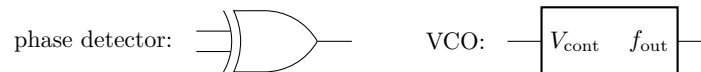
Recall that a phase-locked loop based on a Type I detector is sensitive to the duty cycle of the input signal.

(a) **Briefly**, why?

(b) Suppose you have a signal consisting of a train of (digital) pulses. The *rising* edges occur at regularly, at a well-defined frequency. The *falling* edges, however, occur at irregular times. Show (**draw a schematic and describe your reasoning**) how to use two flip-flops (pick your favorite type) as divide-by-2 counters with a phase-locked loop to create a “cleaned” version of the same signal (i.e., square wave, 50% duty cycle, same frequency as the rising edges, “ignores” the falling edges). Note that the actual output may be phase-shifted compared to what is shown here.



You can use the following schematic symbols for the PLL components in your solution ( $V_{\text{cont}}$  = control voltage;  $f_{\text{out}}$  = oscillator output signal):



You can assume the circuit will work without a low-pass filter.





# Appendices



## Appendix A

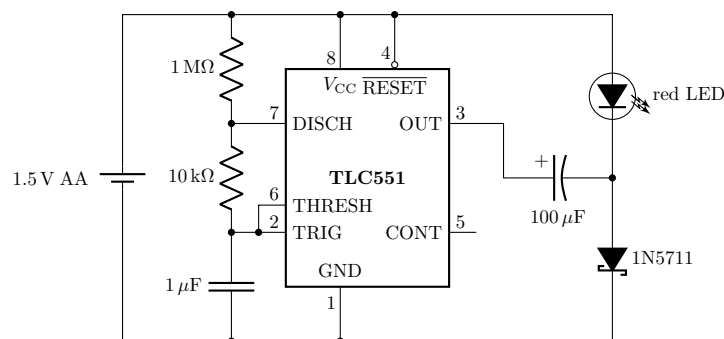
# Homemade Printed Circuit Boards

Although it is not so difficult these days to have a printed circuit board (PCB) fabricated professionally, it is useful to know how to fabricate one yourself (e.g., if you need to build a circuit in a *tout de suite*). This appendix is a guide to fabricating a circuit prototype in the form of a printed circuit board (PCB). As an example, this guide gives as an example a simple, battery-powered, 555 LED blinker circuit on a PCB, and then understand and characterize its operation.

Because this guide serves as a reference for fabricating PCBs on your own, it contains a lot of information (including things like mixing chemicals). As with any involved project where there are critical steps, **read all instructions in this guide *before* you tackle PCB fabrication.**

### A.1 LED Blinker Circuit

The circuit you will construct is shown in the schematic below.



This is a standard 555 astable multivibrator circuit (see Section 15.1), but with some adaptations to make it run and drive an LED on a 1.5 V battery

- The TLC551 is a CMOS variant of the 555, which operates with a power supply down to 1 V.
- A red LED (here, Vishay TLDR5400) typically requires a 1.8 V forward-voltage drop (up to 2.2 V). This is too much for an AA battery to light the LED directly. (Actually, the LED will light with 1.5 V, but it will be quite dim.) Hence the addition of a capacitor and diode in the “voltage doubler” output network. How does this work? (If it’s not obvious, you’ll figure this out later.)
- The 1N5711 is a Schottky-barrier diode, which in this application is desirable because it has a low forward-voltage drop (about 0.25 V at 20 mA). However, with a fresh battery this circuit also works with a normal diode (e.g., 1N914).

Look through the schematic and understand how all the components are connected. If there is anything you’re not sure about, you can always breadboard the circuit *before* you solder it!

In the case that you aim to fabricate this circuit using the PCB design below, a more detailed parts list is:

- IC: TLC551 (CMOS, 1 V 555).
- Resistors: 10 k $\Omega$ , 1 M $\Omega$  (1/4 W).
- Capacitors: 1  $\mu$ F ceramic (TDK FK18X5R1E105K), 100  $\mu$ F electrolytic (Panasonic ECE-A0GKS101).
- Diodes: 1N5711 (Schottky), red LED (Vishay TLDR5400).
- Misc: AA battery; AA battery holder (Keystone Electronics 2460); presensitized PC board; PCB fabrication equipment and chemicals as detailed below.

## A.2 PCB Fabrication

The standard for permanent or semi-permanent prototyping is to fabricate a **printed circuit board (PCB or PC board)**. Fabrication proceeds in a few steps, but begins with a **copper-clad board**, basically a sheet of FR-4 fiberglass (typically 1/16" thick) with a thin layer of copper plated on one or both sides (referred to as "single-sided" or "double-sided" boards, respectively). The most common copper thickness is "1 oz.," which means 1 ounce of copper per square foot of board (per side); this translates to a thickness of 1.344 mils (1 mil = 0.001").

The first step is to remove unwanted copper, leaving behind a set of connections in the copper foil between the components to be placed. This is generally done via either mechanical or chemical means. The mechanical method involves using a tool, such as a high-speed router (something like a Dremel tool) to remove copper. This would be hard to do by hand, but there are computer-controlled routing machines for this purpose. They are noisy and expensive, but handy.

The far more common method is to etch unwanted copper away using nasty chemicals, and that is what we will do. The copper areas of the board to be kept must be protected by some covering, and the board is submerged in etching chemicals, removing the unshielded copper. The protection is not difficult; even a mark with a (permanent) Sharpie marker<sup>1</sup> is enough to inhibit etching, and in principle you can draw your entire design by hand right on the board. In practice it's hard to do this neatly, and so there are some more elegant tricks for laying out your PCB design. For a one-off board you can buy rub-on transfer patterns for electronics. These kits include pads for ICs and other components, as well as lines for connections between components. Thin rolls of black masking tape are also available for making lines on the board.

Using computers, you can more quickly make sophisticated prototypes. Basically, you draw out the copper pattern in most any drawing program (it helps to have one where things can be positioned precisely; CAD programs are great, Illustrator can work), or there are special layout programs for this purpose (Osmond PCB is a commercial program only for OS X, but it is very well done; Eagle is commercial and popular and runs on just about anything; gEDA is an open-source suite<sup>2</sup>). To get the designs onto the copper, there are two main methods. One is to print the design onto special "iron-on" paper using a laser printer. The idea is that when you press it onto the circuit board and run a hot clothing iron over it, the toner releases and transfers to the copper surface (note that you have to print out the mirror image of the design onto the paper). I haven't had much luck with this method; basically extra material on the paper also transfers over, and you have to spend a lot of time carefully cleaning up the board with an X-acto knife.

The most reliable method for getting a computer design onto a board is to use a photographic technique, which we will cover below.

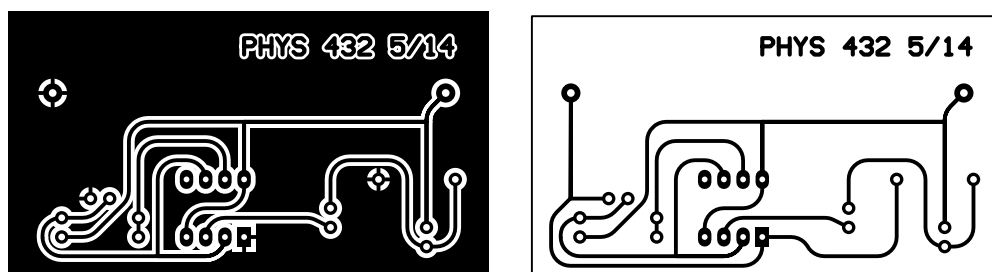
### A.2.1 LED-Blinker PCB

The PCB that we will fabricate is shown below, printed at actual size. This is a view **from the bottom of the PCB** (i.e., the "trace side" or "copper layer;" this will be a single-sided PCB, meaning there will only

<sup>1</sup>actually, Sharpie's are adequate under the right etching conditions; a much better (more robust) choice is the "Staedtler Lumocolor Permanent Special 319," and also the old "TDK CD-Writer Permanent Marking Ink," if you can still find them.

<sup>2</sup>see <http://www.evilmadscientist.com/2007/laying-out-printed-circuit-boards-with-open-source-tools/>

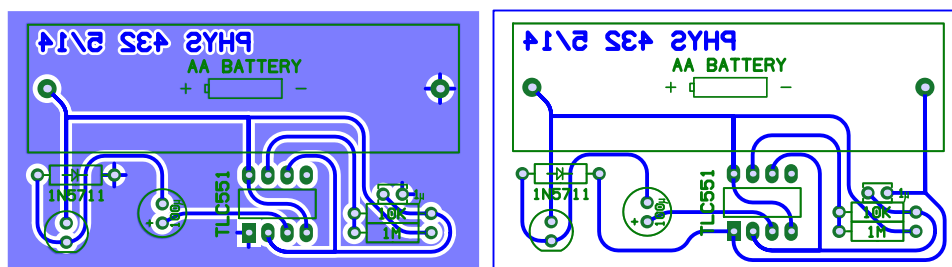
be copper on one side). Anywhere dark is where the copper will stay in the completed circuit board.



Note that there are larger “pads,” where component wires will be soldered, with holes to help guide you when drilling. The pads are connected by paths, or “traces.” You should also be able to pick out where the TLC551 (IC, 8-pin DIP) goes; note that one pin has a square pad, in contrast to the other oval or circular pads; the square pad marks the location of pin 1. Note also that there is lettering in the copper layer; this is a good idea to keep you from accidentally flipping the exposure mask and fabricating a backwards board.

There are two different but equivalent boards shown here. The one on the right is a bit more obvious, and makes straightforward connections between pins, and is shown for clarity. The one on the left uses a **ground plane**. This means that the unused parts of the board are “flooded” with copper, which is connected to ground (battery negative terminal). Other ground connections just connect directly to the ground plane, instead of running traces between the ground connections. (Note that for these connections there are “thermal pads” with connections to the ground plane instead of just a drill hole; this is to make soldering tidier and easier.) In this PCB, the main advantage of the ground plane is to save etchant and speed up the etching process, because we don’t need to remove nearly as much copper. In high-speed and low-noise circuits, the ground plane acts as something of a shield for interference and crosstalk between signal lines. So while a ground planes usually require specialized software to draw easily, it is a useful technique.

The components are not marked in this particular PCB; while it is possible to put fine lettering in the copper layer, it doesn’t necessarily develop and etch well unless you have the process down pat. So that you can sort out the component locations when soldering, here is a view **from the top of the PCB** (i.e., the “component side,” the side on which you will “stuff” the components), with component labels.



Note that the copper layer is mirrored from this point of view; again, you can tell from the text. In a professionally fabricated board (i.e., not the kind you do yourself), these labels would be silk-screened onto the component side over a (typically green) protective layer, called the “solder mask.”

**Take care to get component orientations right when you are soldering. Removing soldered components without damaging the PCB or component is a hassle.** As the adage goes, check twice, solder once.

A couple of notes on recognizing the orientations of components in the LED blinker board:

- The electrolytic capacitor (black cylinder) is *polarized*—it only works in one direction. Note that the *positive* terminal is marked in the PCB diagram, but the manufacturer marks the *negative* terminal on this particular cap.
- The diode has a black band on one end that marks the *cathode* (negative terminal). This is the terminal towards which the “arrowhead” points on the diode symbol on the PCB diagram.

- The LED cathode is marked by a flat spot on the edge, as shown in the PCB diagram.
- We already discussed pin 1 on the IC. The other components can go in either way, but in general you should stuff them in the *same* direction so they are readable (in this board, mostly for the resistors—make the bands “read” from left to right). On the subject of resistors, you should be able to tell the resistors apart by looking at the bands; if you can’t do this, now is a good time to look up how to read resistors!

### A.2.2 Presensitized PCBs

**Presensitized** copper-clad boards come with a layer of *photoresist* coated onto the copper. The photoresist is sufficient to inhibit etching. With a *positive* photoresist, you make an optical mask that is opaque where you want the copper to stay and transparent where you want the copper to etch away. Then you place the mask over the presensitized surface, and expose the masked board to UV light. Then you can chemically “develop” the board, such that the resist is removed wherever it was exposed, and you can proceed to etch away the copper. There are also boards with negative photoresists (i.e., you expose the parts you want to keep), but positive boards are simple, cheap, and easily available.

Some sources of positive, presensitized boards that work well are:

- Ever-Muse boards, available from Circuit Specialists.<sup>3</sup>
- MG Chemicals boards, available from Circuit Specialists or Mouser Electronics.<sup>4</sup>

These boards are available either single- or double-sided, and in various sizes. The relative prices of the two brands change often; it’s fine to buy whichever is cheaper at the moment.

### A.2.3 Cutting to Size

In general, the presensitized board won’t be the same size as your desired finished circuit. Obviously, you should start out with a larger board than you need, and then cut it down. If you are only using a small part of a presensitized board, then now is the time to cut it, so you can save the rest of the board to use later. If you only need to cut down the dimensions a bit, or you are fabricating multiple copies of the same board, you can defer cutting until after the board is developed and etched.

One problem with cutting a board with something like a band saw is that the fiberglass is very abrasive and hard on the saw blade. This also generates fiberglass dust, which is harmful to breathe long-term. If you go this route, make sure to wear a dust mask and protect your lungs.

A more convenient solution is to get your hands on a very old, heavy-duty paper cutter (the green kind with a large, hinged blade; find these at garage sales or flea markets for cheap). This kind of cutter goes right through PCBs without generating significant dust. Newer ones may work, but you really need a sturdy, hinged blade. *Don’t* go slicing off your fingertip—be *careful*.

Note that whichever method you use, you should leave the protective plastic layer on the board while cutting, in order to protect the delicate photosensitive coating beneath. No matter how careful you are, the photosensitive coating will likely be damaged near the cut edge. This may require some repair before etching, or better yet cut the board a bit oversize, so the damaged area is away from any critical areas of your PCB layout. Oversize-cutting is a good idea anyway to make sure any critical traces on the PCB are well away from any edge.

### A.2.4 Exposure

In the exposure process, you subject certain areas of the photosensitive layer to UV light. That means you want to be careful *not* to subject the layer to any *inadvertent* UV light. That means, wait to peel off the protective plastic until right before you are ready to expose, and only handle the exposed board under fairly dim, incandescent lighting (no white-LED or fluorescent lighting).

---

<sup>3</sup><http://www.circuitspecialists.com>

<sup>4</sup><http://www.mouser.com>

You will need a **positive exposure mask**, which shows your “artwork” (PCB design). Typically, this is a transparency printed via a laser printer, and is dark wherever you want the copper to stay. For good results, it is often necessary to “stack” two transparencies to make sure the black-printed areas are dark enough. Align the two transparencies together and use some double-sided tape to hold them together in proper alignment.

Now peel off the protective plastic layer of the PCB, starting at one corner and slowly peeling it by pulling the corner towards the opposite corner. You should see a green, shiny layer—the photosensitive layer. **Don’t touch the green layer**; it’s delicate. Lay the board, green side up, on the table where you will do the exposure. Lay the exposure mask on top of the board, and align it appropriately (e.g., you may want to put the edges of the mask against the cleanest edges of the board).

The mask must stay flat and pressed against the board during the entire exposure. A sheet of plexiglass or window glass works fine for this. Place the glass on top of the mask, and double-check the mask alignment.

Now you’re ready to expose. This can be done with a specialized UV light, but a simple desktop fluorescent light works well. The board instructions suggest 10–15 W of fluorescent light,  $5 \pm 1$  cm from the board, for 8–10 minutes (to compensate for a larger light–board distance, triple the exposure time for every doubling of the distance). The optimal time depends somewhat on the cover glass and the quality of the transparency mask, but in practice this recommendation is excessively long. A good example choice for a lamp we will use employs a 15-W fluorescent bulb, held about 13 cm above the table top. For this setup, an 8-minute exposure is appropriate, using an object to get the board within about 2” from the light.

Underexposure is bad, as not all of the photosensitive layer will be removed from the board. Small gaps in the copper may not develop properly.

Overexposure is also bad, as the copper traces will come out too thin, and may be broken in some areas. Of course, if extremely overdeveloped, there won’t be any photoresist left at all. This process is not especially sensitive to the exposure time, but it’s important to get it within a reasonable range.

If you are doing a double-sided board, you should now proceed with exposing the other side of the board. These are a bit tricky, as it is hard to keep the artwork for the two sides in alignment with respect to the board. If they’re not well-aligned, you’ll see exactly what the problems are when you start drilling the board. The trick here is to tape the front and back layers together in proper alignment *right at the beginning* (e.g., with double-sided tape). Then slide the board in between the two masks (with protective plastics on both sides removed), and proceed with exposure. The only hard part is making sure the board doesn’t move with respect to the masks when you flip it over to do the other side.

After exposure, if you look at the board under bright (incandescent) light, you may see that the exposed areas have changed to a slightly different shade of green, and you can see a faint “imprint” of the artwork. Don’t waste time now; you should proceed to developing your board immediately.

### A.2.5 Developing

The board is developed in a room-temperature bath of NaOH solution. The bath is not especially caustic, but you should still don a nitrile glove.

Developer solution is available from MG Chemicals (via Mouser or Circuit Specialists), which you dilute according to the instructions. However, this is a simple NaOH solution, at 1% concentration by weight (i.e., mixed into tap water, 10 g NaOH per L of solution). For safety, add the NaOH to the full amount of water to avoid inadvertently dealing with a concentrated, caustic solution. This solution must be thoroughly mixed before use, otherwise you might get uneven results when you develop your board. After developing a board, this solution seems to go flat after a few months (probably due to air exposure). In any case, don’t use old developer unless you know it to be fresh enough; it’s better to make up a new batch for reliability.

The developer is easiest to handle in a wide, polyethylene or polypropylene container (e.g., a wide tupperware-type food-grade container, labeled for non-food use!). NaOH reacts with some metals, so metal containers are out, and glass containers are okay but susceptible to breakage.

Using your gloved hand, dip the board into the developer bath, and gently move it around to keep fresh developer in contact with the board. You should see photoresist beginning to “smoke” off the board immediately, and the exposed areas should develop to shiny copper in less than a couple of minutes. Keep



developing for 30 seconds after it looks like it's done, otherwise there may be a very thin layer of photoresist on some of the shiny areas that will inhibit etching. Then rinse off the board. **Be careful not to drip developer all over the place.**

Carefully inspect the board. Are there areas that didn't develop properly? You may need to "edit" the photoresist layer by carefully scraping bits away using an X-acto knife. Are there areas that have overdeveloped, flaked away, or are there defects (pinholes, etc.) in any of the traces? You can "repair" photoresist by using a permanent Sharpie marker to make small edits. Also, before proceeding to the etching step, if there are large swaths of exposed copper, say near an edge of the board that you don't care about (e.g., because you started off with an oversized board), you should mask off that part of the board (with masking tape) so you don't waste etching solution on things that don't need to be etched.

## A.2.6 Etching

For etching the copper, there are two commonly used chemicals. The first is **ferric chloride**. We will avoid this, because it is pretty nasty stuff that looks like toxic sludge, and it will stain your skin and clothing. Otherwise, it works just fine. You can usually buy it as a concentrated solution or as a powder. Both require mixing with water. Be careful with power in the anhydrous form, as it is very exothermic while dissolving in water—start with the full amount of cold water and add the powder slowly and carefully to the water. Again, a polyethylene container is best for etching and storage, and glass can be used for longer-term storage.

The etchant we will use is **ammonium persulfate**, available from MG Chemicals (P/N 410-1KG, via the usual suspects), in 1 kg bottles of dry crystals. (Sodium persulfate is another, similar alternative etchant.) This is mixed in a ratio of 250 g of crystals into 1 L of tap water, and interestingly is endothermic when you mix it. This chemical has the advantage of being somewhat less nasty than ferric chloride, and it's transparent, so you can see the board while it's etching. (It also works less well than ferric chloride in terms of the Sharpie trick.) However, you have to be careful with it. **Never store the etchant in a metal container. Never use a tightly sealed container (glass or otherwise) for the etchant for long-term storage.** Disobey these and you risk an explosion—not fun. A tall polyethylene container is good for etching, because it's inert and shatter-resistant, and taller containers meant for cereal are good because they have tops that will pop open in case of pressure build-up.

And it should go without saying: **Don't dump these chemicals down the drain. Dispose of them properly.**<sup>5</sup> In particular, ammonium persulfate should not be mixed with acids, alcohols, or flammables (or anything else, just in case) for disposal. Also, in case of spilled etchant, wipe it up with **wet** paper towels (and make sure they *stay* wet)—etchants are oxidizers, which can pose a fire hazard.

Any etchant works faster at an elevated temperature (40°C is plenty), but room temperature is fine if you're patient. Agitation is also essential to a good etch; this can include periodically change the orientation of the board and moving the board around, stirring the bath with a stir plate, or bubbling air through the bath with an aquarium pump.

Make sure to wear nitrile gloves when handling etchant or boards going into the etchant. Dip the board into the etchant with your gloved hand, until the shiny copper areas turn dark and reddish. Pull the board out and look for any remaining shiny copper spots. These are areas where there is still a thin layer of photoresist, and you must fix these. After rinsing the board in water, you can scrape these areas with a knife. Alternately, you may have success with rubbing the area with a dry paper towel—if you do this carefully, this will abrade the thin layer without damaging the thicker areas that are supposed to stay.

When the board is etching correctly, place it in the etchant tank in a way that makes it convenient to retrieve with a gloved hand. You could hang it using some fishing line, or in the lab we have a plastic "rack" for boards with a handle that makes it easy to retrieve. Gently agitate the etchant once in a while—this will speed up the etching, which will take some time (10–30 minutes).

<sup>5</sup>In Lane County, Oregon, for example, household hazardous waste is accepted (at no charge) at Glenwood Central Receiving Station: <http://www.lanecounty.org/departments/pw/wmd/hazwaste/pages/hhwcc.aspx>. On UO campus, hazardous waste is collected by Environmental Health & Safety: <http://ehs.uoregon.edu>.

### A.2.7 Drilling

Before proceeding with drilling, check to see if your board has any sharp copper edges that might cut you. If so, sand them off as described in the next section.

Drilling a PCB is best accomplished with a small drill press. A reasonably good yet economical choice is a drill-press mount for a Dremel tool.

Another good, economical choice for drilling in this setup comes in the form of miniature, high-speed-steel drill bits, available for cheap from Harbor Freight. The typical hole size for small components is 0.8 mm (larger for high-power or large components), but a somewhat larger drill bit can be used to reduce drill breakage. On the example PCB you are drilling, you need a larger drill bit (a mm or a bit over) for the battery-mount holes. These bits fit directly into the collet of the Dremel tool (i.e., no drill chuck is necessary).

Some safety/operation notes:

- Small drill bits (especially carbide bits, but steel bits too) are brittle, and lead to small bits of sharp, flying shrapnel when they break. **Always wear eye protection when drilling PCBs. Flying shrapnel piercing cornea = BAD.**
- Avoid breaking drill bits in the first place. When drilling, use **gentle pressure** on the bit. If you notice any flex in the bit, **STOP** and figure out what's wrong.
- The main cause of bit breakage is movement of the PCB while drilling. Hold the PCB **firmly** against the base of the drill press to avoid (1) general side-to-side movement and (2) upward movement when the bit snags on the way out the bottom side. In general, keeping the piece in place is good drilling technique; while this is unlikely with small bits, a snagged drill bit can result in a dangerous spinning object if you don't have it firmly planted.
- Run the Dremel tool at the low end of its speed range.
- After you finish, rinse off the PCB to get rid of dust. (Drilling doesn't generate a whole lot of dust, but it's best not to blow it into the air.)
- Double-check that every hole is drilled; if you discover an undrilled hole after you have half the board stuffed, it's awkward (to say the least!) to go back and fix it. Work in a pattern (e.g., left-to-right, then front-to-back) to reduce your chances of missing a hole.

### A.2.8 Cleaning Up the Board Dimensions

Now is the time to get the outer dimensions of the board down to size. If you have large areas to trim off, go ahead and do this with a cutter, as described above. To trim down a side by a couple of mm or so, use sandpaper. Wet the sandpaper to eliminate airborne dust (use "wet/dry" sandpaper designed for this). Lay the sandpaper on a flat surface and rub the PCB against it, in the direction of the PCB, until the dimension and surface quality are up to your standards.

Non-rectangular board outlines are also possible, but require other tools such as routers or coping saws—with appropriate dust control, of course.

### A.2.9 Preparing to Solder

Now is the time to remove the green photoresist layer from the PCB. This is easy to do by rinsing off the board with acetone or methanol. Be careful not to get solvent on the table or other things, as they can damage a wide variety of materials. Use a glove to protect your hand from solvent and photoresist.

At this point the exposed copper should be bright and shiny. If not, rub it with steel wool to remove any oxides. Make sure not to touch the board with bare skin, as your skin oils will oxidize and etch the surface.

### A.2.10 Soldering

Soldering is a delicate balance between too little and too much heat. Too little, and the solder won't stick to the PCB and component you're trying to connect (in the worst case, you can get a **cold-solder joint**, which *looks* like a good solder joint, but forms a poor and/or intermittent connect—very hard to debug. Too *much* heat, and you can fry delicate components, or cause the copper foil to delaminate from the PCB (difficult to fix).

The key to successful soldering is to be *fast*—heat things up, make the joint, and then let it cool down. This way you can solder at an appropriately high temperature, and cool the joint down before heat travels enough to damage any components. Here are the details, step-by-step.

- You'll need a soldering iron. Most things will work, from a \$10 soldering pen to a regulated soldering station (the kind with a temperature control on the base). Let the iron heat up for a few minutes until it equilibrates. Solder should melt immediately when you touch it to the iron.
- You also need solder, obviously. Lead-based solder is far easier to work with than lead-free. You should also use **rosin-core** solder (never use acid-core solder on electronics). This solder contains rosin-based **flux**, which chemically removes oxides from metal to be soldered when heated. It helps the solder “wet” the bare copper and the component lead.
- To mount a component, begin by inserting it into the PCB, double-checking its position and orientation, as de-soldering components is a pain. For things with long leads like resistors and capacitors, you can bend the leads a *bit* (no more than 30°—less is better, in case you ever need to replace the component). For things like ICs with short leads, they will ideally stick in place if you've used appropriately sized holes. If not, you can use some masking tape to hold them in place while you solder. A quicker solution when you are more experienced is to use one hand to hold the board, with one finger holding the IC in place. Then quickly solder pins on opposite corners with your other hand, by touching the solder with the iron and quickly transferring it to the solder joints. This is bad soldering technique (the joints won't be properly fluxed, and you can't use very much heat if you don't want to burn your finger), but you can use these temporary joints to hold the IC while you solder the rest of the pins with both hands. Then come back and redo the temporary joints.
- To make a solder joint, begin by touching the iron to *both* the pin and the pad that you're joining. Lay the soldering iron down for the largest possible contact area—the key is heating things up *quickly*. In fact, for this reason, large flat soldering-iron tips are preferable to needle-shaped tips, even for fairly fine solder work. The metal should be hot within a second or two. Now to bring in the solder; you can work with a small piece, or hold the whole roll in your hand for stability. Touch the end of the solder wire to the iron tip. As soon as it begins to melt, move the end of the solder wire from the tip to the joint itself, and start feeding the wire into the joint. (That is, you should be feeding into the intersection between the component and copper foil, not into the soldering-iron tip.) If the wire doesn't feed because it's not melting, it means you don't have enough heat at the junction. Feed the solder until you have enough solder to make a conical blob, pointing along the component wire. As soon as you have enough solder, remove the wire and iron, and blow on the joint to cool it off quickly, to avoid overheating anything. Incidentally, it's good to inhale *before* touching the solder to the iron, and exhale after the joint is done, so you don't breathe the rosin smoke.
- Check your joint: It shouldn't bulge like a sphere, but it should be enough solder to “fill out” the junction. The solder should be nice and shiny. If not, heat it back up and try to correct it; you can use the iron tip to “wipe away” excess solder (you can then use a wet sponge or paper towel to clean excess solder from the iron), and then add fresh solder to make a good joint.

### A.2.11 Stuffing in Stages

When stuffing (inserting components and soldering) a PCB, there are some good rules to follow to save you from later headaches.

- It's easiest to stuff small components (resistors, ICs) first, and then larger components (large capacitors, inductors) later—once the big boys are in place, it's harder to maneuver.
- Even more importantly, you should build the board in logical stages, starting with the any power-supply section, and then **test each section as you go**. It's **much** easier to debug a circuit if you've only stuffed a small part of it, compared to dealing with a whole populated board. Go down the signal chain, stuffing and testing. In the example circuit here, you would start with the battery mount (violating the first rule here), then the IC and resistors and timing capacitor, then see that the circuit is oscillating as expected, and only then add the diode, LED, and boosting capacitor.

### A.2.12 Protection

To save time, we will skip this step in lab, but this is something you can do later if you like. To protect the copper traces from oxidizing, they need a protective coating. In professionally fabricated boards, the copper is plated in nickel before soldering, and this works well. There are chemical kits that allow you to chemically plate your own PCBs, but their performance is variable, and it's usually not worth the expense and hassle.

The alternative is to coat the board after soldering with “conformal coating,” available in several materials. Acrylic conformal coating works well. It's something like a spray paint, but goes on thick enough to flow around components and solder joints, and dries to a clear layer. For repairs, a soldering iron will burn right through it. Note that this layer also insulates, so it must be scraped off for measuring voltages with probes. Also before spraying, make sure to mask off any components that shouldn't get gummed up with the coating material (jacks and connectors, pots, trim pots, edge connectors, etc.).

Before you spray on the coating, you must remove the flux from the soldering step, otherwise the conformal coating will look bad or may not adhere well. You can do this by giving the board a bath in isopropanol (ordinary rubbing alcohol; don't use stronger solvents like acetone, which will eat components). You will also need to use a stiff brush (toothbrush or camel-hair brush) to clean off the most stubborn bits of flux, or use an ultrasonic bath if you have one.

Another trick is to brush on a layer of liquid rosin flux<sup>6</sup> before soldering, and let it dry. It will be kind of sticky for a day or more until it dries thoroughly but after that it makes a nice, hard, solderable layer that resists oxidation. It can be cleaned off using isopropanol to put on a conformal coating later.

## A.3 Characterizing the Circuit

Now that you have the circuit stuffed, put in the battery. You should see the LED flashing. If not, recheck the circuit.

Make sure you understand the circuit:

1. Does the period of the blinking match your expectation from the resistor and capacitor values?
2. Make a timing diagram for the circuit (using an oscilloscope to measure the voltages) at pins 2 and 3 of the TLC551 (comparator input and oscillator output), as well as the LED cathode (reference all your voltages to “ground,” the negative terminal of the battery).
3. Based on your timing diagram, explain: how does the circuit light a 1.8 V LED using a 1.5 V battery?

---

<sup>6</sup>e.g., M.G. Chemicals 835: <http://www.mgchemicals.com/products/solder/fluxes/rosin-flux-835/>



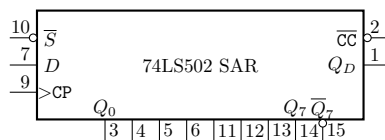
## Appendix B

# State-Machine Emulation of the 74LS502/3 Successive-Approximation Register

A nice example of an application of a general state machine is the emulation of classic, but obsolete, integrated circuits that acts as successive-approximation registers (SARs) for ADC circuits. We will discuss the application of two SARs, the 74LS502 and the 74LS503.<sup>1</sup>

### B.1 74502

The 74LS502 is a 8-bit successive-approximation register (SAR) that is shown schematically below.



It operates as follows. There are three inputs:

- $D$  is the input data (i.e., the output of the comparator feeds this input).
- $CP$  is the clock-pulse input; the SAR action occurs on the rising pulse edges.
- $\bar{S}$  is the “start” input. This is normally HIGH, and brought LOW for one cycle (i.e., during one rising edge of  $CP$ ) to begin conversion.

There are also 11 outputs:

- $Q_0$ – $Q_7$  are the (parallel) digital conversion outputs.
- $\bar{Q}_7$  is an inverted copy of  $Q_7$  (useful for signed-value ADC).
- $\bar{CC}$  is HIGH during conversion, and goes LOW when the ADC operation is complete.
- $Q_D$  is a synchronized copy of  $D$  (i.e.,  $Q_D$  latches the state of  $D$  at the last rising clock edge), which could be used for serial data output.

Pin 8 ground and pin 16 power are not shown in the schematic diagram.

The operational rules are as follows.

---

<sup>1</sup>Related files and more information available at <http://atomoptics.uoregon.edu/~dsteck/teaching/74503>.

- On a LOW start pulse on  $\overline{S}$ , the chip sets  $Q_7Q_6 \cdots Q_0 = 01111111$ , and  $\overline{CC} = 1$ .
- Next clock pulse: set  $Q_7 = D$ ,  $Q_6 = 0$ ,  $Q_D = D$ .
- Next clock pulse: set  $Q_6 = D$ ,  $Q_5 = 0$ ,  $Q_D = D$
- $\vdots$
- Next clock pulse: set  $Q_1 = D$ ,  $Q_0 = 0$ ,  $Q_D = D$ .
- Next clock pulse: set  $Q_0 = D$ ,  $\overline{CC} = 0$ .
- On subsequent clock pulses, we only care that the parallel data  $Q_0$ – $Q_7$  and  $\overline{CC}$  do not change.

We can also summarize this via the truth table (X = “don’t care”):

clock cycle	inputs		outputs									
	$D$	$\overline{S}$	$Q_D$	$Q_7$	$Q_6$	$Q_5$	$Q_4$	$Q_3$	$Q_2$	$Q_1$	$Q_0$	$\overline{CC}$
0	X	0	X	X	X	X	X	X	X	X	X	X
1	$D_7$	1	X	0	1	1	1	1	1	1	1	1
2	$D_6$	1	$D_7$	$D_7$	0	1	1	1	1	1	1	1
3	$D_5$	1	$D_6$	$D_7$	$D_6$	0	1	1	1	1	1	1
4	$D_4$	1	$D_5$	$D_7$	$D_6$	$D_5$	0	1	1	1	1	1
5	$D_3$	1	$D_4$	$D_7$	$D_6$	$D_5$	$D_4$	0	1	1	1	1
6	$D_2$	1	$D_3$	$D_7$	$D_6$	$D_5$	$D_4$	$D_3$	0	1	1	1
7	$D_1$	1	$D_2$	$D_7$	$D_6$	$D_5$	$D_4$	$D_3$	$D_2$	0	1	1
8	$D_0$	1	$D_1$	$D_7$	$D_6$	$D_5$	$D_4$	$D_3$	$D_2$	$D_1$	0	1
9	X	1	$D_0$	$D_7$	$D_6$	$D_5$	$D_4$	$D_3$	$D_2$	$D_1$	$D_0$	0
10	X	1	X	$D_7$	$D_6$	$D_5$	$D_4$	$D_3$	$D_2$	$D_1$	$D_0$	0

### B.1.1 General Emulation Notes

Now we will review the logic to implement this chip as a synchronous state machine. We will need (flip-flop) register outputs for  $Q_0$ – $Q_7$ ,  $Q_D$ , and  $\overline{CC}$ . We will also use three extra register bits  $C_2$ ,  $C_1$ , and  $C_0$ , where  $C_2C_1C_0$  gives (in binary) the next parallel-output bit to latch. That is, the circuit should set  $C_2C_1C_0 = 111$  on clock cycle 1 in the truth table, and then count backwards to 000.

To get started, first note that  $\overline{Q}_7$  can just be implemented with an extra NOT gate at the output of  $Q_7$ , and does not require its own register output. Next, note that we can implement

$$D_{Q_5} = \overline{\overline{S}} + (C_2\overline{C_1}C_0D + \overline{C_2\overline{C_1}C_0Q_5})\overline{C_2C_1\overline{C_0}}. \quad (\text{B.1})$$

Read this as follows. If the start pulse is 0 (so  $S = \overline{\overline{S}} = 1$ ), then force  $Q_5 = 1$ , which is accomplished by the first term. The second term has an overall multiplier that forces the expression to 0 (unless there is an override by  $\overline{S}$ ) when the  $C_2C_1C_0$  is at 6. The  $D$  term then stores the input data  $D$  when  $C_2C_1C_0$  is at 5, and the  $Q_5$  term allows  $Q_5$  to persist for other values of  $C_2C_1C_0$ .

As an example of a counter bit, note that  $C_1$  should change when  $C_2C_1C_0$  is X10 or 100 (remember we are counting backwards in binary). In either case it simply toggles, so

$$D_{C_1} = \overline{\overline{S}} + \overline{C_1\overline{C_0}}C_1 + C_2\overline{C_1}\overline{C_0}. \quad (\text{B.2})$$

The first term forces this bit to 1 on the start pulse, the last term forces the bit to 1 on  $C_2C_1C_0 = 100$ , and the middle term is zero on  $C_2C_1C_0 = \text{X}10$ , and allows  $C_1$  to persist otherwise.

The complete set of logic expressions is as follows.

$$\begin{aligned}
D_{Q_7} &= \overline{S}(C_2C_1C_0D + \overline{C_2C_1C_0}Q_7) \\
D_{Q_6} &= \overline{S} + (C_2C_1\overline{C_0}D + \overline{C_2C_1\overline{C_0}}Q_6)\overline{C_2C_1C_0} \\
D_{Q_5} &= \overline{S} + (C_2\overline{C_1}C_0D + \overline{C_2\overline{C_1}C_0}Q_5)\overline{C_2C_1C_0} \\
D_{Q_4} &= \overline{S} + (C_2\overline{C_1}\overline{C_0}D + \overline{C_2\overline{C_1}\overline{C_0}}Q_4)\overline{C_2C_1C_0} \\
D_{Q_3} &= \overline{S} + (\overline{C_2}C_1C_0D + \overline{\overline{C_2}C_1C_0}Q_3)\overline{C_2C_1C_0} \\
D_{Q_2} &= \overline{S} + (\overline{C_2}C_1\overline{C_0}D + \overline{\overline{C_2}C_1\overline{C_0}}Q_2)\overline{C_2C_1C_0} \\
D_{Q_1} &= \overline{S} + (\overline{C_2}\overline{C_1}C_0D + \overline{\overline{C_2}\overline{C_1}C_0}Q_1)\overline{C_2C_1C_0} \\
D_{Q_0} &= \overline{S} + (\overline{C_2}\overline{C_1}\overline{C_0}\overline{CC}D + \overline{\overline{C_2}\overline{C_1}\overline{C_0}\overline{CC}}Q_0)\overline{C_2C_1C_0} \\
D_{\overline{CC}} &= \overline{S} + \overline{C_2C_1C_0}\overline{CC} \\
D_{C_2} &= \overline{S} + \overline{C_2C_1C_0}C_2 \\
D_{C_1} &= \overline{S} + \overline{C_1C_0}C_1 + C_2\overline{C_1}\overline{C_0} \\
D_{C_0} &= \overline{S} + \overline{C_2C_1C_0}\overline{C_0} \\
D_{Q_D} &= D \\
\overline{Q_7} &= \overline{(Q_7)}
\end{aligned} \tag{B.3}$$

To understand all these in some detail:

- The easiest ones to understand are the last two:  $\overline{Q_7}$  is just an inverted copy of  $Q_7$ , and  $Q_D$  just latches  $D$ .
- On all the others, the start pulse forces the bits to 1, except for  $Q_7$ , which is forced to 0. Note that the start pulse overrides all other information.
- Then  $Q_7$  loads  $D$  when  $C_2C_1C_0 = 111$  and persists otherwise. We covered the similar logic for  $Q_5$  already;  $Q_0$ – $Q_4$  and  $Q_6$  follow the same idea in persisting, except changing to 0 and storing  $D$  at the right stages. However,  $Q_0$  is slightly more complicated since we stop the counter on  $C_2C_1C_0 = 000$ , so we must also use the  $\overline{CC}$  bit to make sure the state of  $Q_0$  persists when conversion is complete.
- For  $\overline{CC}$ , this should be forced to zero on  $C_2C_1C_0 = 000$  (i.e., after the last comparison), otherwise it persists. The factor of  $\overline{CC}$  is not strictly necessary on the second term.
- We covered the logic of the  $C_1$  counter bit already. Then  $C_0$  is simple in toggling on each clock pulse, except that once  $C_2C_1C_0 = 000$ , we will force it to stay at 0. For  $C_2$ , the MSB should only change on  $C_2C_1C_0 = 100$ , so we detect this case explicitly and use a NOT operation to reset to 0, otherwise  $C_2$  persists.

### B.1.2 22V10 Emulation

To emulate the 74LS502, we will choose the 22V10 SPLD (e.g., the ATF22V10C from Atmel). This has a 10-bit register (i.e., 10 D-type flip-flops) and 10 outputs, plus plenty of logic for sum-of-product logic and plenty of inputs. However, there is a problem: we have 13 registered outputs, but only 10 bits in the register. To handle this, note that the counter-logic bits  $C_0$ – $C_2$  are actually redundant—and that is, although they are conceptually useful states in writing down the state-machine logic, they are not needed as register variables in the sense that they can be inferred from the other register variables  $Q_0$ – $Q_7$  and  $\overline{CC}$ . In particular, note that



- $C_2C_1C_0 = 111$  is determined by  $Q_7 \cdots Q_0 \overline{CC} = 01111111$
- $C_2C_1C_0 = 110$  is determined by  $Q_7 \cdots Q_0 \overline{CC} = X0111111$
- $C_2C_1C_0 = 101$  is determined by  $Q_7 \cdots Q_0 \overline{CC} = XX011111$
- etc.

Thus,  $C_0$  should trigger on an even number of ones after the leading zero:

$$\begin{aligned}
 C_0 &= \overline{Q_7}Q_6Q_5Q_4Q_3Q_2Q_1Q_0\overline{CC} + \overline{Q_5}Q_4Q_3Q_2Q_1Q_0\overline{CC} + \overline{Q_3}Q_2Q_1Q_0\overline{CC} + \overline{Q_1}Q_0\overline{CC} \\
 &= (((\overline{Q_7}Q_6Q_5 + \overline{Q_5})Q_4Q_3 + \overline{Q_3})Q_2Q_1 + \overline{Q_1})Q_0\overline{CC} \\
 &= (((\overline{Q_7}Q_6 + \overline{Q_5})Q_4 + \overline{Q_3})Q_2 + \overline{Q_1})Q_0\overline{CC}.
 \end{aligned} \tag{B.4}$$

We used  $A + \overline{A}B = A + B$  in the last step. Similarly  $C_1$  triggers when there are 8, 7, 4, or 3 ones after the leading zero:

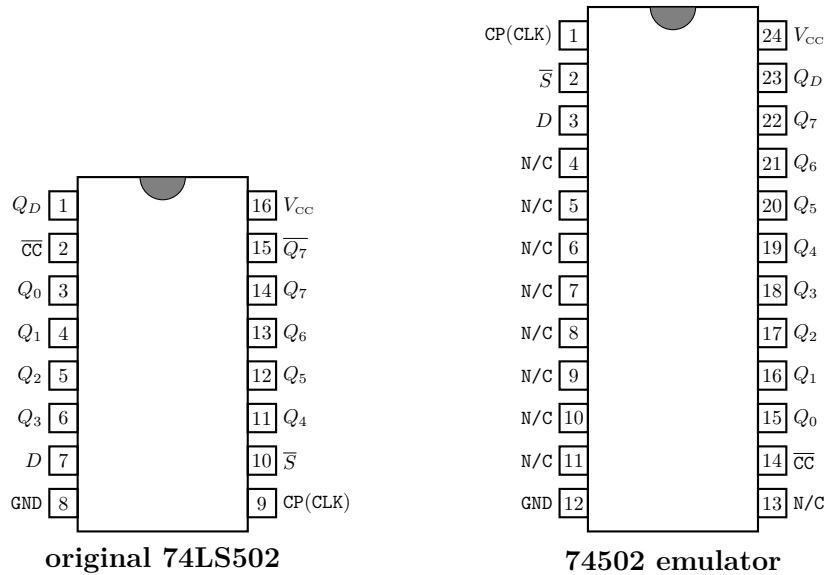
$$\begin{aligned}
 C_1 &= \overline{Q_7}Q_6Q_5Q_4Q_3Q_2Q_1Q_0\overline{CC} + \overline{Q_6}Q_5Q_4Q_3Q_2Q_1Q_0\overline{CC} + \overline{Q_3}Q_2Q_1Q_0\overline{CC} + \overline{Q_2}Q_1Q_0\overline{CC} \\
 &= (((\overline{Q_7}Q_6 + \overline{Q_6})Q_5Q_4Q_3 + \overline{Q_3})Q_2 + \overline{Q_2})Q_1Q_0\overline{CC} \\
 &= ((\overline{Q_7} + \overline{Q_6})Q_5Q_4 + \overline{Q_3} + \overline{Q_2})Q_1Q_0\overline{CC}.
 \end{aligned} \tag{B.5}$$

Finally,  $C_2$  triggers when there are 8, 7, 6, or 5 ones after the leading zero:

$$\begin{aligned}
 C_2 &= \overline{Q_7}Q_6Q_5Q_4Q_3Q_2Q_1Q_0\overline{CC} + \overline{Q_6}Q_5Q_4Q_3Q_2Q_1Q_0\overline{CC} + \overline{Q_5}Q_4Q_3Q_2Q_1Q_0\overline{CC} + \overline{Q_4}Q_3Q_2Q_1Q_0\overline{CC} \\
 &= (((\overline{Q_7}Q_6 + \overline{Q_6})Q_5 + \overline{Q_5})Q_4 + \overline{Q_4})Q_3Q_2Q_1Q_0\overline{CC} \\
 &= (\overline{Q_7} + \overline{Q_6} + \overline{Q_5} + \overline{Q_4})Q_3Q_2Q_1Q_0\overline{CC}.
 \end{aligned} \tag{B.6}$$

Thus, we are down to 10 registered outputs, and we may proceed.

The pinout for the original 74LS502 and the 22V10-based emulator are shown below. Note that because we only have 10 outputs available, we have chosen to include the  $Q_D$  output but not the  $\overline{Q_7}$  output, but we could easily make the opposite choice, as we will discuss below.



The code to implement the state machine in the **CUPL** (Compiler for Universal Programmable Logic) programming language<sup>2</sup> is shown below.

74502-22V10.pld

```

/*
 * 74502 SAR emulator, on a 22v10
 */
Name          74502-22V10;
Partno        74502;
Revision      01;
Date          5/23/2015;
Designer      Daniel Steck;
Company       University of Oregon;
Location      None;
Assembly      None;
Device        g22v10;

/** inputs **/
pin 1 = CP;          /* clock pulse (trig on rising edge) */
pin 2 = !S;          /* start low */
pin 3 = Din;         /* data */

/** outputs **/
pin 14        = !CC;      /* conversion complete low */
pin [15..22] = [Q0..Q7]; /* 8-bit output */
pin 23        = QD;       /* registered/synchronous copy of data input */

/** intermediate counter variables **/
C0 = (((!Q7 & Q6 # !Q5) & Q4 # !Q3) & Q2 # !Q1) & Q0 & !CC;
C1 = ((!Q7 # !Q6) & Q5 & Q4 # !Q3 # !Q2) & Q1 & Q0 & !CC;
C2 = (!Q7 # !Q6 # !Q5 # !Q4) & Q3 & Q2 & Q1 & Q0 & !CC;

/** register inputs **/
CC.D = !(S # !(C2 & !C1 & !C0) & !CC);
Q7.D = !S & (C2 & C1 & C0 & Din # !(C2 & C1 & C0) & Q7);
Q6.D = S # (C2 & C1 & !C0 & Din # !(C2 & C1 & !C0) & Q6) & !(C2 & C1 & C0);
Q5.D = S # (C2 & !C1 & C0 & Din # !(C2 & !C1 & C0) & Q5) & !(C2 & C1 & !C0);
Q4.D = S # (C2 & !C1 & !C0 & Din # !(C2 & !C1 & !C0) & Q4) & !(C2 & !C1 & C0);
Q3.D = S # (!C2 & C1 & C0 & Din # !(C2 & C1 & C0) & Q3) & !(C2 & !C1 & !C0);
Q2.D = S # (!C2 & C1 & !C0 & Din # !(C2 & C1 & !C0) & Q2) & !(C2 & C1 & C0);
Q1.D = S # (!C2 & !C1 & C0 & Din # !(C2 & !C1 & C0) & Q1) & !(C2 & C1 & !C0);
Q0.D = S # (!C2 & !C1 & !C0 & !CC & Din # !(C2 & !C1 & !C0 & !CC) & Q0) & !(C2 & !C1 & C0);
QD.D = Din;

/** handle flip-flop variables set/preset inputs **/
CC.ar = 'b'0;
Q7.ar = 'b'0;
Q6.ar = 'b'0;
Q5.ar = 'b'0;
Q4.ar = 'b'0;
Q3.ar = 'b'0;
Q2.ar = 'b'0;
Q1.ar = 'b'0;
Q0.ar = 'b'0;
QD.ar = 'b'0;

```

<sup>2</sup>The only compiler realistically available nowadays is WinCUPL, which is Windows-based, crash-prone, and proprietary, but it is freely distributed by Atmel: <http://www.atmel.com/tools/wincupl.aspx>.

```

CC.sp = 'b'0;
Q7.sp = 'b'0;
Q6.sp = 'b'0;
Q5.sp = 'b'0;
Q4.sp = 'b'0;
Q3.sp = 'b'0;
Q2.sp = 'b'0;
Q1.sp = 'b'0;
Q0.sp = 'b'0;
QD.sp = 'b'0;

```

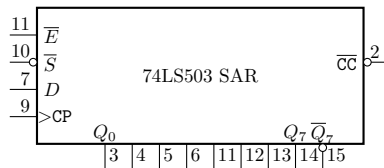
This is a relatively straightforward translation of the equations we have already written down. Note that

- The first block contains obligatory header information, most of which is merely informational, but the “Device” declaration to “g22v10” means we have selected this device (which also covers the ATF22V10C variant).
- The next two blocks declare input and output pin assignments and variables. Note that a NOT is denoted by “!”; for example,  $\overline{S}$  is denoted !S.
- The next block gives the expressions (B.4)–(B.6) for the counter variables  $C_0$ – $C_2$ . Note that the OR operation is represented by a hash (#), and the AND operation is represented by an ampersand (&).
- The next block gives expressions for all the register inputs, as in Eqs. (B.3) (except for counter register variables). The CUPL notation is that the input for the register variable Q7 is Q7.D (i.e., this is what we call  $D_{Q_7}$ ).
- Finally, in the last block, we make sure to tie the other flip-flop controls to default values. Here, the 22V10 flip-flops have asynchronous-reset (AR) and synchronous-preset (SP) inputs; we simply tie all of them to logical 0 (written as “binary 0” or 'b'0 in CUPL).
- To have  $\overline{Q_7}$  as an output instead of  $D_7$ , we can change the “pin 23 = QD;” declaration to now read “pin 23 = notQ7;”, and then change the line “QD.D = Din;” to “notQ7 = !Q7;” (note that this output would then no longer be registered, but “combinatorial”).

The compiler then simplifies the logical expressions and figures out how to make the fused connections in the configurable-logic section of the 22V10. A separate programmer is necessary to then “burn” the chip.

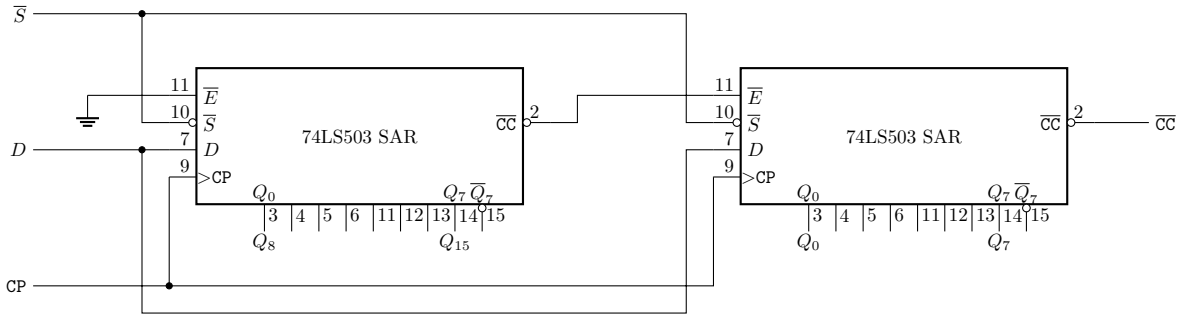
## B.2 74503

The 74LS503 is basically the same as the '502, but it dispenses with  $Q_D$  and adds an enable-LOW input  $\overline{E}$ .



The  $\overline{E}$  operates as follows. If it is held LOW, then the chip behavior is essentially identical to the '502. If it is taken HIGH—the intent is for this to happen after the start operation but before any data acquisition occurs—then  $Q_7$  is *asynchronously* forced HIGH, and the chip does not accept any data from  $D$ . When  $\overline{E} = 0$  again, the acquisition process proceeds as in the '502.

The idea behind the  $\overline{E}$  input in the '503 is that two '503's can be “stacked” to realize a 16-bit SAR. The idea is that the  $\overline{CC}$  of the most-significant chip (byte) drives the  $\overline{E}$  of the least-significant chip, so that when the first chip is finished, acquisition continues on the second chip. The connections are shown in the data sheet for the 74LS503. The idea is to share the data, clock, and start lines, and chain the  $\overline{CC}$  of the MS chip to the enable of the enable of the LS chip.



### B.2.1 ATF750C Emulation

To emulate the 74LS503, we will choose the ATF750C CPLD (from Atmel). Note that when the  $C_j$  are eliminated as in Eqs. (B.4)–(B.6) in the case of the '503, the expressions for  $D_{Q_0}$ – $D_{Q_6}$  must have “+E” tacked on to each expression, since these variables track the counting state. However, this extra addition appears to make the logic too complicated to fit in either the 22V10 or the ATF750C.

Fortunately, the ATF750C has 10 extra register bits that are present “internally” (i.e., they can not be connected directly to outputs as are the other 10 register bits). Thus we can implement  $C_0$ – $C_2$  as register variables.

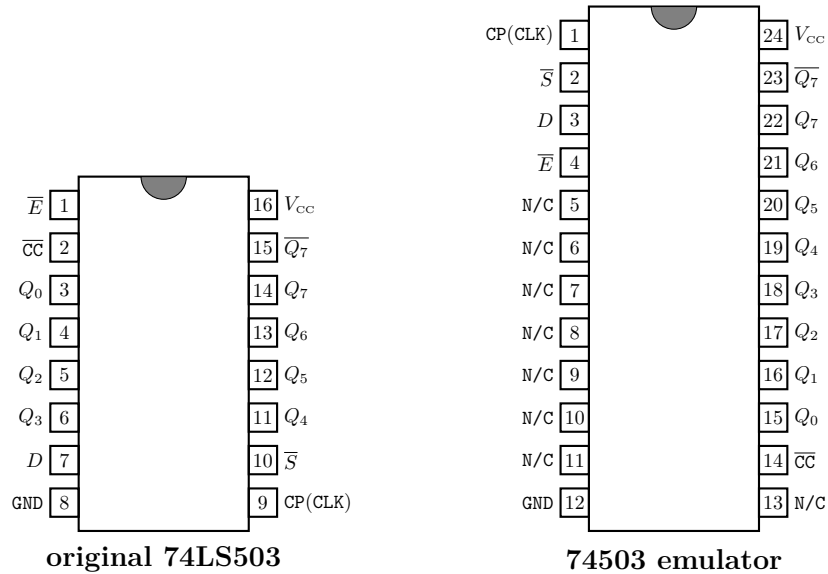
Now we need to show how to modify the '502 logic to accommodate this new input. Note that since the effect on  $Q_7$  is asynchronous,  $Q_7$  can't any more be a register output. So for the sake of notation, let  $P_7$  be a register output, and let  $\bar{Q}_7$  be a Boolean function of  $P_7$  and  $\bar{E}$ . The other outputs  $Q_0$ – $Q_6$  and  $\bar{C}\bar{C}$  can still be register outputs. The summary of logic expressions is below.

$$\begin{aligned}
 D_{P_7} &= \bar{S}(C_2 C_1 C_0 D + \bar{C}_2 \bar{C}_1 \bar{C}_0 P_7) \bar{E} \\
 Q_7 &= \bar{E} + P_7 \\
 D_{Q_6} &= \bar{S} + (C_2 C_1 \bar{C}_0 D + \bar{C}_2 \bar{C}_1 \bar{C}_0 Q_6) \bar{C}_2 \bar{C}_1 \bar{C}_0 \\
 D_{Q_5} &= \bar{S} + (C_2 \bar{C}_1 C_0 D + \bar{C}_2 \bar{C}_1 C_0 Q_5) \bar{C}_2 \bar{C}_1 \bar{C}_0 \\
 D_{Q_4} &= \bar{S} + (C_2 \bar{C}_1 \bar{C}_0 D + \bar{C}_2 \bar{C}_1 \bar{C}_0 Q_4) \bar{C}_2 \bar{C}_1 \bar{C}_0 \\
 D_{Q_3} &= \bar{S} + (\bar{C}_2 C_1 C_0 D + \bar{C}_2 C_1 C_0 Q_3) \bar{C}_2 \bar{C}_1 \bar{C}_0 \\
 D_{Q_2} &= \bar{S} + (\bar{C}_2 C_1 \bar{C}_0 D + \bar{C}_2 C_1 \bar{C}_0 Q_2) \bar{C}_2 \bar{C}_1 \bar{C}_0 \\
 D_{Q_1} &= \bar{S} + (\bar{C}_2 \bar{C}_1 C_0 D + \bar{C}_2 \bar{C}_1 C_0 Q_1) \bar{C}_2 \bar{C}_1 \bar{C}_0 \\
 D_{Q_0} &= \bar{S} + (\bar{C}_2 \bar{C}_1 \bar{C}_0 \bar{C}\bar{C} D + \bar{C}_2 \bar{C}_1 \bar{C}_0 \bar{C}\bar{C} Q_0) \bar{C}_2 \bar{C}_1 \bar{C}_0 \\
 D_{\bar{C}\bar{C}} &= \bar{S} + \bar{C}_2 \bar{C}_1 \bar{C}_0 \bar{C}\bar{C} \\
 D_{C_2} &= \bar{S} + \bar{C}_2 \bar{C}_1 \bar{C}_0 C_2 + \bar{E} \\
 D_{C_1} &= \bar{S} + \bar{C}_1 \bar{C}_0 C_1 + C_2 \bar{C}_1 \bar{C}_0 + \bar{E} \\
 D_{C_0} &= \bar{S} + \bar{C}_2 \bar{C}_1 \bar{C}_0 \bar{C}_0 + \bar{E} \\
 \bar{Q}_7 &= (\bar{Q}_7)
 \end{aligned} \tag{B.7}$$

Note that we are forcing  $P_7 = 0$  on a disable cycle, so that conversion happens correctly afterwards (otherwise spurious data could be loaded). Then the expression for  $Q_7$  allows  $P_7$  to be overridden by  $\bar{E}$  asynchronously. Also, although overkill, the counter bits are all forced to 1 to ensure conversion occurs correctly. (Really, this should only be needed for  $C_0$ ). We only need correct behavior if the chip is disabled right after a start

pulse, so the other bits should be okay. However, it is okay to add “+ $E$ ” to the other register inputs if desired. We have also dispensed with  $D_{Q_D}$ .

The pin diagrams for the original chip and emulator are shown below. Note that we can now accommodate every output from the original chip.



The code to implement the state machine in CUPL is shown below.

74503-F750C.pld

```

/*
 * 74503 SAR emulator, on an ATF750C
 */
Name          74503-F750C;
Partno        74503;
Revision      01;
Date          5/23/2015;
Designer      Daniel Steck;
Company       University of Oregon;
Location      None;
Assembly      None;
Device        v750c;

/** inputs */
pin 1 = CP;          /* clock pulse (trig on rising edge) */
pin 2 = !S;          /* start low */
pin 3 = Din;         /* data */
pin 4 = !E;          /* enable low */

/** outputs */
pin 14 = !CC;        /* conversion complete low */
pin [15..22] = [Q0..Q7]; /* 8-bit output */
pin 23 = !Q7copy;    /* inverted copy of Q7 */

/** internal nodes */
node P7;
node [C0..C2];

/** intermediate counter variables */

```

```

C0.D = S # !(C2 & !C1 & !C0) & !C0 # !E;
C1.D = S # !(C1 & !C0) & C1 # C2 & !C1 & !C0 # !E;
C2.D = S # !(C2 & !C1 & !C0) & C2 # !E;

/** register inputs */
CC.D = !(S # !(C2 & !C1 & !C0) & !CC);
P7.D = !S & (C2 & C1 & C0 & Din # !(C2 & C1 & C0) & P7) & E;
Q6.D = S # ( C2 & C1 & !C0 & Din # !( C2 & C1 & !C0) & Q6) & !( C2 & C1 & C0);
Q5.D = S # ( C2 & !C1 & C0 & Din # !( C2 & !C1 & C0) & Q5) & !( C2 & C1 & !C0);
Q4.D = S # ( C2 & !C1 & !C0 & Din # !( C2 & !C1 & !C0) & Q4) & !( C2 & !C1 & C0);
Q3.D = S # (!C2 & C1 & C0 & Din # !(C2 & C1 & C0) & Q3) & !( C2 & !C1 & !C0);
Q2.D = S # (!C2 & C1 & !C0 & Din # !(C2 & C1 & !C0) & Q2) & !(C2 & C1 & C0);
Q1.D = S # (!C2 & !C1 & C0 & Din # !(C2 & !C1 & C0) & Q1) & !(C2 & C1 & !C0);
Q0.D = S # (!C2 & !C1 & !C0 & !CC & Din # !(C2 & !C1 & !C0 & !CC) & Q0) & !(C2 & !C1 & C0);

/** combinatorial outputs */
Q7 = !E # P7;
Q7copy = Q7;

/** handle flip-flop variables set/preset inputs */
CC.ar = 'b'0;
C2.ar = 'b'0;
C1.ar = 'b'0;
C0.ar = 'b'0;
P7.ar = 'b'0;
Q6.ar = 'b'0;
Q5.ar = 'b'0;
Q4.ar = 'b'0;
Q3.ar = 'b'0;
Q2.ar = 'b'0;
Q1.ar = 'b'0;
Q0.ar = 'b'0;

CC.sp = 'b'0;
C2.sp = 'b'0;
C1.sp = 'b'0;
C0.sp = 'b'0;
P7.sp = 'b'0;
Q6.sp = 'b'0;
Q5.sp = 'b'0;
Q4.sp = 'b'0;
Q3.sp = 'b'0;
Q2.sp = 'b'0;
Q1.sp = 'b'0;
Q0.sp = 'b'0;

/** flip-flop-clock multiplexer (use input clock pin) */
CC.ckmux = CP;
C2.ckmux = CP;
C1.ckmux = CP;
C0.ckmux = CP;
P7.ckmux = CP;
Q6.ckmux = CP;
Q5.ckmux = CP;
Q4.ckmux = CP;
Q3.ckmux = CP;
Q2.ckmux = CP;

```

```
Q1.ckmux = CP;
Q0.ckmux = CP;
```

Again, this is a relatively straightforward translation of the Boolean-algebraic equations. Note that

- In the first block, we now declare the more powerful chip (“v750c”).
- In the next two blocks, we declare  $\overline{E}$  and  $\overline{Q_7}$  (the latter by defining the “copy” variable `Q7copy`).
- In the next block, we implement the counter variables  $C_0$ – $C_2$  as register variables, as in Eqs. (B.7).
- The next block gives expressions for all the register inputs, as in Eqs. (B.7) (except for the counter register variables we already implemented). The subsequent block implements the combinatorial outputs  $Q_7$  and  $\overline{Q_7}$ .
- Finally, in the last block, the ATF750C has a “clock multiplexer” control on the flip-flop inputs. The upshot is that we must declare the flip-flop clock inputs to be connected to the clock-input pin `CP`.
- To have  $\overline{Q_7}$  as an output instead of  $D_7$ , we can change the “`pin 23 = QD;`” declaration to now read “`pin 23 = notQ7;`”, and then change the line “`QD.D = Din;`” to “`notQ7 = !Q7;`” (note that this output would then no longer be registered, but “combinatorial”).

The compiler then simplifies the logical expressions and figures out how to make the fused connections in the configurable-logic section of the 22V10. A separate programmer is necessary to then “burn” the chip.

### B.3 Testing the State Machines

The WinCUPL package also allows simulation tests. We define the test values in another file (i.e., define a sequence of input and expected output values). The simulator will simulate the chip and ensure that it passes the test values. These test values can also be embedded in the code to be sent to the programmer, so the programmer can test the actual chip.

Test files for both emulators are attached below; it’s a good exercise to look through and understand these. The header block here matches that of the `.pld` file. The `ORDER` declaration gives a sequence of pins (variables) for consideration. In the `VECTOR` block, we give a bunch of input/expected-output states, in the order of the `ORDER` declaration. The notation for the values is:

- 0 and 1 are the logical input values.
- L and H are corresponding logical values, but used for expected outputs.
- c is equivalent to a 0, then a 1, and then a 0 (i.e., a clock pulse).

```
74502-22V10.si
```

```
/*
 * 74502 SAR emulator, on a 22v10
 */
Name          74502-22V10;
Partno        74502;
Revision      01;
Date          5/23/2015;
Designer      Daniel Steck;
Company       University of Oregon;
Location      None;
Assembly      None;
Device        g22v10;

ORDER: CP, !S, Din, QD, Q7, Q6, Q5, Q4, Q3, Q2, Q1, Q0, !CC;
```

## VECTORS:

```
c 00 LLHHHHHHHH
c 10 LLLHHHHHHH
c 10 LLLLHHHHHH
c 10 LLLLLHHHHH
c 10 LLLLLLHHHH
c 10 LLLLLLLHHH
c 10 LLLLLLLLHH
c 10 LLLLLLLLLH
c 10 LLLLLLLLLL
```

```
c 01 HLHHHHHHHH
c 11 HHLHHHHHHH
c 11 HHHLHHHHHH
c 11 HHHHLHHHHH
c 11 HHHHHLHHHH
c 11 HHHHHLHHHH
c 11 HHHHHHLHHH
c 11 HHHHHHLLHH
c 11 HHHHHHLLHH
c 11 HHHHHHLLHH
c 11 HHHHHHLLHH
```

```
c 01 HLHHHHHHHH
c 10 LLLHHHHHHH
c 11 HLHLHHHHHH
c 10 LLHLHHHHHH
c 11 HLHLHLHHHH
c 10 LLHLHLHHHH
c 11 HLHLHLHHHH
c 10 LLHLHLHLHH
c 11 HLHLHLHLHH
c 10 LLHLHLHLHH
c 11 HLHLHLHLHH
```

```
c 01 HLHHHHHHHH
c 11 HHLHHHHHHH
c 10 LHLHHHHHHH
c 11 HHLHLHHHHH
c 10 LHLHLHHHHH
c 11 HHLHLHLHHH
c 10 LHLHLHLHHH
c 11 HHLHLHLHHH
c 10 LHLHLHLHHH
c 11 HHLHLHLHHH
c 10 LHLHLHLHHH
c 11 HHLHLHLHHH
```

```
74503-F750C.si
```

```
/*
* 74503 SAR emulator, on an ATF750C
*/
Name          74503-F750C;
Partno        74503;
Revision      01;
Date          5/23/2015;
Designer      Daniel Steck;
Company       University of Oregon;
Location      None;
Assembly      None;
```



Device v750c;

ORDER:

CP, !E, !S, Din, C2, C1, C0, Q7, Q6, Q5, Q4, Q3, Q2, Q1, Q0, !CC;

VECTORS:

c 0 00 HHH LHHHHHHHH  
 c 0 10 HHL LLHHHHHHH  
 c 0 10 HLH LLLHHHHHH  
 c 0 10 HLL LLLLHHHHH  
 c 0 10 LHH LLLLHHHHH  
 c 0 10 LHL LLLLHHHHH  
 c 0 10 LLH LLLLHHHHH  
 c 0 10 LLL LLLLHHHHH  
 c 0 10 LLL LLLLHHHHH  
 c 0 10 LLL LLLLHHHHH

c 0 01 HHH LHHHHHHHH  
 c 0 11 HHL HLHHHHHHH  
 c 0 11 HLH HHLHHHHHH  
 c 0 11 HLL HHLHHHHHH  
 c 0 11 LHH HHLHHHHHH  
 c 0 11 LHL HHLHHHHHH  
 c 0 11 LLH HHLHHHHHH  
 c 0 11 LLL HHLHHHHHH  
 c 0 11 LLL HHLHHHHHH  
 c 0 11 LLL HHLHHHHHH

c 0 01 HHH LHHHHHHHH  
 c 0 10 HHL LLHHHHHHH  
 c 0 11 HLH LHLHHHHHH  
 c 0 10 HLL LHLHHHHHH  
 c 0 11 LHH LHLHHHHHH  
 c 0 10 LHL LHLHHHHHH  
 c 0 11 LLH LHLHHHHHH  
 c 0 10 LLL LHLHHHHHH  
 c 0 11 LLL LHLHHHHHH  
 c 0 10 LLL LHLHHHHHH

c 0 01 HHH LHHHHHHHH  
 c 0 11 HHL HLHHHHHHH  
 c 0 10 HLH HLLHHHHHH  
 c 0 11 HLL HLHLLHHHH  
 c 0 10 LHH HLHLLHHHH  
 c 0 11 LHL HLHLLHHHH  
 c 0 10 LLH HLHLLHHHH  
 c 0 11 LLL HLHLLHHHH  
 c 0 10 LLL HLHLLHHHH  
 c 0 11 LLL HLHLLHHHH

0 0 11 LLL HLHLLHHHH  
 1 0 00 HHH LHHHHHHHH  
 1 1 00 HHH HHHHHHHHH  
 1 1 01 HHH HHHHHHHHH  
 1 0 10 HHH LHHHHHHHH  
 0 0 10 HHH LHHHHHHHH  
 c 0 10 HHL LLHHHHHHH

```

c 0 10 HLH LLLHHHHHH
c 0 10 HLL LLLLHHHHH
c 0 10 LHH LLLLHHHHH
c 0 10 LHL LLLLHHHHH
c 0 10 LLH LLLLHHHHH
c 0 10 LLL LLLLHHHHH
c 0 10 LLL LLLLHHHHH
c 0 10 LLL LLLLHHHHH

```

```

0 0 00 LLL LLLLHHHHH
1 0 01 HHH LHHHHHHHH
1 1 01 HHH HHHHHHHHH
1 1 00 HHH HHHHHHHHH
1 0 11 HHH LHHHHHHHH
0 0 11 HHH LHHHHHHHH
c 0 11 HHL HLHHHHHHH
c 0 11 HLH HHLHHHHHH
c 0 11 HLL HHLHHHHHH
c 0 11 LHH HHHHLHHHH
c 0 11 LHL HHHHLHHHH
c 0 11 LLH HHHHLHHHH
c 0 11 LLL HHHHLHHHH
c 0 11 LLL HHHHHHHHL
c 0 11 LLL HHHHHHHHL

```

```

0 0 01 LLL HHHHHHHHL
1 0 01 HHH LHHHHHHHH
1 1 01 HHH HHHHHHHHH
1 1 00 HHH HHHHHHHHH
1 0 01 HHH LHHHHHHHH
0 0 11 HHH LHHHHHHHH
c 0 10 HHL LLHHHHHHH
c 0 11 HLH LHLHHHHHH
c 0 10 HLL LHLHHHHHH
c 0 11 LHH LHLHLHHHH
c 0 10 LHL LHLHLHHHH
c 0 11 LLH LHLHLHHHH
c 0 10 LLL LHLHLHLHH
c 0 11 LLL LHLHLHLHL
c 0 10 LLL LHLHLHLHL

```

```

0 0 00 LLL LHLHLHLHL
1 0 01 HHH LHHHHHHHH
1 1 01 HHH HHHHHHHHH
1 1 00 HHH HHHHHHHHH
1 0 01 HHH LHHHHHHHH
0 0 11 HHH LHHHHHHHH
c 0 11 HHL HLHHHHHHH
c 0 10 HLH HLLHHHHHH
c 0 11 HLL HLHLHHHHH
c 0 10 LHH HLHLHHHHH
c 0 11 LHL HLHLHHHHH
c 0 10 LLH HLHLHHHHH
c 0 11 LLL HLHLHLHLH
c 0 10 LLL HLHLHLHLH
c 0 11 LLL HLHLHLHLH

```



## Appendix C

# Gallery of Characteristic Curves

This appendix presents a haphazard collection of characteristic curves for discrete semiconductor components that are common, useful, or interesting (or some combination). Data sheets contain this information (and more), but parameters are specified in various ways, and plots are sometimes drawn in bizarre ways. The characteristic curves, however, contain a lot of information at a glance, and having a collection of similarly plotted curves makes it easy to visually compare similar components. (How many people can explain the difference between a 2N3904 and a 2N4401? They're almost interchangeable, but the curves show that they have potentially important differences, such as the larger collector resistance of the 2N4401.)

The characteristic curves here are computed from model parameters used in SPICE simulators. Devices with similar functionality are grouped together to make it easier to compare their properties. The calculations (in lieu of data), besides being easier to generate than real data, are useful in capturing the “typical” device behaviors, according to manufacturers.

The collection of components is in no way complete or representative of any kind of good engineering practice.

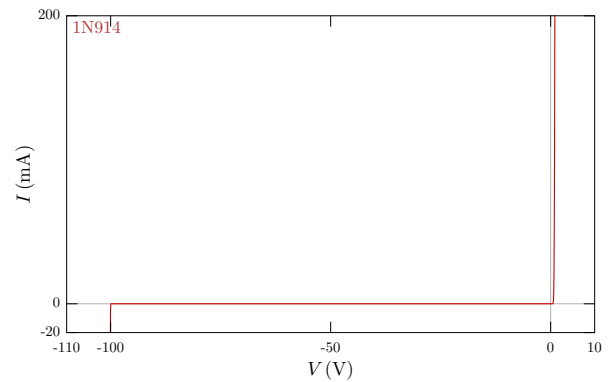
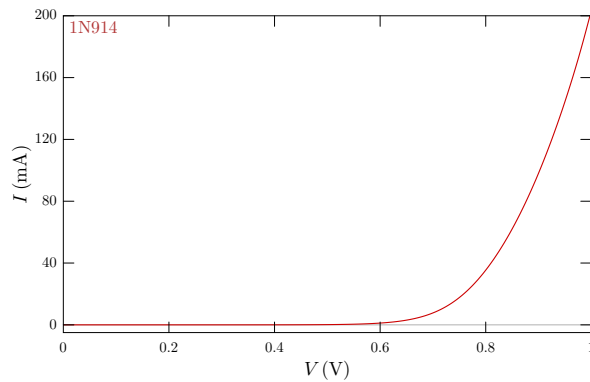
Now on to the pretty pictures...

## C.1 Diodes

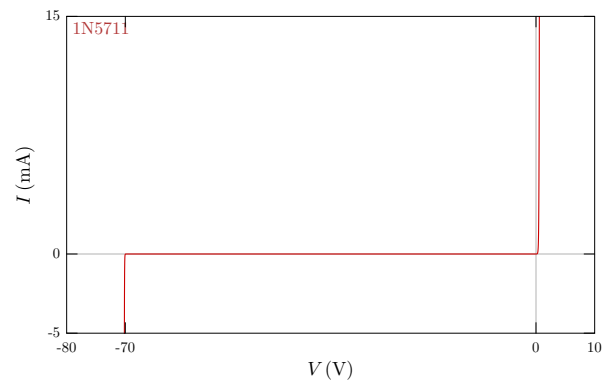
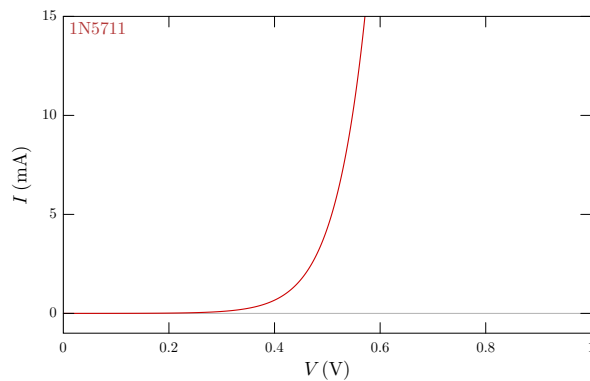
### C.1.1 Small-Signal Diodes

**1N914** (fast, high conductance, Si)<sup>1</sup>

The venerable 1N914 (good as a default small-signal diode) has many very similar cousins, including the 1N914A, 1N914B, 1N916, 1N916A, 1N916B, 1N4148, and the 1N4448. (To wit: all of these share the same data sheet.<sup>2</sup>)



**1N5711** (Schottky)<sup>3</sup>



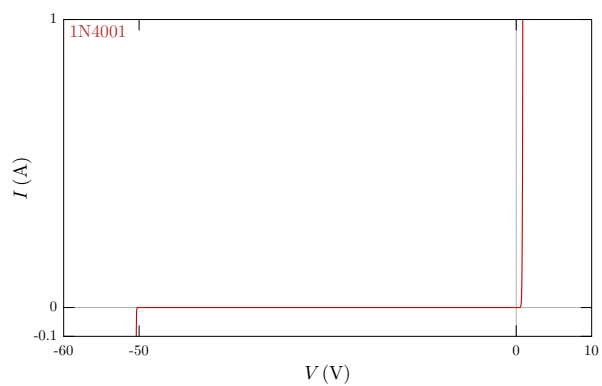
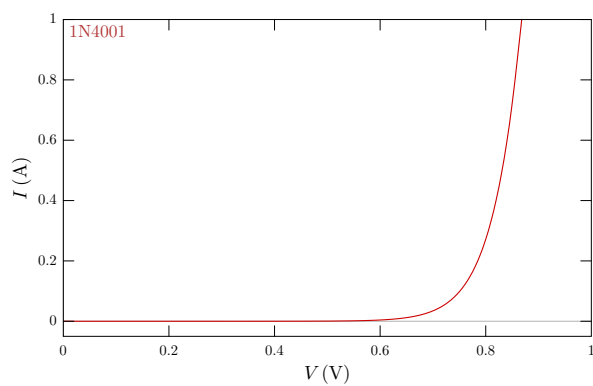
<sup>1</sup>Central Semiconductor model: <https://www.centralsemi.com/product/partpage2.php?part=1N914>

<sup>2</sup><http://www.onsemi.com/pub/Collateral/1N914-D.PDF>

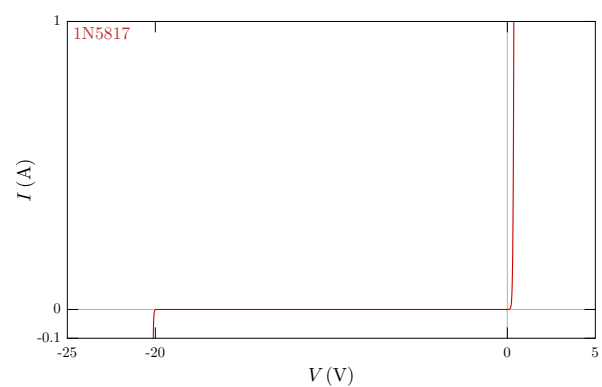
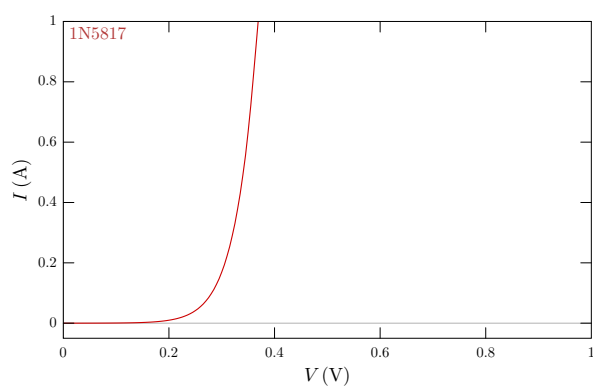
<sup>3</sup>Diodes Inc. model: [http://ltwiki.org/files/LTspiceIV/Vendor%20List/Diodes%20Incorporated/Spice/spicemodels\\_schottky\\_diodes.txt](http://ltwiki.org/files/LTspiceIV/Vendor%20List/Diodes%20Incorporated/Spice/spicemodels_schottky_diodes.txt)

### C.1.2 Medium-Current Rectifiers

1N4001 (50 V, 1 A Si rectifier)<sup>4</sup>



1N5817 (20 V, 1 A Schottky rectifier)<sup>5</sup>

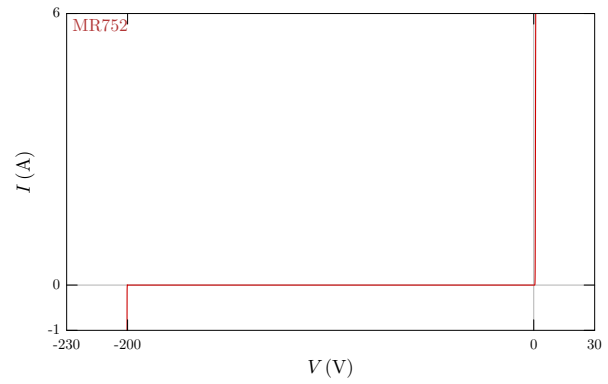
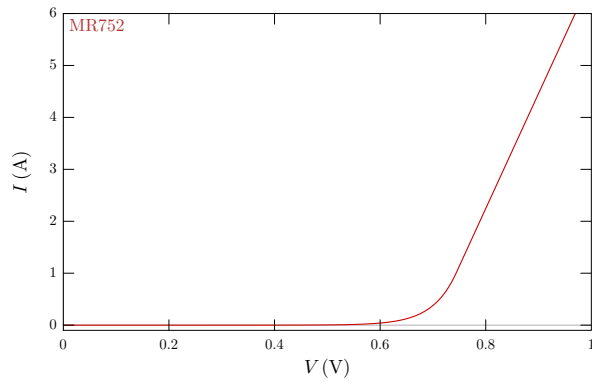


<sup>4</sup>ON Semiconductor model: <http://www.onsemi.com/PowerSolutions/product.do?id=1N4001>

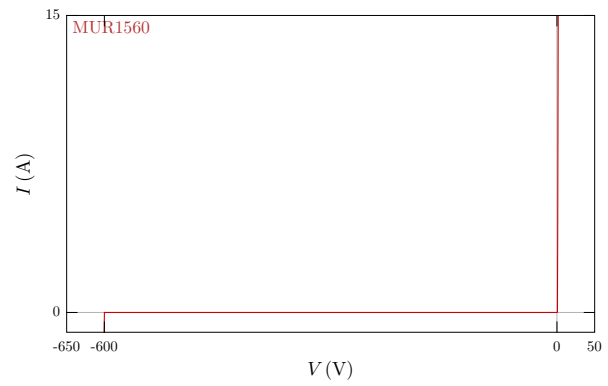
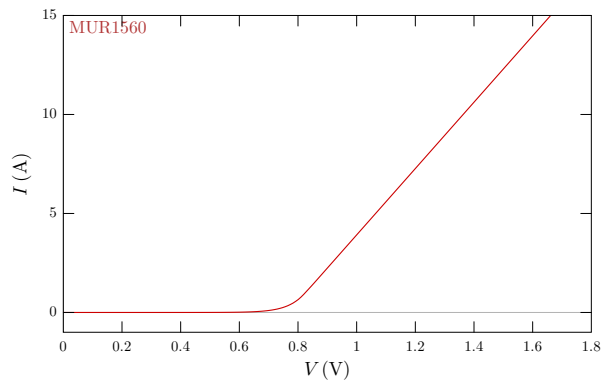
<sup>5</sup>ON Semiconductor model: <http://www.onsemi.com/PowerSolutions/product.do?id=1N5817>

### C.1.3 High-Current Rectifiers

MR752 (6 A, 200 V power rectifier)<sup>6</sup>



MUR1560 (15 A, 600 V power rectifier)<sup>7</sup>



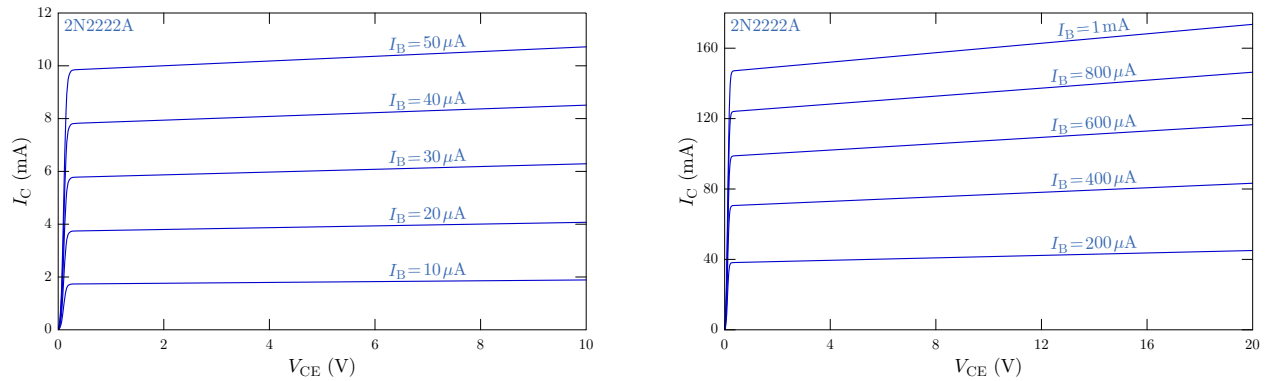
<sup>6</sup>ON Semiconductor model: <http://www.onsemi.com/PowerSolutions/product.do?id=MR752>

<sup>7</sup>ON Semiconductor model: <http://www.onsemi.com/PowerSolutions/product.do?id=MUR1560>

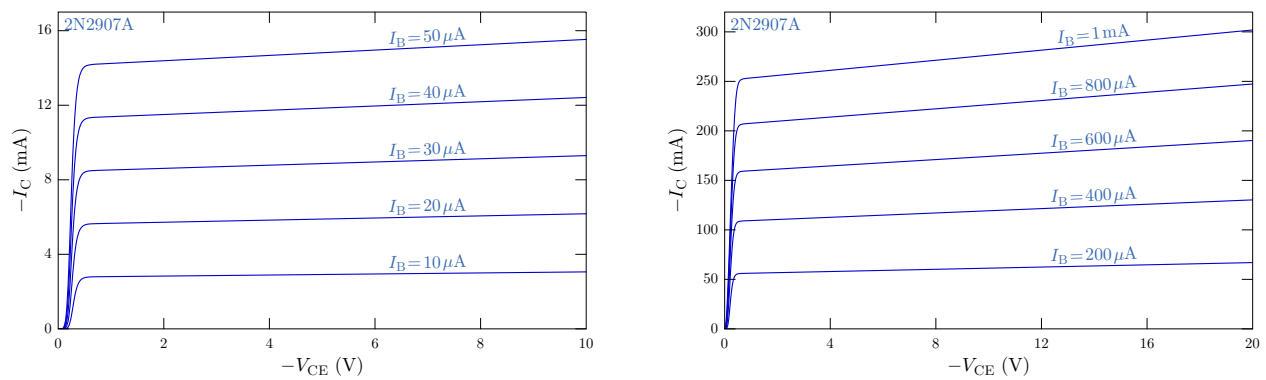
## C.2 Bipolar Junction Transistors

### C.2.1 Small-Signal BJTs

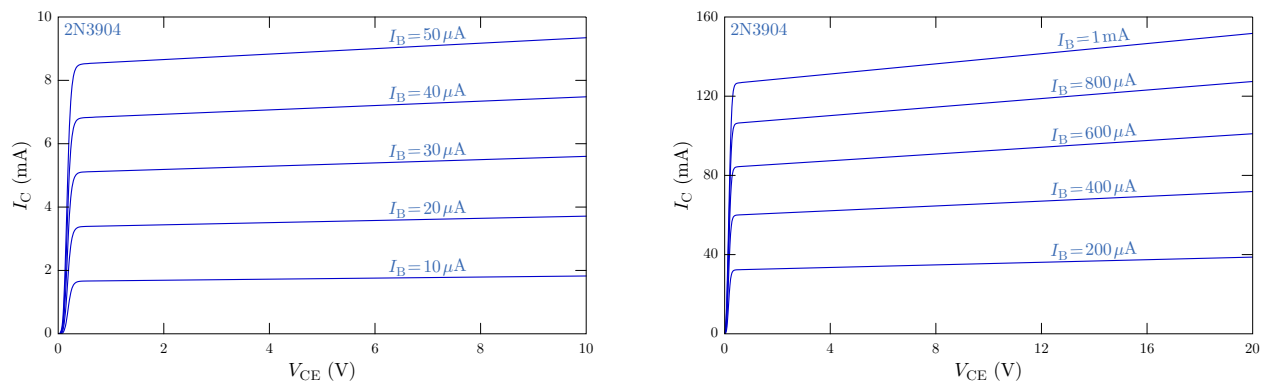
**2N2222A** (NPN, general purpose)<sup>8</sup> The “improved” 2N2222; also PN2222A and P2N2222A.



**2N2907A** (PNP, general purpose; complement to 2N2222A)<sup>9</sup>



**2N3904** (NPN, general purpose)<sup>10</sup>

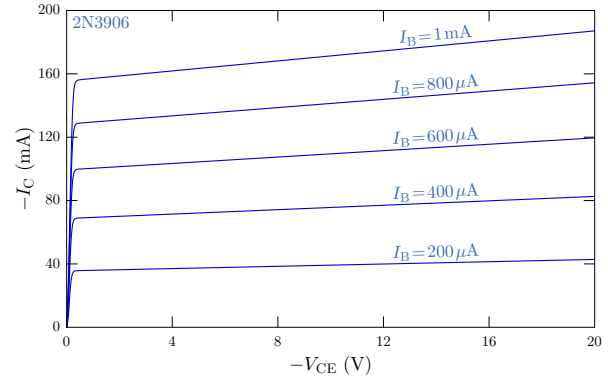
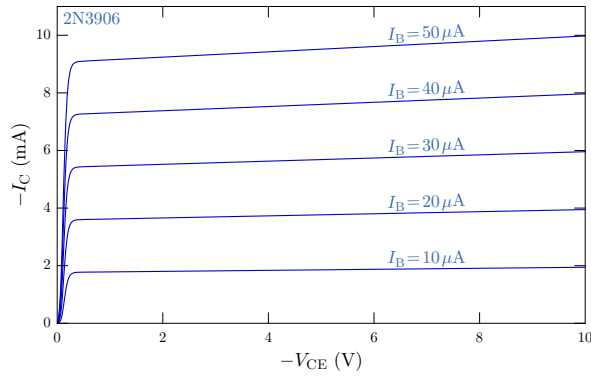
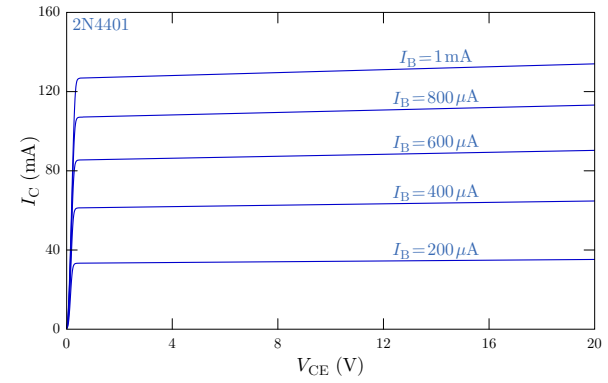
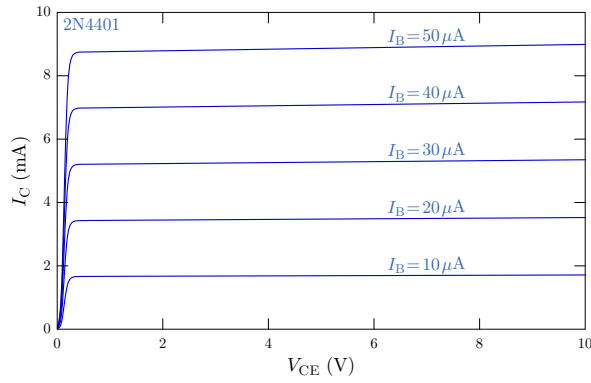
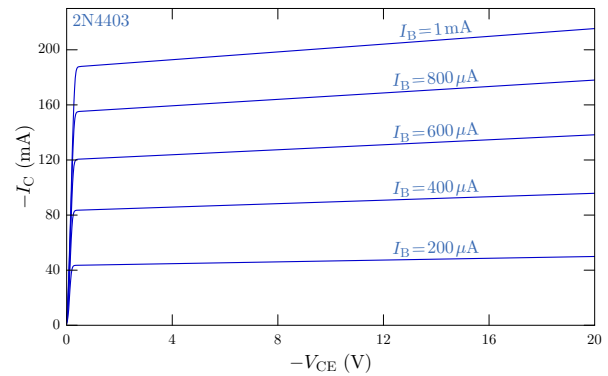
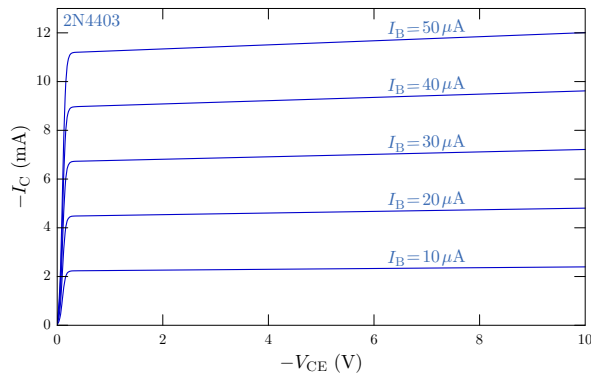


<sup>8</sup>Central Semiconductor model: <https://www.centralsemi.com/product/partpage2.php?part=2N2222A>

<sup>9</sup>Central Semiconductor model: <https://www.centralsemi.com/product/partpage2.php?part=2N2907A>

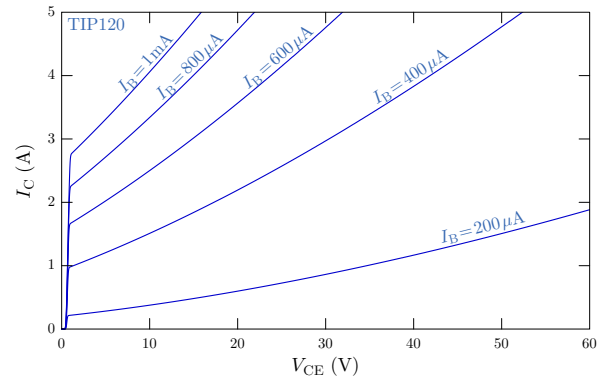
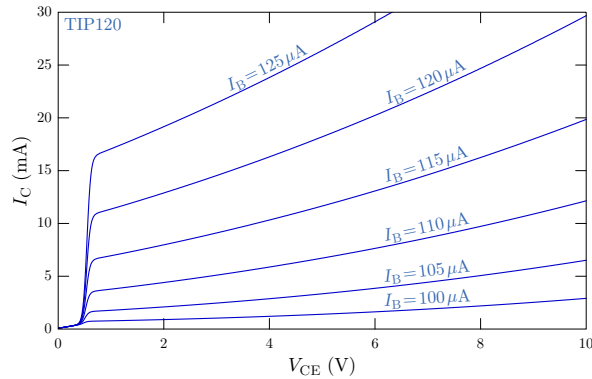
<sup>10</sup>Central Semiconductor model: <https://www.centralsemi.com/product/partpage2.php?part=2N3904>



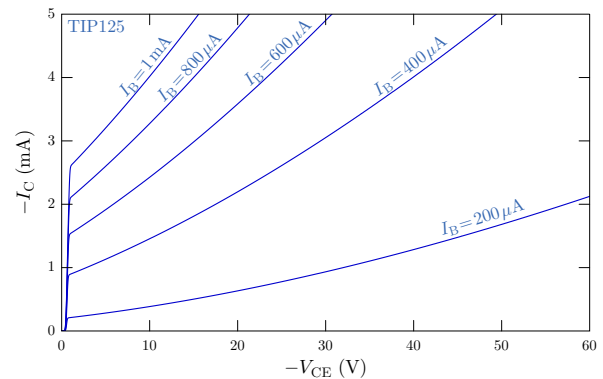
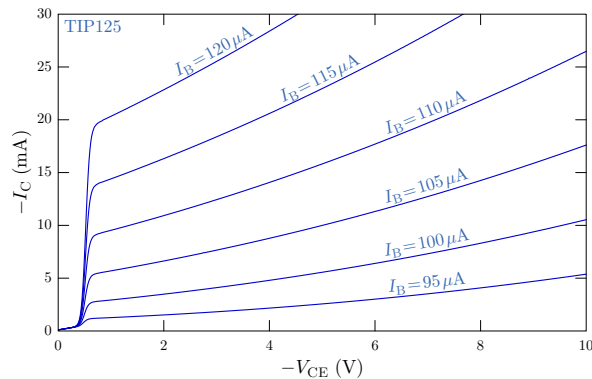
**2N3906** (PNP, general purpose; complement to 2N3904)<sup>11</sup>**2N4401** (NPN, general purpose)<sup>12</sup>**2N4403** (PNP, general purpose; complement to 2N4401)<sup>13</sup><sup>11</sup>Central Semiconductor model: <https://www.centralsemi.com/product/partpage2.php?part=2N3906><sup>12</sup>LTSpice library: <http://www.zen22142.zen.co.uk/ltspice/standard.bjt><sup>13</sup>LTSpice library: <http://www.zen22142.zen.co.uk/ltspice/standard.bjt>

### C.2.2 Power BJTs

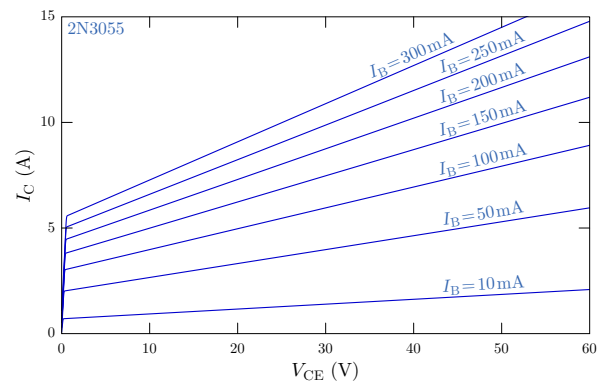
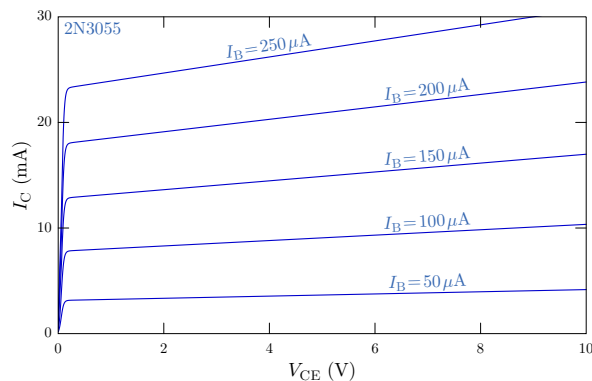
**TIP120** (NPN, 5 A, 60 V, 65 W Darlington power transistor)<sup>14</sup> Power transistors with  $\beta \sim 2500$  at  $I_C = 4$  A, in a TO-220 package.



**TIP125** (PNP Darlington power transistor; complement to TIP120)<sup>15</sup>



**2N3055** (NPN, 15 A, 60 V, 115 W power transistor)<sup>16</sup> Classic<sup>17</sup> power transistors with  $\beta = 20$ –70 at  $I_C = 4$  A, in a TO-3 package.



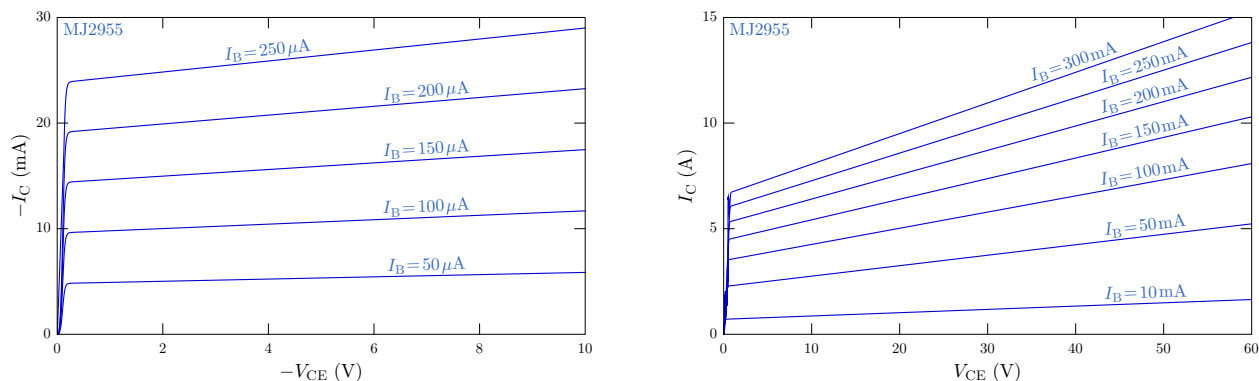
<sup>14</sup>ON Semiconductor model: <http://www.onsemi.com/PowerSolutions/product.do?id=TIP120>

<sup>15</sup>ON Semiconductor model: <http://www.onsemi.com/PowerSolutions/product.do?id=TIP125>

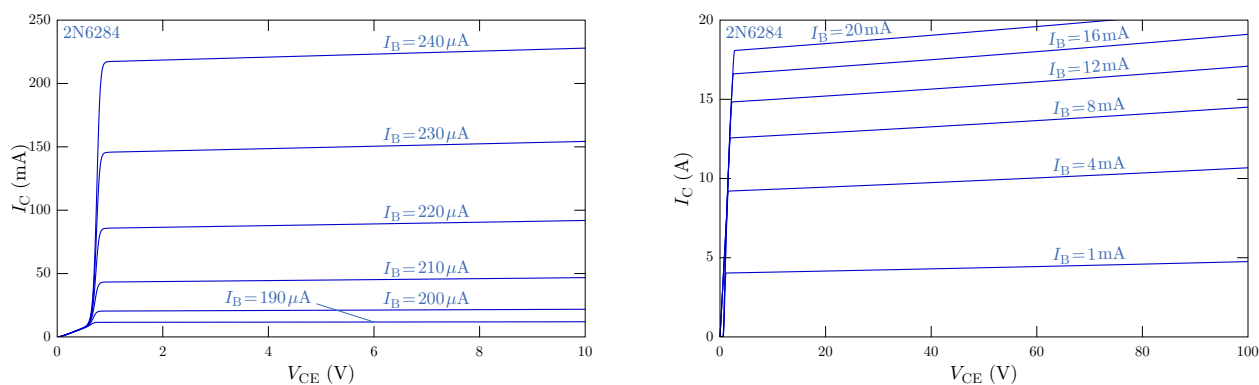
<sup>16</sup>ON Semiconductor model: <http://www.onsemi.com/PowerSolutions/product.do?id=2N3055>

<sup>17</sup>For historical information, see John N. Ellis, "The 2N3055: A Case History," *IEEE Transactions on Electron Devices* **48** 2477 (2001) (doi: 10.1109/16.960371).

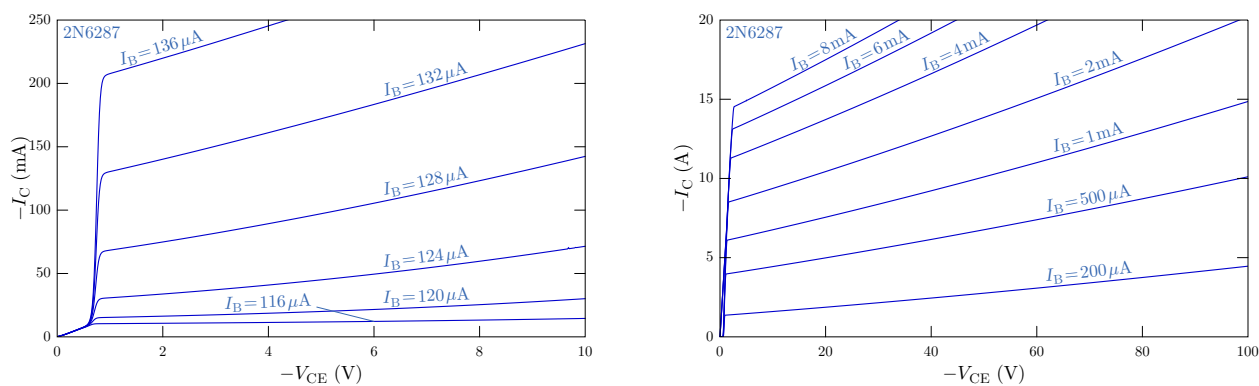
**MJ2955** (PNP, 15 A, 60 V, 115 W power transistor; complement to 2N3055)<sup>18</sup>



**2N6284** (NPN, 20 A, 100 V, 160 W Darlington power transistor)<sup>19</sup> Power Darlington transistors in a TO-3 package.



**2N6287** (PNP Darlington power transistor; complement to 2N6284)<sup>20</sup>



<sup>18</sup>ON Semiconductor model: <http://www.onsemi.com/PowerSolutions/product.do?id=MJ2955>

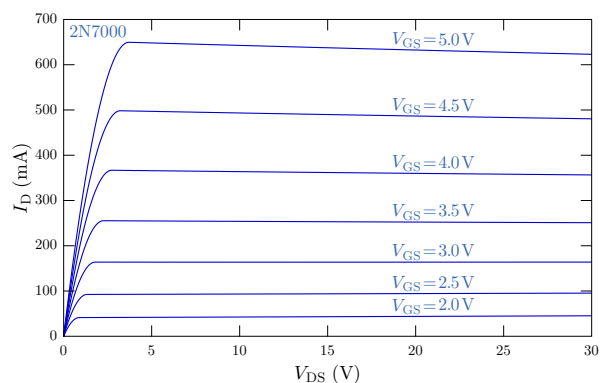
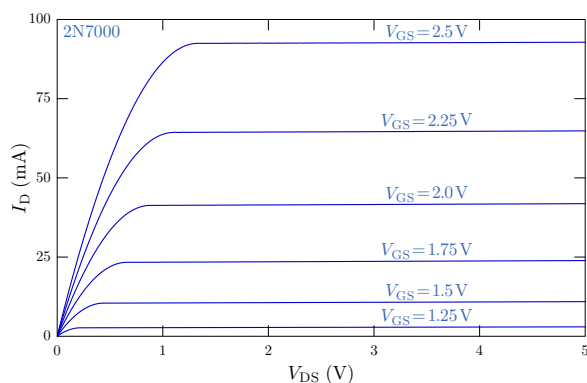
<sup>19</sup>ON Semiconductor model: <http://www.onsemi.com/PowerSolutions/product.do?id=2N6284>

<sup>20</sup>ON Semiconductor model: <http://www.onsemi.com/PowerSolutions/product.do?id=2N6287>

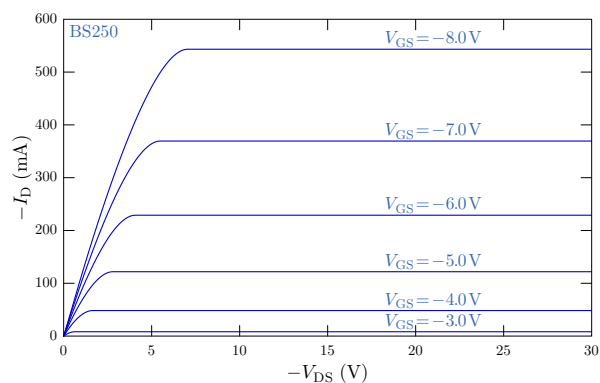
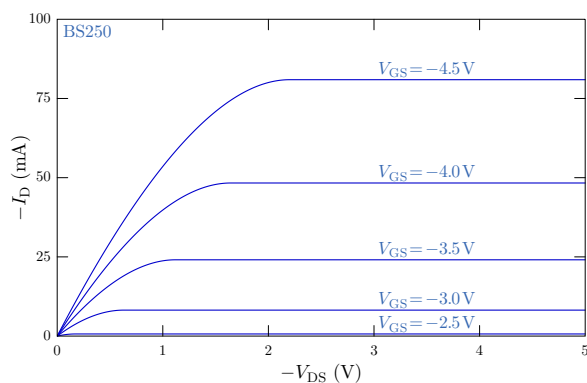
## C.3 MOSFETs

### C.3.1 Small-Signal MOSFETs

**2N7000** (N-channel, enhancement-mode, 200 mA, 60 V, small-signal MOSFET)<sup>21</sup>



**BS250** (P-channel, enhancement-mode, small-signal MOSFET, sometimes used as unofficial complement to 2N7000)<sup>22</sup>



<sup>21</sup>Model from Supertex, Inc., using “level 3” equations as given for example at <https://www.freedra.org/doc/elements/mosn3.pdf>.

<sup>22</sup>Model from Philips, Inc., using “level 3” equations as given for example at <http://web.rfoe.net:8000/ziliaoxiazai/PHILIPS/models/spicespar/data/bs250.html>.



# Index

- $1/f$  noise, 228
- $LC$  frequency, 66
- $Q$  factor, 66–68
- $RC$  time, 56
- $\bar{U}$ , 142
- AND gate, 272
  - 3-input, 273
- NAND gate, 272
- NOR gate, 272
- NOT gate, 271
- OR gate, 272
- XNOR gate, 273
- XOR gate, 272–273, 278–279
  - in terms of NAND gates, 279
- 1.5KE400A, 158
- 1N4001, 80
- 1N4733, 82
- 1N5711, 80
- 1N5819, 158
- 1N914B, 80
- 2's complement convention, 270–271
- 22V10, 335, 409
- 2N3904, 119–123
- 3-dB point, 60
- 4046 (PLL), 388
- 555 timer, 361–364, 367–368
  - astable multivibrator, 362–363
  - frequency modulation, 363–364
  - monostable multivibrator, 365–366
- 6116 SRAM, 330
- 74121 monostable multivibrator, 366–367
- 741C, 189, 198, 200, 206
- 74138, 310
- 74139, 316
- 74150, 310
- 74151, 309–310
- 74154, 310
- 74251, 311
- absolute-value amplifier
  - op-amp, 248–249
- absorption theorems, 278, 284
- active rectifier
  - op-amp, 241
- AD594, 313
- AD985L (DDS), 389
- ADC, 373, 376–380
- address bus, 330
- aliasing, 374
- amplifier
  - common-cathode, 170–177
- analog computer
  - op-amp, 235–236
- analog switch, 312
  - MOSFET, 150–151, 200
  - op-amp, 255
- analog-to-digital conversion, 373, 376–380
- analog-to-digital converter, 353
- anode, 77
- anti-aliasing filter, 374
- associative, 277
- astable multivibrator, 362–363
  - made from one-shots, 368–369
- ATF22V10C, 335, 409
- ATF750C, 336, 413
- band-pass amplifier
  - op-amp, 253–254
- bandwidth, 216–221
- beam power tube, 178
- bias network, 100–102
- BiCMOS, 301
- binary arithmetic, 269–271
  - power of 2, 274
- binary logic, 269
- binary operation, 277
- bipolar junction transistor, 91–137, 294
  - as switch, 95–96
  - current-control model, 108
  - Darlington pair, 136–137
  - forward-active mode, 96
  - hybrid model, 123
  - push-pull pair, 116–119
  - saturation, 95–96, 130
- bistable, 319
- bit, 269
- bit-shift operator, 275–276

- Boltzmann constant, 81
- Boolean algebra, 277–289
- bootstrapping, 129
- BUF634, 219
- buffer
  - level-shifting, 304–305
- buffer gate, 271
- capacitance, 53
- capacitor, 53–57, 73
- cascode amplifier, 128, 137
- cathode, 77, 167–168
  - directly heated, 167
  - filamentary, 167
  - heated, 167
  - indirectly heated, 167
- characteristic curves
  - bipolar transistor, 119–123
- charge, 23
- Child–Langmuir Law, 168
- chip select, 330
- Clapton, Eric, 127
- class-A amplifier, 117
- class-AB amplifier, 117
- class-B amplifier, 117
- clipping, 104
- closed-loop mode, 190
- CM600HA-24A, 157
- CMOS, 142, 300–303
- CMOS logic, 300–303
- CMOS switch, 312
- coaxial cable, 204
- Cockroft–Walton multiplier, 86–87
- common-cathode amplifier, 170–177
- common-emitter amplifier, 102–104, 114–116, 131–134
- common-mode gain factor, 107
- common-mode rejection ratio (CMRR), 107–108, 150, 201, 208
- common-mode signal, 105
- common-source amplifier, 147–149
- commutative, 277
- comparator, 189–190, 221–224
  - open-collector output, 222, 354
  - Schmitt trigger, 222–223, 355–358
- compensation, 219–221
- complex notation, 57–58
- complex PLD, 336
- compliance, 100
- conductance, 30–32, 41, 50–51
- control
  - integral, 264–265
  - PI, 265
  - PID, 265–266
  - proportional, 262–264
- control grid, 169
- control theory
  - linear, 261–266
- controller, 261
- counter
  - asynchronous, 323–324
  - backwards divide-by-3, 338–339
  - divide-by-2, 323
  - divide-by-2-or-3, 339–340
  - divide-by-3, synchronous, 326–327
  - divide-by-3-with-hold, 327–328, 332–333
  - divide-by-4, up/down, 328, 344
  - ripple, 323–324
- coupled oscillators, 68–71
- coupling coefficient
  - for inductors, 69
- CPLD, 336
- crossover distortion, 117, 218
- current, 23
- current mirror, 110–112, 133–134
  - improved Wilson, 129
- current source
  - for laser diode, 255–257
  - Howland, 247–248
  - JFET, 146
  - op-amp, 246, 255–259
  - transistor, 99–102
- DAC, 373–376, 380
- data bus, 330
- DDS, 389
- De Morgan’s theorems, 278
- debounced switch, 320–321
- delta-star transformation, 41–43
- demultiplexer, 310–316
  - analog, 311–315
- DEMUX, 310–316
- depletion zone, 79
- desert island, 316
- DG407, 312
- DG412, 312
- DIAC, 153–155
- differential amplifier
  - JFET, 149–150
  - op-amp, 194, 201, 246–247
  - transistor, 104–108
- differential gain factor, 106
- differential signal, 105
- differentiator, 73
  - op-amp, 194–197
- digital electronics, 269
- digital logic, 269

- digital-to-analog conversion, 373–376, 380
- dimmer
  - for lighting, 154–155
- diode, 77–90
  - capacitance, 88
  - forward voltage drop, 80
  - forward-biased, 77
  - ideal, 77
  - ideality factor, 81
  - reverse-biased, 77
  - reverse-breakdown voltage, 81
  - reverse-leakage current, 81
  - Schottky, 85–86
  - semiconductor, 77–82
  - TVS, 158
  - vacuum, 77, 167–169
  - vacuum full-wave rectifier, 168–169
  - Zener, 82–83, 88, 158
- diode law, 81–82
- diode logic, 292–293
- direct digital synthesis, 389
- distributive, 277
- DL, 292–293
- DRAM, 331
- droop, 264
- duty cycle, 364, 367–368
- dynamic RAM, 331
  
- Early effect, 113, 120–121, 124, 137
- Early voltage, 113
- Ebers–Moll model, 108–114, 119–123
- EEPROM, 331
- electromotive force (EMF), 23
- electronically erasable PROM, 331
- emitter degeneration, 115
- emitter follower, 96–99
- EPROM, 331
- erasable PROM, 331
- Eric Clapton Stratocaster, 127
- error, 262
- exponential amplifier
  - op-amp, 252–253
  
- Farad, 53
- feedback
  - in BJT amplifier(, 136
  - in BJT amplifier), 136
- feedback control, 261–266
- feedback signal, 262
- Fender Musical Instruments, 127
- FET, 139–151, 162–165
- FGA60N65SMD, 157
- field-effect transistor, 139–151, 162–165
  - CMOS, 142
  - IGFET, 141
  - JFET, 139–141, 162–164
  - MOSFET, 141–142
  - threshold voltage, 140
- filter
  - high-pass, 61–62
  - low-pass, 59–61, 63
- filters
  - op-amp, 194–200
- first-order loop, 388
- fixed-point notation, 269
- flicker noise, 228
- flip-flop, 319–323, 336–338
  - clocked, 321
  - D-type, 321–322, 328
  - D-type, edge-triggered, 322
  - debounced switch, 320–321
  - JK, 322–323
  - memory, 324–325
  - pulse-area stabilizer, 337–338
  - register, 324–325
  - ripple counter, 323–324
  - SR, 319
- floating-point notation, 269
- fluence, 200
- flyback transformer, 127
- frequency, 57
- frequency modulation
  - in 555 timer, 363–364
- full width at half maximum , 67
- full-wave rectifier, 84–85, 88
- fundamental charge, 81
  
- gain, 102
- gain factor
  - common-mode, 107
- gain factor differential, 106
- gain–bandwidth product (GBWP), 217
- GI754, 80
- goal, 261
- Gray code, 280
- grid, 169
  - control, 169
  - screen, 177
  - suppressor, 178
- ground loop, 204
- ground plane, 210
- grounded-emitter amplifier, 114–116
- grounded-emitter amplifier(, 135
- grounded-emitter amplifier), 136
- guitar preamp
  - op-amp, 238–241



- Gummel–Poon model, 121–123
- gyrator
  - op-amp, 236–238
- half-wave rectifier, 83–84, 88
- harmonic oscillator
  - damped, forced, 265
- hexagon from hell, 44
- high-pass filter, 61–62
  - cascaded, 74–75
- hole, 78
- Howland current source, 247–248
- hysteresis, 319
- IGBT, 155–161
- IGFET, 141
- impedance, 59
- impedance-matching condition, 30
- inductive load
  - transistor switch, 124–125
- input bias current, 198–200, 206–208
- input impedance, 29
  - through transistor, 97–99
- input offset current, 208
- input offset voltage, 200
- instrumentation amplifier
  - ac-coupled, hi- $Z$  input, 205–206
  - differential receiver, 204–205
  - op-amp, 200–206
  - thermocouple, 203–204
- insulated-gate bipolar transistor, 155–161
- integrating factor, 55
- integrator, 54–56
  - op-amp, 195–200
- intrinsic emitter resistance, 108, 110
- inverter, 158–161, 271
  - Tesla coil, 160–161
- inverting amplifier
  - op-amp, 191–192, 214–216
- IRF1405, 157
- IXYN80N90C3H1, 157, 158
- JFET, 139–141, 162–164
  - common-source amplifier, 147–149
  - current source, 146
  - differential amplifier, 149–150
  - source follower, 146–147
- Johnson–Nyquist noise, 226–228
- joule thief, 125–126
- junction, 78
- Karnaugh map, 280–285, 288–289
- Kirchoff’s laws, 24–25
- Lagrange multiplier, 34
- Lambert  $W$  function, 90
- level-shifting buffer, 304–305
- LF411, 189
- linear algebra, 30–32, 50–51
- little- $h$  notation, 123–124
- LM311, 222, 223
- LM399, 82
- logarithmic amplifier
  - op-amp, 247
- logic
  - three-state, 311
- logic gates, 271–273, 291–307
- long-tailed pair, 105
- low-pass filter, 59–61
  - in PLL, 387
  - inductor, 73–74
  - phase, 63
  - two-pole capacitor–inductor, 75
- maximum-value amplifier
  - op-amp, 249–250
- memory, 324–325, 328–331, 339
  - state machines with, 331–336
- mho ( $\Omega$ ), 142
- microprocessor, 333
- Miller effect, 113–114, 134–135, 137
- Millman’s theorem, 43
- mixer, 389
  - rf, 387
- monostable multivibrator, 364–367
  - 74121, 366–367
- MOSFET, 141–142, 297–303
- MR752, 80
- multiplexer, 309–316
  - analog, 311–315
- multivibrator
  - astable, 362–363
  - monostable, 364–367
- mutual inductance, 68
- MUX, 309–316
- n-type carrier, 78
- n-type semiconductor, 78
- negative feedback, 104, 190
- negative-feedback mode, 190
- negative-impedance converter
  - op-amp, 250–251
- NMOS, 298–299
- NMOS logic, 298–299
- noise immunity, 104
- non-Ohmic, 77
- noninverting amplifier

- op-amp, 192–193, 211–214
- nonlinear, 77
- Norton's theorem, 37
- notch filter, 76
- Nyquist frequency, 373
- octave, 60
- Ohm's law, 23–24
- ohmic contact, 85
- one-shot, 364–367
- op-amp, 112
  - current feedback, 208, 212
  - golden rules, 190
  - voltage feedback, 208, 212
- op-amps, 127–129, 163–164, 189–259
  - stability, 257–259
  - transistor internals of AD829, 127–129
  - transistor internals of TL082, 163–164
- OPA111B, 189
- OPA602C, 198, 200
- open-collector output, 222, 354
- open-loop gain, 189
  - finite, 211–216
- open-loop mode, 189
- operational amplifier
  - absolute-value amplifier, 248–249
  - active rectifier, 241
  - analog computer, 235–236
  - band-pass amplifier, 253–254
  - current source, 246, 256–259
  - differential amplifier, 194, 201, 246–247
  - differentiator, 194–197
  - exponential amplifier, 252–253
  - filters, 194–200
  - guitar preamp, 238–241
  - gyrator, 236–238
  - Howland current source, 247–248
  - instrumentation amplifier, 200–206
  - integrator, 195–200
  - inverting amplifier, 191–192, 214–216
  - logarithmic amplifier, 247
  - maximum-value amplifier, 249–250
  - negative-impedance converter, 250–251
  - noninverting amplifier, 192–193, 211–214
  - phase-shift oscillator, 224–225
  - photodiode amplifier, 245–246
  - pulse-area stabilizer, 241–243
  - relaxation oscillator, 223–224
  - single-supply, 209
  - stability, 192
  - summing amplifier, 193–194
  - transimpedance amplifier, 245–246, 252
  - unity-gain buffer, 190–191
- operational amplifiers, 127–129, 163–164, 189–259
  - bandwidth, 216–221
  - comparator, 221–224
  - compensation, 219–221
  - noise analysis, 225–234
- output admittance, 113, 123, 133
- output enable, 330
- output impedance, 29
  - through transistor, 97–99
- output swing, 189
- p-n junction, 78
- p-type carrier, 78
- p-type semiconductor, 78
- P6KE20CA, 158
- PAL, 326
- path to ground
  - dc, 199
- pentode
  - triode connection, 183–184
- phase detector, 386
  - type I, 386
  - type II, 386
- phase shift, 62–63
  - and power, 63–65
- phase-locked loop, 385–392
- phase-shift oscillator
  - op-amp, 224–225
- photodiode, 245–246
  - amplifier instability, 252
- PID control, 261–266
- pink noise, 228
- PLA, 326
- plant, 261
- plate
  - vacuum-tube, 167
- plate resistance, 173
- PLD, 326, 335–336
- PLL, 385–392
- PMOS, 299–300
- PMOS logic, 299–300
- potential, 23
  - difference, 23
- power, 24
- power factor, 65
- power-supply rejection ratio (PSRR), 209
- printed circuit board (PCB), 210
- programmable array logic, 326
- programmable logic array, 326
- programmable logic devices, 326, 335–336
- programmable ROM, 331, 333–335
- PROM, 331, 333–335
- proportional–integral (PI) control, 265

- proportional–integral–derivative (PID) control, 265
- pulse-area stabilizer
  - op-amp, 241–243, 337–338
- pulse-width modulation, 364
- push–pull amplifier, 116–119
- push–pull amplifier, 217–219
- quiescent current, 111
- race condition, 282–283, 289
- RAM, 328–331
- random-access memory, 328–331
- reactance
  - capacitive, 58
  - inductive, 58–59
- read enable, 330
- read-only memory, 331
- rectifier
  - active, 241
  - full-wave, 84–85, 88
  - half-wave, 83–84, 88
- register, 324–325
  - shift, 325
- relaxation oscillator
  - op-amp, 223–224
- resistance
  - plate, 173
- resistor network
  - matrix formalism, 30–38, 50–51
- resistor-transistor logic, 293–295
- resistors, 24
  - parallel, 25–26
  - series, 25
  - voltage divider, 26–29
- resonant circuit, 65–68, 72
  - coupled, 68–71
- resonant frequency, 66
- ripple counter, 323–324
- rms, 64
- ROM, 331
- RTL, 293–295
- S401E (SCR), 152
- sample, 269
- sampling rate, 373
- sampling theorem, 373
- saturation current, 81
- schmactivated, 76
- schmapacitor, 75–76, 245
- schmesistor, 47, 244
- schmesonator, 76
- Schmitt trigger, 222–223, 355–358
- Schmohm’s law, 47, 244
- Schottky barrier, 85
- SCR, 151–153
- screen grid, 177
- second-order loop, 388
- secondary emission, 178
- Seebeck effect, 203
- semiconductor-controlled rectifier, 151–153
- sequential logic, 325–328, 331–336, 338–340
- shift register, 325
- Shockley, William, 81
- siemens, 142
- sign-magnitude convention, 270
- signed integer, 270–271
- simple PLD, 335
- slew rate, 217
- source follower
  - JFET, 146–147
- SPICE, 81
- SPLD, 335
- SR flip-flop, 319
- SRAM, 328–331
- state diagram, 327–328
- state machine, 325–328, 331–336, 338–340
- static RAM, 328–331
- summing amplifier
  - op-amp, 193–194
- superposition theorem, 37–38
- suppressor grid, 178
- switch, 291–292
  - debounced, 320–321
  - SPDT, 292
  - SPST, 291
- tank circuit, 65
- Tesla coil, 68, 72
- tesla coil
  - solid-state, 126–127
- Thévenin’s theorem, 27–30, 32–40
  - proof, 32–38
- thermal runaway, 118
- thermal voltage, 108
- thermistor, 253
- thermocouple, 312–315
- three-state logic, 311
- thyatron, 153
- thyristor, 151–155
- transconductance, 109, 141, 142
- transfer function
  - feedback-control loop, 262–263
- transformer, 68
- transient-voltage suppressor, 158
- transimpedance amplifier
  - op-amp, 252

- transistor
  - $\alpha$  parameter, 93
  - $\beta$  parameter, 93
  - bipolar junction, 91–137, 294
  - Darlington pair, 136–137
  - forward current ratio, 123
  - hybrid model, 123
  - input impedance, 123
  - insulated-gate bipolar, 155–161
  - MOSFET, 297–303
  - output admittance, 123
  - reverse voltage ratio, 123
- transistor switch
  - inductive load, 124–125
- transistor-transistor logic, 295–297
- transistor-transistor logic (TTL), 269
- TRIAC, 153–155
- triode
  - vacuum, 169–177
- truth table, 271
- TTL, 295–297
- tube, vacuum, 167–187
- TVS, 158
- two's complement convention, 270–271
- type I phase detector, 386
- type II phase detector, 386
  
- UCC3732x, 158
- uninterruptible power supply, 159
- unity-gain bandwidth, 217
- unity-gain buffer
  - op-amp, 190–191
- unsigned integer, 269–270
  
- vacuum tube, 167–187
  - diode, 167–169
  - triode, 169–177
- VCO, 387
- voltage, 23
- voltage divider, 26–29, 41
- voltage follower
  - op-amp, 190–191
- voltage multiplier
  - Cockroft–Walton, 86–87
- voltage-controlled oscillator, 387
  
- Wheatstone bridge, 253
- white noise, 227
- work function, 85
- write enable, 329, 330
  
- XKCD, 50–51
  
- Zener diode, 82–83, 88, 158







767180T  
Hnn9210PE

74LS27N 8038  
5A

MTS 7  
8120

MC14520B  
CP 0731

SN74LS42N  
CPKG0025

49AVFCM  
MAX232N

CD4002AE  
RCA 440

MC14020  
BCP RR8346

7606  
SN7427N

M9024  
2706M  
N6339N

R56-1D.5-12  
NTE ELECTRONICS, INC.

425  
3008038-00

MC1496L  
7522  
KOREA

LM5622AN  
CD4022AE

LT  
0749  
LTC7541AKN  
U67033

74HC74N  
801170T  
Hnn9210PE

GS104KOREA  
GL324

P3068  
945

M9024  
LM339N

LM565CN  
340

P3068  
945

RC8120MP  
RAY 7426

142 09  
274 006

6CA7EXT  
TLO44ACN

LM339N

142 09  
274 006

AZ 7602  
SN74LS190N

949  
CD40208CN  
MM5620BN

SN74LS04N  
XAH9440

4G3 T  
HD74LS175P

1008

MC14011B