Giovanni Colombo (Ed.)

CIME Summer Schools

# Modern Questions of Celestial Mechanics

43

## Bressanone, Italy 1967

Springer

Giovanni Colombo (Ed.)

# Modern Questions of Celestial Mechanics

Lectures given at a Summer School of the
Centro Internazionale Matematico Estivo (C.I.M.E.),
held in Bressanone (Bolzano), Italy,
May 21-31, 1967

Springer

FONDAZIONE
**CIME**
ROBERTO CONTI

CENTRO INTERNAZIONALE MATEMATICO ESTIVO

(C. I. M. E.)

1° - Ciclo - Bressanone dal 21 al 31 maggio 1967

"MODERN QUESTIONS OF CELESTIAL MECHANICS"

Coordinatore: G. COLOMBO

# SUL PROBLEMA DELLE AURORE BOREALI
## MOTO DI UN CORPUSCOLO ELETTRIZZATO IN PRESENZA DI
## UN DIPOLO MAGNETICO

di

### C. AGOSTINELLI

(Università di Torino)

Introduzione.

1 --- Quei meravigliosi fenomeni celesti che vanno sotto il nome di auro
re boreali, le quali, con una fantasmagoria di luci e di colori appaiono e si
contemplano nelle lunghe notti polari, sono determinati come si sa da scia-
mi di particelle elettrizzate provenienti da regioni cosmiche e principal-
mente dalle macchie solari.

Come è ben noto le macchie del Sole sono immense cavità formatesi
nella fotosfera dell'altro e dalle profondità di queste voragini vengono lan-
ciati all'esterno getti potentissimi di gas di idrogeno, elio , calcio ecc.,
e insieme ad essi, essendo tali gas altamenti ionizzati per la loro ele-
vata temeratura, vengono lanciati elettroni liberi, cioè corpuscoli elettriz-
zati negativamente, con una velocità che in condizioni medie si aggira in-
torno ai 1600 km/$_{sec.}$

Questi corpuscoli possono, con quella velocità, percorrere in circa
trenta ore il tragitto dal Sole alla Terra, la cui distanza è di circa 150
milioni di Km.

Una volta che detti corpuscoli elettrizzati abbiano raggiunta l'atmo-
sfera tèrrestre avviene la scarica con emissione di luce, e questo costi-
tuisce appunto il fenomeno dell'aurora boreale.

Ora, se la Terra non esercitasse nessuna azione sul movimento di
questi corpuscoli elettrizzati, essi potrebbero colpirla indifferentemente
in tutte le regioni della sua superficie. Ma, essendo la Terra dotata di un
campo magnetico, quando essi giungono in vicinanza del nostro globo so-
no costretti ad abbandonare il loro cammino rettilineo per descrivere del-
le traiettorie a forma di spirale, che ordinariamente convergono verso il

polo boreale, essendo quei corpuscoli carichi di elettricità negativa .

Soltanto in casi eccezionali, quando i corpuscoli sono estremamente veloci, ciò che avviene quando il Sole è in un periodo di grande attività, le spirali si allargano e allora la Terra può essere colpita anche in latitudini piuttosto basse. Ciò spiega perchè la grande aurora boreale che si verificò nella notte del 25 gennaio 1938 fu visibile in molta parte di Europa ed anche qui in Italia.

I periodi di emissione più intensa dei corpuscoli elettrizzati si susseguono a distanza di 11 anni e sono legati ai periodi di maggiore attività delle macchie solari, il cui numero e la cui estensione diviene massima, come si sa, ogni undici anni.

Un fenomeno analogo a quello che determina le aurore boreali ha luogo per le radiazioni cosmiche,consistenti anch'esse in particelle elettrizzate lanciate con velocità altissime. Anche queste particelle subiscono l'azione del campo magnetico terrestre, sebbene in misura minore, a causa della loro maggiore velocità . Le radiazioni cosmiche possono pertanto osservarsi in ogni luogo della Terra, pur essendo più intense nelle regioni polari.

I primi ad interessarsi di questi fenomeni furono i due astronomi scandinavi Störmer [1] e Birkeland , che effettuarono una serie di osservazioni e di esperienze molto interessanti.

2.--- Ammesso dunque che i fenomeni di aurore boreali siano determinati da corpuscoli elettrizzati provenienti da regioni attive del Sole, lo studio rigoroso di essi è subordinato alla risoluzione analitica del problema del moto di un corpuscolo elettrizzato in presenza di un dipolo magnetico, quale è appunto il dipolo costituito dalla coppia dei poli magnetici terrestri, è subordinato cioè all'integrazione delle equazioni differenziali che reggono un tale movimento.

C. Agostinelli

Queste equazioni si possono integrare completamente nel caso di un solo polo magnetico [3], caso che interessa lo studio dei raggi catodici, e in questo caso le traiettorie sono le geodetiche di un cono di rotazione col vertice nel polo. Ma esso si può utilizzare anche per lo studio del moto dei corpuscoli elettrizzati in prossimità del polo boreale ad una distanza sufficientemente piccola in confronto della distanza dei due poli terrestri [5].

Ma l'integrazione completa delle equazioni del moto nel caso del dipolo si presenta estremamente difficile.

Alla risoluzione numerica di quelle equazioni furono dallo Störmer e dai suoi allievi dedicati molti lavori, e dopo una mole ingente di calcoli, durati circa 5000 ore, e consumati quintali di carta, egli riuscì a costruire alcuni modelli, in filo di ferro, delle curve spirali descritte dai corpuscoli elettrizzati lanciati dalle macchie solari.

Ma i calcoli istituiti dallo Störmer sono fondati però sulla considerazione di un campo magnetico dovuto a un magnete elementare, cioè di lunghezza infinitesima, posto nel centro della Terra, coll'asse coincidente coll'asse magnetico terrestre, avente un momento magnetico uguale a quello della Terra, cioè di circa $8,52.10^{25}$ unità magnetiche. Se questa semplificazione può portare a risultati sufficientemente approssimati per il moto dei corpuscoli elettrizzati lanciati dalle macchie solari, quando la distanza di quei corpuscoli dalla Terra ecceda, come afferma lo stesso Störmer, un milione di Km, essa non è più evidentemente accettabile, quando si voglia spingere la conoscenza del moto di quei corpuscoli in prossimità della Terra.

Di qui la necessità di una indagine basata su una schematizzazione più aderente alla realtà, come è quella fornita dalla considerazione di un dipolo magnetico di lunghezza finita, quale è appunto il dipolo terrestre.

Partendo da questo punto di vista, in seguito alla segnalazione fat-

C. Agostinelli

fa dal Prof. Armellini su "Scienze e Tecnica" $\begin{bmatrix} 2 \end{bmatrix}$ , in occasinne della
grande aurora bpreale del 1938, mi sono occupato anch'io della questione
in diversi lavori $\begin{bmatrix} 4 \end{bmatrix}$ , e mi propongo ora di riferire sui risultati più es-
senziali in essi conseguiti.

### Equazioni del moto di un corpuscolo elettrizzato nel campo di un dipolo magnetico. Integrali primi ed equazioni ridotte.

3. --- Supponendo che i corpuscoli elettrizzati non siano soggetti ad altre for-
ze che a quelle di un campo magnetico , trascurando l'azione del peso e la
loro mutua influenza, l'equazione differenziale vettoriale del moto di uno di
tali corpuscoli P, elettrizzato negativamente, è notoriamente

(1)
$$m \frac{d\vec{v}}{dt} = e \; \vec{v} \wedge \vec{H}$$

dove $\underline{m}$ è la massa del corpuscolo, $\vec{v}$ la velocità ; $\underline{e}$
l'intensità della sua carica elettrica ed $\vec{H}$ l'intensità del campo magnetico.

Nel caso di un dipolo, indicando con $\varphi$
il potenziale magnetico, risulta

(2)
$$\vec{H} = grad \; \varphi, \quad \varphi = k \left( \frac{1}{r_1} - \frac{1}{r_2} \right)$$

essendo $\underline{k}$ la costante del dipolo ed $r_1$
$r_2$ le distanze del corpuscolo dai due po-
li $O_1$ , $O_2$ .

La (1) ammette il noto integrale delle forze vive

(3)
$$v^2 = v_0^2 \; (cost.)$$

Inoltre, se si indica con $\psi$ la funzione che uguagliata a costante for-
nisce nel semipiano $P(O_1 \; O_2)$ le linee di forza magnetica, cioè le tra-
iettorie ortogonali delle linee equipotenziali $\varphi$ = cost. , sussiste

C. Agostinelli

il nuovo integrale primo

(4)
$$y^2 \dot{w} = c - \mu \Upsilon \, , \quad \dot{w} = \frac{dw}{dt} \, ,$$

che ho stabilito per la prima volta in una nota dell'Accademia delle Scienze di Torino del $1^o$ giugno 1938. In esso y è la distanza del punto P dell'asse polare $O_1 O_2$, w è l'angolo che il semipiano mobile $P(O_1 O_2)$ forma con un semipiano fisso passante per lo stesso asse, $\mu$ è una costante data da

$$\mu = \frac{2 e \mathcal{K}}{m}$$

infine <u>c</u> è la costante d'integrazione. L'integrale (4) esprime la velocità areolare del corpuscolo intorno all'asse polare in funzione della sua posizione, e pertanto la costante <u>c</u> si può chiamare <u>costante delle aree.</u>

4. --- Se si assume come origine delle coordinate il punto medio O del segmento $O_1 O_2$, l'asse polare come asse x, positivo nel verso da $O_1$ ad $O_2$, e si indica con <u>a</u> la semidistanza polare (raggio terrestre), in coordinate x, y risulta

(5)
$$\Upsilon = \frac{1}{2} \left( \frac{x + a}{r_1} - \frac{x - a}{r_2} \right), \quad r_1 = \sqrt{(x + a)^2 + y^2} \, ,$$

e in tutto il campo del moto si ha

$$r_2 = \sqrt{(x - a)^2 + y^2}$$

(6)
$$0 \leqslant \Upsilon \leqslant 1$$

In virtù dell'integrale (4) la risoluzione del problema si riduce all'integrazione delle equazioni differenziali

(7)
$$\frac{d\dot{x}}{dt} = -\frac{1}{2} \frac{\partial}{\partial x} \left( \frac{c - \mu \Upsilon}{y} \right)^2 \, , \quad \frac{d\dot{y}}{dt} = -\frac{1}{2} \frac{\partial}{\partial y} \left( \frac{c - \mu \Upsilon}{y} \right)^2$$

che sono le equazioni del moto relativo del corpuscolo P nel semipiano

C. Agostinelli

$P(O_1 O_2)$ , ed esprimono che detto moto relativo equivale a quello del mo-
to di un punto, nello stesso semipiano, soggetto a una forza posizionale
e conservativa derivante dal potenziale

(8)
$$\mathcal{U} = - \frac{1}{2} \frac{(c - \mu \psi)^2}{y^2}$$

Le equazioni (7) ammettono l'integrale delle forze vive

(9)
$$\dot{x}^2 + \dot{y}^2 + \frac{(c - \mu \psi)^2}{y^2} = v_o^2 \quad (cost.) ,$$

il quale non è altro che l'integrale (3) quando si tenga conto dell'integra-
le (4) delle aree.

Dalla (9) risulta che durante il moto è sempre

(10)
$$\frac{(c - \mu \psi)^2}{y^2} \leqslant v_o^2$$

e questa, avendo riguardo alla (8), dimostra che le traiettorie relative, cor-
rispondenti a un dato valore della costante c, sono contenute in un campo
$\sigma_c$ del semipiano P x, descritto dalle linee di livello $U = - \frac{1}{2} k^2 v_o^2$ ,
con $k^2$ parametro costante compreso fra zero ed uno, .

Ne nasce quindi l'opportunità dello studio dei campi $\sigma_c$ , e delle
linee di livello che li limitano, al variare della costante c da $-\infty$ a $+\infty$ .
Detti campi si presentano con caratteristiche profondamente diverse a se-
conda che sia

$$c < o , \qquad o \leqslant c \leqslant \mu , \qquad c > \mu .$$

Le linee limiti di questi campi generano intorno all'asse polare del-
le superficie di rotazione che possono essere toccate dai corpuscoli ma
non attraversate, e che corrispondono, nel caso delle radiazioni cosmiche
alle superficie riflettenti di questi raggi .

5. In virtù della (6) si ha che per

(11)              $c < o$       , e per $c > \mu$ ,

è in tutto il campo del moto

(11')
$$c - \mu \ddot{y} \neq o .$$

In questi casi la (10) dimostra che durante il moto si ha sempre
$$y > 0 ,$$

cioè il corpuscolo non taglia mai in alcun punto l'asse polare e le equazioni (7) sono regolari in tutto il capo del moto.

In base a questo risultato, nei casi in cui ora ho accennato, se si considerano nelle equazioni (7) gli elementi incogniti $x, y, \dot{x}, \dot{y}$ come funzioni analitiche del tempo $\underline{t}$ , riguardando questo come variabile complessa, e si costruiscono, prendendo il valore iniziale del tempo come centro, quelle figure che il Mittag-Leffler chiama stelle , si può dimostrare che l'asse reale dei tempi giace nell'interno di quelle stelle. Pertanto si possono ottenere senz'altro , col metodo di Mittag-Leffler, gli elementi incogniti sviluppati per tutti i valori reali del tempo e ciò in infiniti modi e colla sola conoscenza delle condizioni iniziali del moto.

Però per le serie che così si ottengono non è assicurata l'uniforme convergenza fintantochè si sappia solo che l'asse reale dei tempi è contenuto nelle stelle degli elementi incogniti.

Ma in una Memoria dell'Accademia delle Scienze di Torino [4] ho dimostrato che si può costruire una striscia di larghezza finita, contenente l'asse reale dei tempi, limitata da due rette ad esso parallele e tutta contenuta nelle stelle suddette.

Gli sviluppi del Mittag-Leffler restano perciò, nel nostro problema, molto semplificati, ed ho potuto pertanto assegnare le funzioni incognite sviluppate in serie di potenze, uniformemente convergenti per tutti i valori reali del tempo da $-\infty$ a $+\infty$ .

6--- Nel caso invece in cui la costante $\underline{c}$ delle aree è interna allo intervallo $( 0, \mu )$ , cioè si abbia

(12)
$$o \leqslant c \leqslant \mu ,$$

per la (6) può risultare durante il moto

(12') $$c - \mu z \rho = 0 \, ,$$

e pertanto la (10) non esclude che in qualche istante si abbia  $y = 0$ .

L'Analisi   riguardante i campi   $\mathfrak{S}_c$   in cui si possono svolgere le traiettorie relative, dimostra appunto che nel caso in cui la costante   $\underline{c}$ è compresa nell'intervallo   (12)  le traiettorie relative possono tagliare l'asse polare. In tal caso le equazioni (7) del moto risultano irregolari per  $y = 0$ , e   pertanto non si può a   priori ammettere l'esistenza di una soluzione analitica di esse.

Ma con un'opportuna trasformazione le dette equazioni si possono regolarizzare e nel caso in cui   $o < c < \mu$ , pur   estendendosi il campo   $\mathfrak{S}_c$ in cui si   possono svolgere le traiettorie relative fino ai poli   $O_1, O_2$ , ho mostrato che il corpuscolo non può colpire nessuno di  detti poli.

Nel caso limite di   $c = \mu$ , il campo   $\mathfrak{S}_c$ , o una parte di esso, si estende fino   al seguento  $O_1 O_2$  che congiunge i due poli. In tal caso ho dimostrato come sia possibile che delle traiettorie intersechino   l'asse polare in un punto del segmento'  $O_1 O_2$  .

Nel caso infine di   $c = 0$  il campo   $\mathfrak{S}_c$   è limitato da una linea che ha per  asintoto   l'asse polare , e pertanto in tal caso il corpuscolo può toccare l'asse polare soltanto  nei   punti all'infinito  di esso .

## Linee di livello e campi in cui si svolgono le traiettorie

### relative

7. --- Le traiettorie  effettive del corpuscolo saranno evidentemente contenute nello spazio di rivoluzione   $S_c$  , ottenuto dalla rotazione del campo  $\mathfrak{S}_c$  intorno all'asse polare   $x$  .

Tutte le traiettorie possibili si potranno quindi distribuire in una se-

rie infinita di famiglie, ciascuna delle quali corrisponde a un dato valore della costante $c$ , e le traiettorie di una di queste famiglie non escono dallo spazio $S_c$ corrispondente .

I punti di una di queste famiglie di traiettorie , corrispondenti a un dato valore della costante $c$ , per i quali è

(13) $$\frac{(c - \mu \psi)^2}{y^2} = k^2 v_0^2, \quad (0 \leqslant k^2 \leqslant 1)$$

sono situati su una superficie di rotazione intorno all'asse x, che è generata dalla linea del semipiano Px che ha per equazione la (13) , la quale non è altro che la linea di livello corrispondente al valore $-\frac{1}{2} k^2 v_0^2$ del potenziale U che compete al moto relativo del corpuscolo.

Le linee di livello (13) , corrispondenti a uno stesso valore della costante $c$ , per le quali è $0 \leq k^2 \leq 1$ , descrivono evidentemente il campo $\mathfrak{S}_c$ innanzi definito.

Ciò premesso , chiameremo curve limiti di livello quelle corrispondenti ai valori estremi $k^2 = 0$ e $k^2 = 1$ del parametro $k^2$ , e che vengono a limitare il detto campo, nel quale, ripeto, sono contenute le traiettorie relative corrispondenti a uno stesso valore della costante $c$ .

La curva limite di livello corrispondente al valore $k^2 = 1$ , risulta :

(13') $$\frac{(c - \mu \psi)^2}{y^2} = v_0^2$$

e da questa, per l'equazione (9) delle forze vive, si deduce che se una traiettoria corrispondente a un dato valore della costante $c$ , ha un punto M$_a$ comune con la curva limite di livello (13') , in quel punto è $\dot{x} = o$ , $\dot{y} = o$ , cioè esso è un punto nel quale la velocità relativa è nulla, mentre la velocità assoluta del corpuscolo è, nel punto M, diretta normalmente al piano Mx . Si ha inoltre che la traiettoria relativa uscente dal punto M è ortogonale in M alla curva limite di livello (13') . In-

C. Agostinelli

Invero, essendo nel punto M, $\dot{x} = 0$, $\dot{y} = 0$ , la traiettoria relativa è in quel punto tangente alla linea di forza del moto relativo, e perciò ortogonale alla linea di livello corrispondente.

8. --- Per vedere ora l'andamento delle linee di livello (13), e in particolare delle curve limiti di livello (13') , nei casi di

(14)        $c \lessgtr 0$      ,    $0 < c < \mu$   ,        $c \geqslant \mu$ ,

poniamo per semplicità

(15)     $\gamma = \dfrac{c}{\mu}$ ,     $\nu_0 = \dfrac{a \, v_0}{\mu}$ ,  $\xi = \dfrac{x}{a}$ , $\eta = \dfrac{y}{a}$

con

       $-\infty < \xi < +\infty$   ,  $\eta \geqslant 9$ ,

e osserviamo che, avendo le costanti $c$ e $\mu$ le stesse dimensioni, risultano $\gamma$ e $\gamma_0$ numeri puri, con $\gamma$ variabile da $-\infty$ a $+\infty$, e $\nu_0 > 0$ .

La (13) diventa allora

(16)      $(\gamma - \psi)^2 = k^2 \nu_0^2 \, \eta^2$

che si scinde nelle due equazioni

(17)          $\gamma - \psi = k \, \nu_0 \, \eta$

(17')        $\psi - \gamma = k \, \nu_0 \, \eta \, , \; 0 \leqslant k \leqslant 1 \, , \; \nu_0 > 0 \, , \; \eta \geqslant 0$

e in queste si dovranno considerare separatamente i casi

(14')      $\gamma \leqslant 0$  ,  $0 < \gamma < 1$  ,  $\gamma \geqslant 1$

E' facile vedere che nel piano $(\xi, \eta)$ le due curve (17) e (17') sono curve algebriche del $12^0$ ordine, una simmetrica dell'altra rispetto all'asse $\xi$ .

a) Nel caso in cui è

(18)      $\gamma = -\gamma_1 < 0$ ,  $(c < 0)$

C. Agostinelli

il ramo (17) della linea di livello (16) non ha punti reali, e resta da consi-
derare soltanto il ramo definito dalla (17'), che ora diventa

(19) $$Y + \gamma_1 = k \, v_0 \, \eta, \quad (0 \leqslant k \leqslant 1)$$

Questa linea di livello è , nel piano $(\xi, \eta)$ , simmetrica rispetto all'asse
$\eta$ , e si riconosce che essa non ha alcun punto comune coll'asse $\xi$
e pertanto le traiettorie del corpuscolo, come ho già detto, non toccheran-
no l'asse polare.

La stessa linea am-
mette come asintoto la
retta

(20) $$\eta = \frac{\gamma_1}{k \, v_0}$$

$Y = -0,5$
$v_0 = 0,2$

parallela all'asse $\xi$.

Le curve limiti di li-



Fig. 1

vello sono quelle che si otten-
gono dalla (19) per k=o e k=1 ; per k=o si ha come linea di livello la
retta $\eta = \infty$ (retta all'infinito del semipiano $\eta > 0$) , e per k=1 si ha la
linea di livello

(19') $$Y + \gamma_1 = v_0 \, \eta$$

Il cui asintoto dista dall'asse $\xi$ di $\gamma_1 / v_0$ , e si trova che questa linea
è tutta situata al disopra di detto asintoto. Essa ha l'andamento della fig. 1,
che è stata costruita per $Y = -\gamma_1 = -0,5$; e $v_0 = 0,2$ .

6) Se ora supponiamo che sia $Y = 0$ , (c=o) , la (17) è soddisfatta soltan-
to per $\gamma = 0$ e $\eta = 0$ , mentre la (17') diventa

(21) $$\psi = k \, v_0 \, \eta$$

In questo caso le linee di livel-
lo sono tutte asintotiche all'asse
$\xi$. Esse hanno l'andamento
della fig. 2 , dove è rappresen-
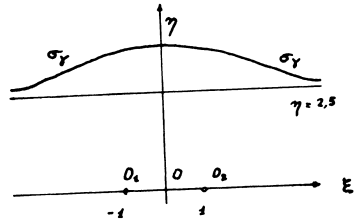tata una linea limite di livello

$Y = 0$
$v_0 = 0,1$


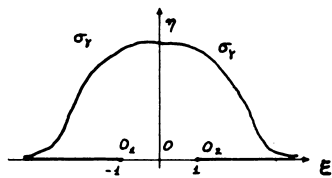
fig. 2

(k=1) , per $\gamma$= o e $\gamma_o$ = 0 , 1 .

Le traiettorie relative del corpuscolo, svolgendosi nella parte del semipiano $\eta \geqslant 0$, esterno alla curva limite $2\gamma = \gamma_o \eta$, possono avere in comune coll'asse polare $\xi$ soltanto i punti all'infinito di esso .

E' da osservare che la (21) è soddisfatta anche per $2\gamma$ = o e $\eta$ = o , e pertanto questa curva comprende anche la parte dell'asse polare esterna al dipolo, nei cui punti è $|\xi| \geqslant$, 1 , $\eta$ = 0 , cioè $|x| \geqslant a$ , y = 0 . Questo è dovuto al fatto, come del resto si deduce dalle equazioni del moto, che per c = 0 , ($\gamma$= o) , è possibile il moto del corpuscolo lungo l'asse polare, partendo da distanza infinita verso il polo positivo, oppure partendo dal polo negativo e allontanandosi indefinitamente.

9--- Consideriamo ora il caso in cui sia

$$\gamma > 1 , \quad ( c > \mu ) .$$

In questo caso il ramo (17') delle linee di livello (16) non ha punti reali, e resta da considerare soltanto il ramo

(17) $$\gamma - \gamma = k \nu_o \eta$$

Le linee di livello che sono rappresentate da questa equazione sono simmetriche rispetto all'asse $\eta$ , ed hanno ancora come asintoti le rette

$$\eta = \frac{\gamma}{k \gamma_o} , \quad (o \leqslant k \leqslant 1 ) ,$$

le quali vanno allontanandosi indefinitamente dall'asse $\xi$ col tendere di k a zero.

La linea limite corrispondente al valore zero del parametro $\underline{k}$ è la retta all'! all'infinito del semipiano $\eta > 0$ .

La linea limite

(22) $$\gamma - \gamma = \nu_o \eta$$

che si ha per k = 1 , può presentare andamenti differenti, come mostrano

le figure 3, 4, 5, 6, 7, dipendente-
mente dai valori che assumono
le costanti $\gamma$ e $\nu_o$ , e quindi dai
punti d'intersezione che quella
linea può avere coll'asse $\eta$
che sono dati dalle radici della
equazione

$$(23 \quad F(\eta) \equiv \nu_o \, \eta + \frac{1}{\sqrt{1+\eta^2}} - \gamma = 0$$

In questo caso, essendo il primo
membro della (22) sempre maggio-
re di zero, si ha che nel campo
in cui si svolgono le traiettorie
è sempre $\eta > 0$ , e quindi le tra-
iettorie del corpuscolo non pos-
sono toccare in alcun punto l'as-
se polare.

E' notevole il caso rappre-
sentato dalla figura 5 in cui la curva
limite di livello (22) ha un ramo
aperto asintotico alla retta $\eta = \gamma/\nu_o$ ,
e un ramo chiuso posto tra il pre-
cedente e l'asse polare.

Entrambi i rami, simmetrici
rispetto all'asse $\eta$ , vengono a li-
mitare due campi entro i quali si
possono svolgere le traiettorie re-
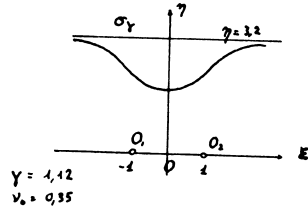lative del corpuscolo.

Il primo campo $\sigma_\gamma'$ si estende
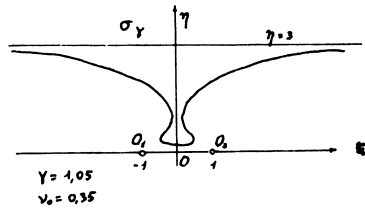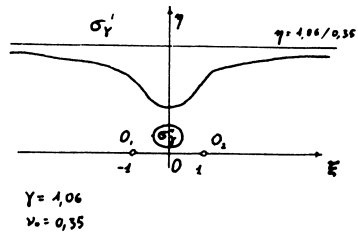


$\gamma = 1,12$
$\nu_o = 0,35$

Fig. 3



$\gamma = 1,05$
$\nu_o = 0,35$

Fig. 4



$\gamma = 1,06$
$\nu_o = 0,35$

Fig. 5

all'infinito mentre l'altro $\mathfrak{S}_\gamma''$
è un campo chiuso finito. A
quest'ultimo campo corrisponde
nello spazio, un campo $S_\gamma''$ anula-
re che è di rotazione intorno all'as-
se polare, entro il quale sono realiz-
zabili dei moti periodici del corpuscolo
elettrizzato, moti messi in evidenza
dalle esperienze di Villard.



$\gamma = \frac{3\sqrt{2}}{4} \qquad \nu_o = \frac{\sqrt{2}}{4}$

Fig. 6

10. --- Passiamo ora a considerare il
caso in cui si abbia

$$o < \gamma < 1 , \ (o < c < \mu) .$$

In questo caso i due rami (17) e (17') di
una linea di livello, corrispondente ad
un dato valore di k compreso fra zero e
uno, sono entrambi reali e tangenti nei
poli $O_1, O_2$ alla linea di forza magnetica
$\mathcal{Y} = \gamma$, che passa anch'essa per quei punti
(fig. 8).



$\gamma = 4\sqrt{2} / (3\sqrt{3})$
$\nu_o = 2/(3\sqrt{3})$

Fig. 7

Inoltre il ramo (17) ha come asinto-
to la retta di equazione $\eta = \gamma / (k \, \gamma_o^2)$ , e
in ogni punto di questo ramo è $\eta \leqslant \gamma / (k \gamma_o^2)$
cioè esso si svolge tutto al disotto dello
asintoto corrispondente.

Il ramo (17') è posto invece interna-
mente alla linea di forza magnetica $\mathcal{Y} = \gamma$.

Le traiettorie in questo caso posso-
no toccare l'asse polare soltanto nei poli
$\Theta_1, O_2$ .



$\gamma = 0,4$
$\nu_o = 0,2$

Fig. 8

C. Agostinelli

Anche qui si possono presen-
tare diverse conformazioni dei campi
in cui si svolgono le traiettorie
relative a seconda dei valori delle
costanti $\gamma$ e $\gamma_o$ . Nelle figure 8,
9 e 10  sono rappresentati i tre
tipi  di campi $\sigma_\gamma$ che si possono
avere.

Quando si presentano configu-
razioni  della forma delle figure 9
e 10 , si possono realizzare dei
moti periodici del corpuscolo le
cui traiettorie possono toccare lo
asse polare  nei poli $O_1$ e $O_2$ .

11. --- Consideriamo infine il caso
in cui $\gamma = 1$ , (c $= \mu$) . Se ciò avvie-
ne le (17) e (17') per  k = 1 diventa-
no

(24)        $1 - 2f = \gamma_o \eta$

(24')     $2f - 1 = \gamma_o \eta$, $(\eta \geqslant o)$ .

Essendo $o \leqslant 2f < 1$ , la (24') è sod-
disfatta soltanto per $2f = 1$ e $\eta = o$ , la
qual cosa sussiste  per i punti del segmento
mento $O_1O_2$ che congiunge i due poli.

Questa  congiungente nel caso di
$\gamma = 1$ fa perciò parte  delle curve limi-
ti di livello.



$\gamma = 0,35$
$\gamma = 0,2$

Fig. 9



Fig. 10



$\gamma = 1$
$\nu_o = 0,5$

Fig. 11

La (24) è soddisfatta anco-
ra per $\frac{2}{3} = 1$ e $\eta = 0$ (segmento $O_1$
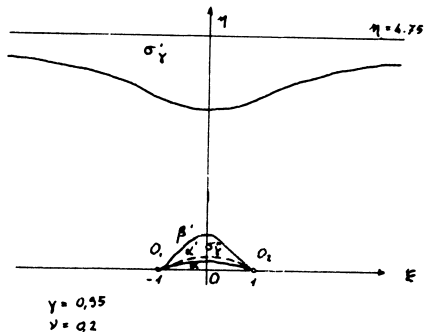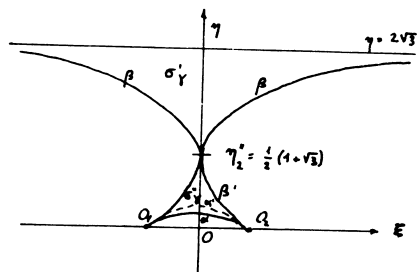$O_2$) , ma essa rappresenta inoltre
due rami uscenti rispettivamente dai
poli $O_1$, $O_2$, ivi tangenti all'asse pola-
re, simmetrici rispetto all'asse $\eta$ ,
e asintotici alla retta $\eta = 1/\gamma_0$ .

In particolare se $\gamma_0 >$
$> 2/( 3\sqrt{3} )$ il campo $\sigma_\gamma$ si presen-
ta come nella figura 11, si estende
cioè dall'infinito fino al segmento
$O_1 O_2$ dell'asse polare e perciò le
traiettorie del corpuscolo, provenien-



Fig. 12

ti da distanza infinita possono intersecare l'asse polare nei punti della
congiungente i due poli.


Se invece è $\gamma_0 < 2 / (3\sqrt{3})$ la linea $\beta$ si può spezzare in un
ramo aperto $\beta'$, asintotico alla retta $\eta = 1/\gamma_0$ , al disotto della
quale esso si svolge, e un ramo $\beta''$ che termina nei poli $O_1 O_2$
ed ivi tangente all'asse polare . Il campo $\sigma_\gamma$ si scompone in due :
il primo $\sigma_\gamma'$ va dall'infinito fino alla linea $\beta'$ ; il secondo $\sigma_\gamma''$
è finito ed è compreso fra la linea $\beta''$ e il segmento $O_1 O_2$ . Le traiet-
torie relative che si svolgono nel campo $\sigma_\gamma''$ possono in questo caso
intersecare l'asse polare. I due campi $\sigma_\gamma'$ , $\sigma_\gamma''$ risultano separati
da un campo intermedio entro il quale non è possibile il moto re-
lativo del corpuscolo. Detti campi hanno in questo caso una confi-
gurazione analoga a quella della figura 9, colla sola differenza che qui
la linea $\alpha$ coincide col segmento $O_1, O_2$ .

C. Agostinelli

In particolare, per un opportuno valore della costante $\gamma_0$ , le due linee $\beta'$ , $\beta''$ possono avvicinarsi fino a formare una cuspide in un punto dell'asse $\eta$ , come mostra la fig. 12.

Per quanto concerne l'applicazione alla teoria delle aurore boreali, in base al valore medio della velocità dei corpuscoli lanciati dal Sole verso la Terra e al valore del momento magnetico terrestre, si trovano per le costanti $\gamma$ e $\gamma_0$ i valori medi

$$0 < \gamma \approx 35 \cdot 10^{-6} < 1 \ , \quad \gamma_0 \approx 1,5 \cdot 10^{-6} < 2/(3\sqrt{3}).$$

Applicando allora le considerazioni svolte si ha che nel caso delle aurore boreali il campo $\mathfrak{S}_\gamma$ , entro il quale si svolgono le traiettorie relative, ha la forma della fig. 8, dove l'asintoto della linea $\beta$ ha per ordinata $\eta = \gamma/\gamma_0$ = 23,3 , e l'ordinata del punto d'intersezione della linea $\alpha$ coll'asse $\eta$ risulta uguale a circa 816 . Cioè queste ordinate risultano rispettivamente uguali a circa 23,3 e 816 volte il raggio terrestre.

### Risoluzione analitica del problema. Caso in cui la costante $\underline{c}$ delle aree è esterna all'intervallo$(0,\mu)$

12.--- Nel caso in cui la costante $\underline{c}$ delle aree è esterna all'intervallo $(0,\mu)$ , se cioè risulta

(25) $\qquad c < 0 , \qquad$ oppure $\qquad c > \mu ,$

ricordando che la funzione $\mathcal{2f}$ è in tutto il campo del moto compresa fra zero e 1, si ha sempre durante il moto

$$c - \mu \, \mathcal{2f} \neq 0 .$$

Allora, dovendo essere per l'equazione delle forze vive $\dfrac{c - \mu \, \mathcal{2f}}{y}$ una quantità finita, la $y$ non potrà mai annullarsi, e dalla condizione

C. Agostinelli

$$\frac{(c - \mu \dot{y})^2}{y^2} \leqslant v_o^2$$

si ricava

(26) $$y \geqslant \frac{|c - \mu \dot{y}|}{v_o} \ .$$

In queste condizioni nella mia Memoria dell'Accademia delle Scienze di Torino ho dimostrato che esiste un numero reale positivo $\tau_o$ tale che se nel piano della variabile complessa $\underline{t}$ si costruisce una striscia di larghezza $2\tau_o$ , limitata da due rette parallele all'asse, reale e distanti da esso da ambo le parti della quantità $\tau_o$, questa striscia risulta tutta interna alle <u>stelle</u> di Mittag-Leffler e relativa alle variabili $x, y, \dot{x}, \dot{y}$ . Queste variabili sono cioè , in questa striscia, delle funzioni olomorpfe di $t$ senza alcuna singolarità .

Il numero $\tau_o$ è dato da

$$\tau_o = \frac{y_1}{3 v_o (1 + k)}$$

dove

$$y_1 = -\frac{c}{v_o} \qquad , \qquad y_1 = \frac{c - \mu}{v_o}$$

e $\underline{k}$ è la radice positiva dell'equazione

$$3 v_o^2 y_1^2 k (k + 1) - 17 (|c| + 4\mu)^2 = 0$$

Con la trasformazione conforme

(27) $$T = \frac{e^{\frac{\pi}{2} \frac{t}{\tau_o}} - 1}{e^{\frac{\pi}{2} \frac{t}{\tau_o}} + 1}$$

la striscia anzidetta si trasforma, nel piano $T$, in un cerchio di raggio unitario col centro nell'origine e gli sviluppi delle funzioni $x, y$ risultano della forma :

(28)
$$x = A_0 + A_1 T + \frac{1}{2!} A_2 T^2 + \dots$$

$$y = B_0 + B_1 T + \frac{1}{2!} B_2 T^2 + \dots$$

dove i coefficienti $A_i$, $B_i$ , i = 0, 1, 2, ... sono determinabili, e perciò noti, in base alle condizioni iniziali del moto.

Così ad esempio si ha

$$A_0 = (x)_{T=0} = (x)_{t=0} = x_0 \ , \quad A_1 = \left(\frac{dx}{dt}\right)_{t=0} \cdot \left(\frac{dt}{dT}\right)_{t=0} , \ldots$$

$$B_0 = (y)_{T=0} = (y)_{t=0} = y_0 \ , \quad B_1 = \left(\frac{dy}{dt}\right)_{t=0} \cdot \left(\frac{dt}{dT}\right)_{t=0} , \ldots$$

Gli sviluppi (28), tenendo conto della (27), sono uniformemente convergenti per valori reali del tempo da t = = - ∞ a t = + ∞, e rappresentano,
nel caso considerato, l'integrale generale delle equazioni del moto relativo
del corpuscolo. Determinati x ed y in funzione di t, l'angolo w che il
semipiano mobile P x forma con un semipiano fisso, passante anch'esso per l'asse polare, si ottiene infine con una quadratura utilizzando l'integrale delle aree; si ha cioè

$$w = w_0 + \int_0^t \frac{c - \mu y}{y^2} \, dt$$

Il problema del moto di un corpuscolo elettrizzato in presenza di un
dipolo magnetico, nel caso in cui la costante c delle aree soddisfa a una
delle condizioni (25) , è così , dal punto di vista analitico completamente
risolto.

Regolarizzazione delle equazioni del moto e risoluzione del problema nel caso in cui la costante delle aree
è compresa nell'intervallo (0, $\mu$).

13.--- Abbiamo visto come nel caso in cui si abbia

(29)
$$0 \leqslant c \leqslant \mu$$

in qualche istante del movimento può risultare    $y = 0$, cioè il corpusco-
lo può toccare l'asse polare.

Essendo le equazioni  del moto relativo singolari per $\dot{y} = 0$, o più
in particolare par $r_1 = 0$, oppure $r_2 = 0$, non si può  più senz'altro ammet-
tere l'esistenza di una  soluzione analitica del problema, quando la costan-
te $\underline{c}$ soddisfa alla condizione  (29).

Possiamo però regolarizzare le dette equazioni introducendo una
nuova variabile indipendente $\tau$ definita dalla relazione differenziale

(30)
$$d\tau = dt/y^2$$

Così  facendo si riconosce facilmente che il problema del moto re-
lativo del corpuscolo, definito dalle equazioni differenziali (7), si trasfor-
ma, mediante la (30), in un altro problema dinamico  nelle stesse incogni-
te x, y, nel quale la forza viva $\mathcal{T}$, e il  potenziale $\mathcal{U}$, sono rispet-
tivamente espressi dalle relazioni

(31)
$$\mathcal{T} = \frac{1}{2y^2}(x'^2 + y'^2), \quad \mathcal{U} = \frac{1}{2}y^2\left\{v_0^2 - \frac{(c-\mu\psi)^2}{y^2}\right\}$$

dove
$$x' = \frac{dx}{d\tau}, \quad y' = \frac{dy}{d\tau}$$

e la costante delle forze vive è nulla.

Ponendo ora

(32)
$$p_1 = \dot{x}, \quad p_2 = \dot{y}$$

si ha

(33)
$$x' = p_1 y^2, \quad y' = p_2 y^2$$

e la risoluzione del sistema di equazioni differenziali (7) equivale alla
integrazione del seguente sistema di equazioni differenziali

$$\frac{dx}{d\tau} = p_1 y^2, \qquad \frac{dy}{d\tau} = p_2 y^2$$

(34)
$$\frac{dp_1}{d\tau} = \mu \sqrt{v_o^2 - p_1^2 - p_2^2} \; y \frac{\partial \psi}{\partial x}$$

$$\frac{dp_2}{d\tau} = y \left( v_o^2 - p_1^2 - p_2^2 \right) + \mu \sqrt{v_o^2 - p_1^2 - p_2^2} \; y \frac{\partial \psi}{\partial y}$$

Questo sistema è regolare sia per y=0, e sia per $r_1$ = 0, oppure $r_2$ = 0, come si riconosce osservando che sono regolari i prodotti

$$y \frac{\partial \psi}{\partial x} \; , \quad y \frac{\partial \psi}{\partial y}$$

Ora , nel caso in cui la costante delle aree  c  è compresa  fra zero e  $\mu$  , $(0 < \gamma < 1)$ , estremi esclusi, il campo  $\sigma_\gamma$  si estende, come abbiamo già visto, fino ai poli  $O_1 O_2$, senza avere  altri punti in comune coll'asse polare. Allora se si  ammette che vi sia un istante in cui il corpuscolo vada a colpire per esempio   il polo  $O_1$  , si dimostra che in quell'istante la terza delle equazioni (34) non è soddisfatta. Ciò vuol dire che nel caso  in cui    $0 < c < \mu$  le traiettorie del corpuscolo non possono passare per nessuno dei poli  $O_1, O_2$ . Si ha cioè  il fatto pa-radossale che in questo caso un corpuscolo può avvicinarsi  quanto si vuole a uno dei poli, per poi allontanarsene, senza toccare quel polo.

Non è possibile perciò  determinare in questo caso un limite inferiore per  la y, e quindi non è possibile assegnare una soluzione analitica del problema analoga a quella che si ha per c < 0 , oppure per   $c > \mu$ .

Però data la regolarità dei  secondi membri delle equazioni  (34) in tutto  il campo del  moto, il teorema di  Cauchy  dell'Analisi assicura che fissati dei valori reali arbitrari   $x_o, y_o, p_{10}, p_{20}$  delle funzioni inco-gnite (tali da soddisfare all'equazione delle forze vive) e corrispondenti a un valore reale  $\tau_o$  di  $\tau$  ,  facendo variare    x, y , $p_1$ e $p_2$ rispettiva-

mente in opportuni intorni di $x_o$, $y_o$, $P_{10}$, $P_{20}$ , esiste un intorno $\mathscr{R}$ di $\tau_o$ entro il quale quelle funzioni sono sviluppabili in serie di potenze di $\tau - \tau_o$ , convergenti in $\mathscr{R}$ .

Osserviamo che nel caso limite di $c = \mu$ , ($\gamma = 1$) , avendo il campo $\mathscr{G}_f$ a comune coll'asse polare il segmento $\overline{O_1 O_2}$ , le traiettorie relative possono intersecare l'asse polare in un punto intermedio del segmento che congiunge i due poli $O_1$, $O_2$ .

L'istante $t_1$ in cui ciò avviene corrisponde al valore infinito del parametro $\tau$ definito dalla (30) , e nell'intorno di $\tau = \infty$ le funzioni x,y sono sviluppabili in serie di potenze della forma

$$x = a_o + \frac{a_1}{\tau} + \frac{a_2}{\tau^2} + \cdots$$

$$y = \frac{b_1}{\tau} + \frac{b_2}{\tau^2} + \cdots$$

mentre il tempo sarà dato dalla

$$t - t_1 = \int_\infty^\tau y^2 d\tau .$$

In particolare per $c = \mu$ il corpuscolo può percorrere il segmento $O_1 O_2$ con velocità costante $v_o$ .

E' da rilevare infine che tra le traiettorie del corpuscolo vi è anche la parte dell'asse polare esterna al dipolo, dove è

$$|x| \geqslant a , \qquad y = 0 , \qquad \gamma = 0 .$$

Ciò è possibile nel caso di $c \neq 0$ , e in tal caso il corpuscolo si può muovere sull'asse polare partendo da distanza infinita verso il polo positivo, oppure partire dal polo negativo e allontanarsi indefinitamente.

C. Agostinelli

Bibliografia

[1] C. Störmer, Sur les trajectoires des corpuscules électrisés dans l'espace sous l'action du magnétisme terrestre avec application aux aurores boréales ( "Archives des sciences phisiques et naturelles''. T. XXIV, Ginevra 1907) .
Idem , Les aurores boréales ("Livre du cinquantenaire de la Soc. franc. de Phys." , Paris, 1925) .

[2] G. Armellini, Le aurore boreali ("Scienza e Tecnica" , Rivista della S. I. P. S. , vol. $2^o$ , fasc. $3^o$, marzo 1938)

[3] C. Agostinelli e A. Pignedoli, Meccanica razionale (vol. I, Cap. V, § 2, n. 7 , Zanichelli , Bologna 1961) .

[4] C. Agostinelli, Sul moto di un corpuscolo elettrizzato in presenza di un dipolo magnetico ("Atti della R. Accademia delle Scienze di Torino'' vol. 73, 1937-38) .

Idem , Sulla risoluzione analitica del problema del moto di un corpuscolo elettrizzato in presenza di un dipolo magnetico ("Memorie della R. Accademia delle Scienze di Torino", Serie $2^a$ , Tomo 69, P. I, 1938-39)

Idem , Sul moto di un corpuscolo lettrizzato in un campo magnetico simmetrico rispetto a un asse ecc. ("Atti della R. Accademia delle Scienze di Torino", vol. 74, 1938-39)

[5] A. Pignedoli, Sul problema delle aurore polari. Moto di un corpuscolo elettrizzato in presenza di un dipolo magnetico e in prossimità di uno dei poli ("Atti del Seminario Matematico e Fisico dell'Università di Modena" , vol. I a. 1967 .

CENTRO INTERNAZIONALE MATEMATICO ESTIVO

(C. I. M. E.)

G. COLOMBO

"INTRODUCTION TO THE THEORY OF EARTH'S MOTION ABOUT ITS
CENTER OF MASS"

Corso tenuto a Bressanone dal 21 al 31 maggio 1967

# INTRODUCTION TO THE THEORY OF EARTH'S MOTION ABOUT ITS CENTER OF MASS

by

G. COLOMBO

(Università di Padova)

Introduction

The motion of the earth about its center of mass has become in the recent past and will become in the near future more and more complex with the increasing accuracy of measurement. Besides more and more sophisticated will be the model needed to explain this motion . To begin correctly we have to define exactly : a) what we inted by earthsystem, b) how to choose the reference system, c) how to describe the motion of the earth about its center of mass.

Concerning a) we should first define where we intend to put the external closed boundary of the earth : on the sea and solid earth surface? at an altitude of 200 Km above sea level where the density is almost constant in time and position ? at the boundary (not yet very well defined) of the magnetoshpere? Suppose we assume the boundary at 200 Km above seal level. The center of mass is a very well defined geometrical point, independent of the reference system, but its position is not fixed with respect to any observatory of the earth surface because the earth is not a rigid body .

Concerning b), for defining exactly a reference system there are evidently no difficulties. In fact we may choose a reference system completely attached to an observatory of the earth (geographical system) or to e fixed stars for the orientation of the axes, or we may choose the instantaneous center of mass of the earth and the instantaneous central axes of inertia of the earth (inertial reference system). We meet with almost insuperable difficulties when we have to relate the different systems to each other and the local geographical system with the inertial reference system.

G. Colombo

The general equation

(1)
$$\frac{d\vec{S}_G}{dt} = \vec{M}_G^{(e)}$$

holds, $S_G$ being the total angular momentum vector with respect to G and $\vec{M}_G^{(e)}$ the torque due to all external forces. The reference system has to be nearing to a Newtonian system. In expressing $S_G$ we meet with the difficulty, predicted in c), arising from the relative motion of the different constituents of the earth as defined above (core, mantle, ocean, atmosphere). Since we are not dealing with a rigid body it is impossible to define the axis of instantaneous rotation of the system earth. It is generally assumed that, apart from the motion of ocean and atmosphere relative to the solid part, the motion of the solid crust are of low frequency and small amplitude. Following Jeffreys we may define as the instantaneous angular velocity vector of the earth the vector $\omega$ which minimizes the following expressions

(2)
$$(\vec{\omega} \times \vec{l} - \vec{v})^2 \, \mu \, d\mathcal{C}$$

where $\vec{l}$ is the vector from G to the mass element $d\mathcal{C}$, $\mu$ the density of $d\mathcal{C}$, $\vec{v}$ its velocity with respect to G. We could call $\vec{\omega}$ the mean instantaneous angular velocity.

Suppose we assume a reference system fixed with respect to an observatory O and with the fixed stars. We should write the equation

(3)
$$\frac{d}{dt}\vec{S}_o + \vec{v}_o \times M\vec{v}_G = \vec{M}_o^{(e)}$$

where $\vec{S}_o$ is the angular momentum of the earth, with respect to O, $\vec{v}_o$ and $\vec{v}_G$ are the velocity of the O and G with respect to a Newtonian reference system, M is earth's mass.

Concerning $\vec{M}_o^{(e)}$ or $\vec{M}_G^{(e)}$ we have to face the problem of choo-

G. Colombo

sing between what is called the "torque approach" or the "momentum approach".

If in the system earth we include (momentum approach) ocean and atmosphe-re (up to 200 Km above sea level) we have complications in computing $\vec{S}_G$ or $\vec{S}_o$ , (and $\vec{v}_o$ and $\vec{v}_G$ if we use equation 3) because the motion of the oceans and the atmosphere have no negligible components. Beside there are other oompli-cations arising from the poor knowledge of the core's motion.

The only important external torque is in this case the gravitational torque because the effect of other external torques of electromagnetic origin, or inte-raction with interplanetary matter is found to be negligible in the relatively short time scale (say $10^2$ , $10^3$ years) we are considering here .

We would meet other and perhaps more severe difficulties if we confined the earth to the solid crust boundary (torque approach) since besides the external gravitational torque we should have to compute the contact forces experienced by the ocean and atmosphere on the solid earth.

. . . .    . . . . . . .    . . . .

# 1 - Liouville equation

Since we have to begin a theory in some way, we start as usual assu-ming a reference system $\Sigma^*$ centered in the center of mass G , the $x_3$ axis on an axis  reasonably  well aligned with the axis of figure, or with the axis of maximum momentum of inertia, which fortunately is close to the axis of instantaneous rotation, as defined above, and the $x_1$ axis in the plane through $x_3$ containing, say, the Greenwich Observatory.

Since it is reasonable to think that the $x_3$ axis is quite far from Greenwhich the reference system is,  in  this way, defined. The local geo-graphical systems,  the principal axes of inertia, are all moving with respect to this conventional reference system which is also moving.

We proceed now to give an explicit form to $\vec{S}_G$ . Let us call

$\vec{i}, \vec{j}, \vec{k}$ , the unit vectors along the axes of $\sum^*$ , $\omega_1$, $\omega_2$, $\omega_3$ the component of the angular velocity $\vec{\omega}$ of $\sum^*$ with respect to a newtonian reference system , A, B, C, D, E, F the coefficient of the inertia tensor or the earth, $h_1$, $h_2$, $h_3$ the components of the angular momentum due to the velo-city of the earth's components (oceans, atmosphere, core and even mantle ) with respect to $\sum^*$. Then we have

$$\vec{S}_G = ( A \omega_1 - E \omega_3 - F \omega_2) \vec{i} + (B \omega_2 - F \omega_1 - D \omega_3) \vec{j} +$$

(4)

$$+ (C \omega_3 - D \omega_2 - E \omega_1) \vec{k} + h_1 \vec{i} + h_2 \vec{j} + h_3 \vec{k} .$$

Here A, B, ..., $\omega_i$, $h_i$ are functions of time, and naturally $\vec{i}, \vec{j}, \vec{k}$ are also moving.

We note that the following assumptions are generally accepted on the basis of observations and of the choice of the reference system

All the following functions of t

(5)     $$\frac{A(t) - A_o}{A_o} , \quad \frac{B(t) - B_o}{B_o} , \quad \frac{C(t) - C_o}{C_o} , \quad \frac{\omega_1}{\omega_3} , \frac{\omega_2}{\omega_3} , \quad \frac{h_i}{C_o \omega_3} ,$$

are first order infinitesimal where $A_o$, $C_o$ are constant values.

Since in the equation 1) the first derivative $\frac{t}{dt}$ is supposed to be made with respect to a newtonian reference system the equation 1) has to be written in the form

(6)     $$\frac{d'\vec{S}_G}{dt} + \vec{\omega} \times \vec{S}_G = \vec{M}_G^{(e)}$$

where d' denotes differentiation made with restect to the rotating refe-rence system. With the usual notation and neglecting second order quanti-ties, we have the following system

$$A_o \omega_1 + (C_o - A_o) \omega_2 \omega_3 - E \omega_3 + D \omega_3^2 = - h_1 + \omega_3 h_2 + M_1 \ ,$$

(7)
$$A_o \omega_2 + (A_o - C_o) \omega_1 \omega_3 - D \omega_3 - E \omega_3^2 = - h_2 - \omega_3 h_1 + M_2 \ ,$$

$$C \omega_3 + C \omega_3 = M_3 \ .$$

These are the Liouville equations. We will start from them in order to study

the motion of $\vec{\omega}$ with respect to $\Sigma^*$.

## 2 - The Chandler Wobble

We will consider in this paragrah only short periodic variations of
$\vec{\omega}$ . In an interval of time in which secular variation of $\vec{\omega}$ as well of
A, B, C, ...... are negligible we can write , neglecting second order quantities

$$A_o \dot{\omega}_1 + (C_o - A_o) \omega_2 \omega_o - E \omega_o + D \omega_o^2 = - \dot{h}_1 + \omega_o h_2 + M_1$$

(8)
$$A_o \dot{\omega}_2 + (A_o - C_o) \omega_1 \omega_o - D \omega_o - E \omega_o^2 = - \dot{h}_2 - \omega_o h_1 + M_2$$

$$C_1 \dot{\omega}_o + C_o \dot{\omega}_3^1 = M_3$$

where $A_o$ , $C_o, \omega_o$ are the constant mean values of A, C, $\omega_3$ and $C_1/C_o$
$\omega_3^1/\omega_o$ are first order infinitesimal functions of time.

If $E, D, h_1, h_2, M_1, M_2$, in the relatively short interval of time we
are considering here can be supposed to depend not on the position of $\Sigma^*$ but
only on time and on $\vec{\omega}$ and moreover can be supposed to be a linear function of
the component of $\vec{\omega}$ with possibly time dependent coefficients, the system
(8) is a linear system and the general solution is the combination of a
free and a forced term.

Actually $M_1$, $M_2$, $M_3$ do depend on the position of $\Sigma^*$ as the
gravitational torque does, but in a relatively short time interval, since

G. Colombo

the motion of the symmetry axis is slow and it is close to $x_3$ we may consider $\vec{M}_G^{(e)}$ as a function only of time; the same can be generally said for $h_1, h_2$ which are mostly due to ocean and atmosphere meteorological and tidal motions, to sun and moon tides on the solid crust, and finally to the problematical motions of the core. That $h_1, h_2$ depend on $\omega_1, \omega_2, \omega_3$ at least for the contribution of the core's motion, seems to us out of doubt but the problem expressing this dependence, is so severe that we are forced to neglect it at least for the moment, though we are not convinced this has a negligible effect on the motion, even in this first approximation model.

It is certainly true that E and D do depend on $\omega_1, \omega_2$, since the earth is a deformable body and only a uniform rotational motion about an axis fixed in $\Sigma^*$, and therefore in space, can lead to a constant configuration of the inertia ellipsoid with respect to $\Sigma^*$. This axis must be a principal axis of inertia, for the motion to be dynamically possible, when the external torques are zero. Besides for the rotational motion to be stable the axis must be an axis of maximum momentum of inertia, in the case of a symmetrical body, or in presence of internal dissipation, or both. When the angular velocity changes the centrifugal potential also changes. It is assumed that the changing in the axis of rotation is so slow, that the earth has time to reach the static configuration relative to the centrifugal potential at any instant, before the centrifugal potential significantly changes.

The centrifugal potential, neglecting second order terms, may be written

(9)
$$U = \frac{1}{2} \omega_o^2 (x_1^2 + x_2^2) - \omega_o x_3 (\omega_1 x_1 + \omega_2 x_2) .$$

The first term leads to a small constant contribution to the earth's flattening. The second term is the variable term and is the only one we shall consider in the following.

G. Colombo

The corresponding deformation gives rise to an exterior gravitational potential. This, by definition of the Love number $k$, is

$$(10) \qquad V = - k \frac{R^5}{r^5} \omega_o x_3 (\omega_1 x_1 + \omega_2 x_2)$$

where $R$ is the earth's radius and $r^2 = x_1^2 + x_2^2 + x_3^2$ .

But the corresponding terms (in $x_1 x_3$ and $x_2 x_3$) in the expansion of a gravitational potential departing slightly from spherical symmetry are. (from Mac Cullagh's formula) .

$$(11) \qquad \frac{3G}{r^5} (E x_1 x_3 + D x_2 x_3)$$

G being the gravitational constant .

By comparison of (10) with (11) we have

$$(12) \qquad E = - \frac{k \omega_o \omega_1 R^5}{3G} \quad , \qquad D = - \frac{k \omega_o \omega_2 R^5}{3G}$$

D and E being the measures of the distortion of the earth in the yx plane and xz plane respectively, caused by the changing of $\vec{\omega}$ inside the body.

Before proceeding with the Chandler Wobble study, an observation should be made, concerning the Love number k. We assume k to be the tidal effective Love number. The Love number is in some way a measure of earth's yielding to a centrifugal potential or to any perturbing potential like the Sun's or Moon's differential gravity field. From the earth's tide and from the Chandler Wobble (assuming the present theory is valid) one obtains $k \simeq 0,29$. The same value is obtained from the free oscillations of the earth as an elastic body . The evaluation by different authors differ by 1 part in 20. This means a good agreement which may be misleading. Apart from the significant contrast between the value obtained from the figure of earth (k = 0,96) and the value k = 0.30

which may be explained in several ways (see MunKand Mac Donald pag. 27) there can be no doubt that the response of the earth to perturbing forces is not proportional to the amplitude and is not independent of the frequency of perturbations, as we have assumed , in order to write equation (12) .

Now we proceed to substitute expression (12) in the system (8) . In order to find the free component of $\omega_1, \omega_2$ , we consider the following associated homogeneous linear differential system

(13)
$$A_o \omega_1 + (C_o - A_o) \omega_2 \omega_o + \frac{k R^5}{3G} \omega_o^2 \omega_1 - \frac{k R^5}{E G} \omega_o^3 \omega_2 = 0$$

$$A_o \omega_2 + (A_o - C_o) \omega_1 \omega_o + \frac{k R^5}{3G} \omega_o^2 \omega_2 + \frac{k R^5}{3G} \omega_o^3 \omega_1 = 0$$

Solving the system we have

(14) $$\omega_1 = A \cos \Omega(t - \tau) , \qquad \omega_2 = Q \sin \Omega(t - \tau) ,$$

where Q is an arbitrary constant (amplitude and phase of the equatorial component of $\vec{\omega}$) and $\Omega$ has the following expression

(15) $$\Omega = \omega_o \frac{3G(C_o - A_o) - k \omega_o^2 R^5}{3G A_o + k \omega_o^2 R^5}$$

The component of the constant angular momentum in this free motion are

(16) $$\vec{S}_G = (A_o \omega_1 - E \omega_o) \vec{i} + (A_o \omega_2 - D \omega_o) \vec{j} + C_o \omega_o \vec{k}$$

and since E and D are proportional to $\omega_1$ and $\omega_2$ the vector $\vec{S}_G$ is complanar to $\vec{\omega}$ and $\vec{k}$, as it can be immediately shown .

If we assume $k = 0,29$, $A_o = 8,089 \times 10^{44}$ (c.g.s.), $\frac{C_o - A_o}{A_o} =$

$= \frac{1}{304,8}$ , $\omega_o = 7,2921 \times 10^{-5}$ sec$^{-1}$ we find as a period of the free wobble $T_c$ , the value:

G. Colombo

(17)
$$T_c = \frac{2\pi}{\Omega} = 440 \text{ syderal days}$$

The motion is a regular precessional motion with the axis of precession along $\vec{S}_G$ which is constant in the fixed space. In fig. 1) the polhode and herpolhode are drawn. The angular opening of the polhode, which is described by $\vec{\omega}$ in the fixed space is $0.001''$. It is very important at this point to compare the two angular openings because as we shall see later for the precessional motion forced by the gravitational torques of the moon and sun we have the opposite situation the vector $\omega$ hardly moves in $\Sigma^*$ but the motion in fixed space is quite large. The polhode opening is much smaller than the opening of the herpolhodie.

In some sense this is the main reason why the Chandler wobble can be reasonably well separated from all forced motions.

.... ........ ....

3 - Precession and nutation
─────────────────────────

Before discussing the forced component of the motion of the earth's axis of rotation with respect to the solid earth and the damping and excitation mechanism of this motion, it is advisable to recall the main feature of the motion of this axis in the fixed space. In other words if the forcing terms $M_i$, in the equation (8), are put equal to zero the total angular momentum of the earth is fixed in space, since no external torques are present. Actually there are significant external torques more precisely the gravitational torque due to the differential gravi-

G. Colombo

ty field of the moon and sum acting on a non-spherical earth. The angular momentum vector moves consequently in the fixed space. We are now going to compute this forced motion and evaluate the component $\omega_1$, $\omega_2$ in order to see how large is the corresponding motion of the earth's axis of rotation in the solid earth.

Considering the earth as a rigid body with an axially symmetric inertial ellipsoid with axis $\vec{k}$, the torque exerted by the sun on the obla-the earth is

(18)
$$\vec{M} = \frac{3\ G\ M_o}{E^5}\ (C - A9\ (y\ z\vec{c}_1 + z\ x\vec{c}_2))$$

where $M_o$ is the mass of the Sun, $E$ is the earth-sun distance, $C$, $A$ are respectively the maximum and the equatorial moment of inertia, $x, y, z$ are the coordinates of the sun with respect ot the reference system $c_1, c_2, k$, where $c_1$ is oriented in the direction of the vernal equinox.

If $\varepsilon$ is the obliquity, that he actual inclination of the earth's axis on the ecliptic , we have.

$$\vec{M} = \vec{M}_s + \vec{M}_r = \frac{3\ G\ M_o}{2\ E^3}(C - A)\quad \sin\varepsilon\cos\varepsilon\ \vec{c}_1 -$$

(19)
$$- \sin\varepsilon\cos\varepsilon\cos\ 2\ \Omega\ t.\vec{c}_1 + \sin 2\ \Omega\ t\ .\ \vec{c}_2$$

Here $\Omega\ t\ =\ \dfrac{2\ \pi}{T}t$ is the mean longitude of the sun measured from the vernal equinox, $T = 1$ for the tropical year, $\Omega$ is the mean motion of the earth. $\vec{M}_s$ is the secular part of the torque and $\vec{M}_p$ is a periodic

G. Colombo

term with a semiannual period. Other small effects like eccentricity of the earth are neglected.

It is now easier to consider the equation :

$$(20) \qquad \frac{d\vec{S}_G}{dt} = \vec{M}_s + \vec{M}_p = \frac{3\,G\,M_o}{2\,E^3}(C-A)\,(\vec{k}\cdot\vec{n})\,(\vec{k}x\vec{n}) + \vec{M}_p \quad ;$$

Here $\vec{n}$ is the unit vector normal to the ecliptic, plane. If we write $\vec{S}_G = \vec{k}\,C\,\omega_o$ we have

$$(21) \qquad \frac{d\vec{k}}{dt} = \frac{3(C-A)\Omega^2}{2C\omega_o(1-e^2)^{3/2}}(\vec{k}\cdot\vec{n})\,(\vec{k}x\vec{n}) + \frac{\vec{M}_p}{C}$$

The secular part corresponds to a rotation of the vector $\vec{k}$ about n in the clockwise direction, the corresponding displacement of $\vec{c}_1$ is $15.92''/$ year.

Superimposed to this motion there is a nutation with a semiannual period and an amplitude of $0.5''$.

If we follow the same procedure to find out the effect of the moon's gravitational torque, neglecting a small periodic part, which corresponds to a nutation with a period of $\frac{1}{2}$ month , we have a secular displacement given by

$$(22) \qquad \frac{d\vec{k}}{dt} = \frac{3}{2}\frac{(C-A)\,n^2}{C\omega_o(1-e^2)^{1/2}}(\vec{k}\cdot\vec{N})\,(\vec{k}x\vec{N})$$

where n' is the mean motion of the moon, $\vec{N}$ is the unit vector normal to the moon orbital plane.

Since $\vec{N}$ is rotating about $\vec{n}$ in 18.6 years, it may be decomposed into a fixed component along $\vec{n}$ and magnitude cos i where i is roughly $5^{\circ}$ (inclination on the eccliptic of the moon orbital plane) and a rotating component with a period of 18.6 years.

The angular momentum vector is again rotating about the $\vec{n}$ axis in the same sense as for the sun-torque . The resulting rotation about the $\vec{n}$ axis of the vector $\vec{S}_G$, on a cone of angular opening brings the displacement of the vernal equinox to 50.37" / year (the period is 25,730 years) .

Besides there is a nutational motion of $\vec{S}_G$ with period of 18.6 years on an elliptic cone of semiamplitude 6.87" in the direction tangent to the secular motion of $\vec{S}_G$ , and 9.21" in the normal direction.

We may now evalutate the equatorial components of the angular velocity $\omega_1$, $\omega_2$ which correspond to this motion.

We observe only that the ratio $\omega_e/\omega_o = \sqrt{\omega_1^2 + \omega_2^2}\ \omega_o$ is of the order of the ratio of the two angular velocities 50" sin $_o$/ year and $\omega_o$ = = $360^{\circ}$/ day , which means that the angular velocity vector forms with the $\vec{k}$ axis and with the $\vec{S}_G$ vector an angle smaller than 0.02" .

The main precessional motion may be obtained by rolling the polhode which is a cone of semiamplitude of roughly 0.02" sec rigidly connected with the earth and centered in the axis of figure on the herpolhode which is a cone of roughly $23^{\circ}$ semiamplitude, centered in the pole of the ecliptic.

The rolling period is roughly one day. The nutations correspond to a negligible wiggling of the two Poinsot conical surfaces.

For this reason the nutation of the axis or rotation inside the earth due to the external torque averaged in 1 day is small with

respect to the motion of the pole due to the Chandler wobble and in any case it is precisely determined.

.... ..... .....

## 4 - Forced component and excitation and damping of Chandler wobble.

Considering again equations (8) the first two terms on the right side of these equations, give, in the model we have considered in 2) , the forcing terms due to motion of masses, ocean, atmosphere, with respect to the main body which behaves elastically. The main term is an annual term due mainly to the annual displacement of air masses.

From reduction of the latitude observations data, (I.L.S.) the orientation of the instantaneous axis of rotation with respect to the stations are derived, and consequently with respect to the solid earth.

The observed displacement $m_1 = \dfrac{\omega_1}{\omega_0}$ ; $m_2 = \dfrac{\omega_2}{\omega_0}$ are represented for the interval of time 1900 in Fig. 7.4 taken from M.M.D.

Also from M.M.D. is taken the fig. 7.5. showing the power spectrum of variation of latitude obtained through an analysis of 54.4 year records. From this data the following results are obtained (reported also from M.M.D. book).

(1)  98.5 % of the power is associated with positive (west-to-east motion)

(2)  93 % of the power is contained in the frequency range 0.74 to 1.14 cycles per years.

(3)  22 % of the power in this range is in an annual line without recognizable structure (in the analysis it happens to fall between harmonics 54 and 55) .

(4)    78 % of the power in this range   is contained in the Chandler peak
which  is centered near   0.85 c/year and has a noticeable  band-struc-
ture.

From   observations at Greenwich and Washington the annual term
can be evaluated through   different procedures followed by different authors,
using smoothed or unsmoothed values,  give different results. It   seems qui-
the well established that  the annual ellipse has the semimajor   axis   of
the order   $0.092''$ and  it seems to  be oriented quite well  with the Green-
wich  meridian . Whether  this is an observational effect   or  not it is
questionable. The other semimajor axis   has   amplitude $0.075''$ . If  the
term $h_1$, $h_2$   exciting the annual component  could be found, we would have
a good check of the validity of the model represented by system (8).


Unfortunately this is   not   the case.

Satisfactory interpretation   of   the non seasonal residue correspon-
ding to the Chandler wobble has   not   yet   been achieved. The main   reason
is apparently due to   the   lack   of good observation for a long   period of
time.

The problem   is mainly concerned with   the mechanism of supplying
the energy to maintain the natural mode (Chandler motion) of   14   months
period   in   a body   which is dissipative in nature.

Two types of model are considered in the literature. The first
one is the   "damped model" . The finite spectral width   may   be associated
with   the model   of   a linear dissipative oscillator excited at random. The
damping   is due to   imperfections   of elasticity or dissipation at   the boun-
dary between mantel and core,   if the core is viscous . Irregular variation
of   the   atmosphere is believed to be the most   likely cause  of the needed
random excitation.

The second model   is   a time-variable model , that   is a linear

oscillator with time dependent characteristics. This time dependence of the parameters characterizing the model implies a time dependence of the physical parameters (like rigidity, ellipticity, etc) . Physical considerations speak against a time variable model.

Accepting the linear damped model excited at random a problem arises for the excitation due to irregular atmospheric variations.

An evaluation, from the latitude observation data, of the specific dissipation function

$$(23) \qquad Q = \frac{1}{2\pi E} \oint \frac{dE}{dt} \; dt$$

gives a value of $Q$ of the order of 30 to 40. If we consider the case of a linear damped oscillator

$$(24) \qquad \ddot{x} + 2\alpha \; \dot{x} + \sigma_o^2 \; x = 0 \; ,$$

we have

$$(25) \qquad Q = \frac{\sigma_o}{2\alpha} \quad ,$$

and since $\dfrac{1}{\alpha} = \dfrac{2}{\sigma_o} = \dfrac{2 Q T}{2\pi}$ a value of 30 to 40 correspond to a relaxation or damping time of 10 to 13 periods roughly. Q is also related to the amplitude at resonance and to the sharpness or resonance.

This value of Q requires atmospheric variations somewhat larger than one would like. Irregular motions of the core although not excluded seem to be unlikely.

A larger value of Q (between 100 and 200) is obtained from pole tide observation ; this value would imply a larger relaxation time and would require a more acceptable value for the irregular atmospheric variation needed for exciting the Chandler wobble.

On the other hand an analysis of the possible values of Q bring to

G. Colombo

the following conclusions taken also for M. M. D. (pag. 172).

1) If the  Q  is betwenn 100 and 200, as vaguely suggested by the poletide
   analysis  then solid friction in  the mantle can account for the damping.

     For  a  Q  of  30 to 50, as indicated by the latitude analysis, the-
   re  are many possibilities :

2) The damping can be in the oceans

3) The lower mantle is a possible energy sink; a model involving a Maxwell
   viscosity, is, however, unsupported by laboratory and seismic evidence.

4) Damping by viscosity in the core appears to be ruled out; electromagne-
   tic damping is still a possibility (Jeffreys, 1956).

5) Impulses of a non-random kind (originating in the core, oceans or  atmo-
   sphere) can absorb as well  as  excite the wobble. The computed  Q
   is then not due to damping, but associated with the interaction between
   these loosely-coupled components;

  I would like to emphasize again the weakness  of the damped model
randomly excited in explaining the observed motion of the pole.

  Apart from the peak period of the Chandler wobble band in the
power spectrum, neither the damping nor the energy supply is clear.

  Even accepting as  a first conclusion that the statistical propery
of the latitude time series, up to now available, are those associated with
a damped oscillator excited at random it seems to me that some other
possibility has  to  be carefully examined in the near future when better
and more numerous observations will be available. For instance the possi-
bility that a non-linear excitation mechanism is responsible of the energy
transfer from both high and low frequency periodic excitation (from
diurnal to annual) to the Chandler free wobble.

  The real model both with regard to the elastic property and the
internal dissipation of energy is certainly not linear. Firction internal forces,

G. Colombo

besides the viscous  one , are certainly present. On the other hand excitation of proper mode of low frequency by an external periodic force with high frequency is not an exception in non-linear systems. Period and amplitude of the Chandler wobble seems to be proportional. The correlation coefficient found by Nicolini  is +0.88. The change in period with the amplitude is characteristic in non-linear oscillators.

Before a definite conclusion can be drawn more refined observations for a longer period  of time are badly needed. It is quite  surprising that a regular behaviour of the    $m_2$ component after removal of the annual term, as shown in the second curve of  fig. 7.4, can be only explained as a randomly excited oscillator of a linear one degree of freedom system.

The interpretation of the curve becomes more surprising after reading the last  sentence of page 3 of the book by Munk and Mc Donald, first paragraph:  "The wobble is generated by random impulses of unknown origin, and damped  by some unknown imperfections from elasticity or by some other means."
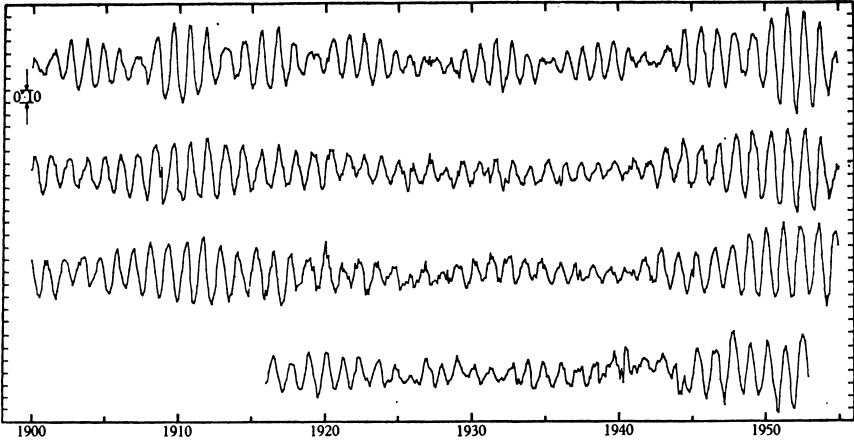
Fig. 7.4. The component, $m_1$, of the unsmoothed ILS observations, before (*top*) and after (second curve) removal of the seasonal variation; the component, $- m_2$, of the unsmoothed ILS observation after removal of the seasonal variation (third curve) and the corresponding non-seasonal variation in the latitude of Washington, as obtained with the PZT (*bottom*).

OBSERVATIONS OF LATITUDE



Fig. 7.5. The spectrum of variation in latitude, according to Rudnick (1956). The upper figure refers to the positive (west-to-east) motion, the lower figure to the negative motion of the pole of rotation (see § 6.7). For both motions the harmonics 40 to 62 are shown, with the corresponding frequency scale, in cycles per year, indicated below. The length of the spectral lines gives the contribution per harmonic toward the mean square radius arm (in units of (0″01)²). The scale for negative motion is ten times that for positive motion. The curve has been fitted by the method of maximum likelihood (see Appendix A.2).

CENTRO INTERNAZIONALE MATEMATICO ESTIVO

(C.I.M.E.)

E.M. GAPOSCHKIN

THE MOTION OF THE POLE AND THE EARTH'S ELASTICITY AS
STUDIED FROM THE GRAVITY FIELD OF THE EARTH BY
MEANS OF ARTIFICIAL EARTH SATELLITES

# THE MOTION OF THE POLE AND THE EARTH'S ELASTICITY AS STUDIED FROM THE GRAVITY FIELD OF THE EARTH BY MEANS OF ARTIFICIAL EARTH SATELLITES[1]

E. M. Gaposchkin (Cambridge, Mass)

## 1. INTRODUCTION

This conference has been entitled "Modern Questions of Celestial Mechanics." Although the emphasis has been on the problems of celestial mechanics that need to be solved, the orientation of my talk is more toward the problem to which celestial mechanics will provide an answer.

In general, the celestial-mechanics problem can be characterized by a statement about the forces acting upon bodies, either planets or earth satellites, and the celestial mechanician attempts to solve the differential equations of motion. If one is really conscientious, then observations of the system being studied are taken and the theory is tested against these observations. Artificial earth satellites have posed the reverse problem, that of finding the forces acting on the satellite.

For several years now at the Smithsonian Astrophysical Observatory (SAO), this investigation has been carried on. Needless to say, a formal celestial-mechanics theory has to be developed in order to determine the forces. Within the last year, an extensive documentation on the final results of SAO has been published, and is available. However, I am not going to discuss that here in any detail. Implied in any set of geodetic parameters defining the earth as a reference system is that the earth is a rigid body with a pole, or an equator, and with an arbitrary meridian as a reference point. This is what we would call a terrestrial, or earth-fixed, system. If we ignore such questions as continental and crustal motions, the fixed points on the terrestrial system are constants. The earth is not a rigid body, however, and undergoes deformation due to the attraction of the sun and the

moon. This tidal deformation was studied by Lord Kelvin, George Darwin, and Harold Jeffreys, among many investigators. The largest tidal effect is of some 60 cm.

Newtonian mechanics is formulated in a so-called inertial system. The terrestrial system is not inertial, and the equations of motion have extra terms to account for this. The sidereal system — that system defined by the stars or the far-distant galaxies — is the most pleasing definition of an inertial reference frame. In any case, the earth has a motion with respect to this special reference frame. It is made up of three kinds of motion, that of precession and nutation, that of the variation in the length of the day (in other words, the irregular changes in the rate of rotation of the earth), and that of the Eulerian nutation or wobble of the earth.

This separation is, to a large extent, natural, since there are essentially different observations used to measure these quantities. Also, the physics involved is to some extent different. The precession and nutation, for instance, are due to the gravitational attraction of the sun and the moon on the earth's equatorial bulge. The variation in the length of the day has one component that is called the secular deceleration. Modern geophysical thinking describes this as the loss of energy in tidal friction. The motion of the pole is a free nutation. The currently accepted geophysical thinking on this is that the period of the free nutation is governed by the elasticity of the earth. It is this latter subject that I propose to discuss.

## 2. DESCRIPTION OF THE PHENOMENON

For many years the rotation of the earth was the best clock that the astronomer had. It became obvious only this century that the earth was not a good timekeeper. For instance, many problems in lunar theory were ascribed to the lack of mathematical rigor. People used to talk about the great empirical term in lunar theory. Only in the past several decades has it become known that the problems with observations of the longitude of the moon were not problems in the theory, but problems in the earth's rotation. Similarly, the free nutation of the earth was a matter of some interest. It was Euler who showed in 1765 that a rigid body could have a stable rotation

about the principal axis of inertia. He also s..owed that a rigid body like the
earth could have a free nutation with a period of A/(C-A) sidereal days, where
A and C are the principal moments of inertia. There were many unsuccess-
ful attempts to observe this 10-month period. It was in 1891 that Chandler,
a prosperous merchant in Cambridge, Massachusetts, discovered the motion
that bears his name. The period of this motion is 14 months. Newcomb
immediately showed that the elastic yielding of the earth would lengthen this
period. Chandler's discovery of this motion led to the establishment of the
International Latitude Service, and it is from about 1900 that measurements
of the position of the pole have been made.

   The Latitude Service essentially measured the angle between the vertical
and the pole by using known reference stars. The pole, defined by preces-
sion and nutation, is then defined as the instantaneous pole. This pole
becomes the spin axis, which is 0.003 arcsec different from the angular
momentum and is observationally indistinguishable. Now, the crust of the
earth moves with respect to this in a classical Poinsot construction, so that
the latitude of any station will change; hence the name, variation of latitude.
The Latitude Service provides plots of the x and y displacement with respect
to an arbitrary reference. The magnitude of this displacement is some 10 m.
The first slide gives the record for 1958 to the present. We notice immediately,
that there is a roughly irregular motion; the maximum excursion is about
6 arc sec, which is approximately equivalent to 20 m on the pole. Immediately,
we come upon one subject of controversy today: whether the pole has a
secular motion or not. I do not plan to address myself to this problem, but
in the next slide I will show the same data plotted with respect to a mean pole
of the year. We see essentially the same size motion, but the motion is
significantly more regular. The latter data are published by the Bureau
International de l'Heure in Paris. The former plot is from the International
Polar Motion Service. Some people have tried to use these data to demon-
strate the fact of the secular variation of the pole. Others have stated that
we can attribute this apparent secular motion to the change in latitude of one
station, the candidate being the Japanese station, which is in a geophysically
very active area, having large tectonic activity, earthquakes, and crustal
motions. We just do not know what to make of this situation.

A few comments on the data are probably appropriate. The number of latitude stations throughout the whole period has been small, between five and seven, with only three participating during the entire interval. They have undergone different administrations. It is also noted that the irregular changes in the apparent secular motion of the pole are correlated with the change of star catalogs. Many of the problems in treating the astronomical data result from these irregularities. This is an area in which I think much careful work remains to be done.

The accuracy of the data can be estimated as follows: The measurement of zenith distance can be made accurate to perhaps 0.1 arcsec. If we make a thousand observations, the reliability reduces to less than a hundredth of an arcsec. This figure of a hundredth of an arcsec is equivalent to 1 foot at the pole. The thousand observations are made over a period of 15 to 20 days, which means that any very short-period variations are averaged out of the data, and we would detect no correlation with short-lived geophysical events, such as daily meteorological changes.

The next slide shows the data of the x component from 1900 to 1960. One of the major questions concerning the Chandler wobble is the phenomenon that maintains it, which has not yet been discovered. Since the earth is elastic and has imperfections, it would be dissipative and the free nutation would have been expected to have died out a long time ago. This is the Chandler wobble, with a 14-month period. Attempts to explain this in terms of atmospheric phenomenon have been unsuccessful (Munk and Hassen, 1961), and as I said, there do not seem to be any suitable candidates for exciting this 14-month period. The mechanism for maintaining the motion is not understood at all today.

Going on quickly to some other properties of this motion, we show a power-spectrum analysis of this latitude variation. We immediately notice a very narrow spike at 1 year. This leads us to look for driving effects of a 1-year period, and this investigation has been moderately successful. An early work in this area was by Harold Jeffreys in the fascinating paper

in 1916. The major effects are the distribution of air masses and the distribution of water in the oceans. You can see a broader peak roughly centered at 0.85 cycles per year.

The next slide shows the annual variation. It is considerably smaller, with a maximum excursion of about 10 feet. This is the determination made by Sir Harold Jeffreys. Another determination by Walker and Young is not in very good agreement with this, although the sizes are not changed significantly. One curious thing about this is the orientation of the semimajor axis along the Greenwich meridian.

## 3. THE DEFORMATION OF THE EARTH

As I mentioned in my introduction, the earth is an elastic body. It therefore undergoes certain deformation. Then, whatever the theory of the constitution of the earth, provided the physical constants such as density and incompressibility are functions of $r$, we can write the effects of a disturbing potential of second degree as follows. The earth surface is lifted by the amount $h(U_2/f)$ and the horizontal displacements will be $(l/f)(\partial U_2/\partial\phi)$ and $[l/(f\cos\phi)](\partial U_2/\partial\lambda)$, where f is the gravitational constant. If we turn our attention to the external gravity field, the deformation of the earth produces the additional external potential:

$$\Delta U = k\left(\frac{a}{r}\right)^5 U_2 \quad ,$$

where a is the radius of the earth. Now the displacements due to rotation around $(l, m, n)$ are the same as those due to the potential of the second degree; and to the first order in the direction cosines $l$, m is

$$\frac{1}{2}\Omega^2(x^2+y^2) - \Omega^2 z (lx + my) \quad .$$

The first term adds to the oblateness, and does not concern us. The second term will give rise to the additional external graviational potential

$$\Delta U = - k\Omega^2 z(\ell x + my)\left(\frac{a}{r}\right)^5 \quad .$$

Now we recall that MacCullagh's formula gives for the gravitational potential due to a deformed earth

$$V = f\left[\frac{M}{r} + \frac{(A+B+C)r^2 - 3(Ax^2+By^2+Cz^2-2Fyz-2Gxz-2Hxy)}{2r^5}\right] \quad .$$

Writing the value of the moment of inertia around the direction to the external point (x, y, z), and comparing terms, we note that the terms in the moment of inertia tensor, corresponding to this elastic deformation, are

$$F = \frac{-k\Omega^2 m a^5}{3f} \quad ,$$

and

$$G = \frac{-k\Omega^2 \ell a^5}{3f} \quad .$$

The products of x and y, and x, y, and z can be written in terms of spherical harmonics of second degree, and the resulting expression for the gravitational potential due to the deformed earth is

$$V = \frac{fM}{r}\left\{1 + \left(\frac{a}{r}\right)^2\left[\frac{A-C}{fMa^2} P_{20}(\cos\phi) + 2 P_{21}(\cos\phi)\left(\frac{G}{fMa^2}\cos\lambda + \frac{F}{fMa^3}\sin\lambda\right)\right.\right.$$

$$\left.\left. + \frac{H}{fMa^2} P_{22}(\cos\phi)\sin\lambda\right]\right\} \quad ,$$

where $P_{\ell m}$ is the usual associated Legendre function. If we put in some approximate values for the radius of the earth, the rotational rate, the gravitational constant, etc., we get a gravity-field harmonic coefficient of the order of $10^{-8}$:

$$C_{21} = 2\sqrt{\frac{3}{5}}\,\frac{G}{Ma^2} = 1.55\,\frac{-k\Omega^2\,\ell\,a^5}{3f\,Ma^2}$$

where

$$\Omega = 0.727 \times 10^{-4}\text{ rad sec}^{-1},$$

$\ell$ = component of nutation $\approx 1$ sec $= 0.5 \times 10^{-5}$ radius,

a = earth radius $6.378 \times 10^{-8}$ cm,

$$fM = 3.986 \times 10^{20}\text{ cm}^3\text{ sec}^{-2},$$

and

$$C_{21} = -k\,0.1 \times 10^{-7}.$$

## 4. DETERMINATION OF THE GRAVITY FIELD OF THE EARTH

For the past decade, SAO has been observing artificial earth satellites with a camera system and using these observations to determine the earth's gravitational field. We have adopted as our reference system the instantaneous pole for the gravity field, and the terrestrial earth as our reference system for the station locations. We have utilized the empirically observed data of the motion of the pole, which I have shown you on a previous slide, and the empirical data on the rotation of the earth in terms of $UT_1$. We have translated the site locations into this reference system to do our orbital dynamics, and we have then solved the problem for the earth's gravity field.

The assumption we have made is that our coordinate system is oriented along an axis of principal moment of inertia, and that the quantities $C_{21}$ and $S_{21}$ have therefore been zero. We have set them arbitrarily for zero. In the past year, we have completed an 8th-degree, 8th-order solution with some higher order terms in the tesseral harmonics, and Dr. Kozai determined the zonal harmonics up to the 14th degree.

The next slide shows the general form of the gravity field used in the equations of motion. The previous equation for the moments of inertia of the earth up to the second degree is given in the same notation.

Table 1 gives the values of some of the coefficients we have determined and the computed uncertainty. I estimate that the uncertainties here are accurate to a factor of 2. The second-degree terms are computed to about 1 part in $10^8$. I give in this table a very high-order term. This term is determined to about 2 parts in $10^{10}$. We are actively improving our procedures, acquiring new and more accurate data, and will perform a more accurate solution within the next year. It is conceivable that the accuracies will be improved by a factor of 5, especially the low-order harmonics. We also intend to include direct measurements of gravity, in our determination of the gravity field, and our future work should be a significant improvement.

Table 1. Selected coefficients of the geopotential (from Gaposchkin, 1967)

$$C_{22} = 2.379 \times 10^{-6} \pm 0.13 \times 10^{-7}$$

$$S_{22} = -1.351 \times 10^{-6} \pm 0.13 \times 10^{-7}$$

$$C_{31} = 1.936 \times 10^{-6} \pm 0.13 \times 10^{-7}$$

$$S_{31} = 0.266 \times 10^{-6} \pm 0.14 \times 10^{-7}$$

$$C_{15\ 13} = -0.058 \times 10^{-6} \pm 0.17 \times 10^{-8}$$

$$S_{15\ 13} = -0.046 \times 10^{-6} \pm 0.17 \times 10^{-8}$$

$$C_{15\ 14} = 0.0043 \times 10^{-6} \pm 0.21 \times 10^{-9}$$

$$S_{15\ 14} = -0.0211 \times 10^{-6} \pm 0.22 \times 10^{-9}$$

Comparing the size of the uncertainty in the accuracy with which we can compute the low-order harmonics, we see that it is well within our grasp to compute the values of $C_{21}$ and $S_{21}$ using our current techniques.

The next slide shows the development of the perturbations in one of the elements, the mean anomaly, based on this theory. This perturbation is the most intricate, but in general is typical of the mathematical form that results. The essential point is the divisor of the form

$$[(\ell - 2p)\dot{\omega} + (\ell - 2p + q)n + m(\dot{\Omega} - \dot{\theta})]^{1,2}$$

If the mean motion of a satellite is nearly an integral number of days, and if we recall that the motion of perigee and the node is small and the sidereal rate is roughly 1, we get a very small number. This leads to the so-called resonant harmonics. In the satellites that we used for the Standard Earth, we had several resonant harmonics, and I presented an example of the determination of one in Table 1.

Other essential points to note are the kinds of functions involved in the perturbations. The function G involves the well-known Hanson coefficients. It is a function that goes as the eccentricity to the absolute value of q. Therefore, the perturbation is largest for values of $q = 0$. This means that $p = \ell/2$. The functions F are polynomials in sine and cosine inclination. For the case of $\ell$ and $m = 2$ and 1, and $p = 1$, F is proportional to $\sin I \cos^2 I$. We can then suggest that the most useful satellite for determining these $C_{21}$ and $S_{21}$ would be a satellite with a mean motion of 1 revolution per day, at an inclination of 35.4°. With such a satellite, the low-order term should be determined to 1 part in $10^{10}$.

Of course, all the degrees from $\ell = 2$ to infinity will have a term with $m = 1$, and therefore the synchronous satellite will excite perturbations from all these harmonics. This is why the other terms $(\ell, 1)$ must be determined as well as possible, so that they can be treated as known in the determination. We note that a synchronous satellite of 1 revolution per day is at a radius of $\approx 6.6$ earth radii. There is a divisor in all the perturbations of $A^{\ell + 3}$. This means that for a large L, this divisor will damp out the perturbation, if we do not get too close to resonance.

## 5. SUMMARY

If we return to slide 1 now, we can see what the procedure will be. We will take all the orbits and all the observations of the satellites, whether they are synchronous or not, and group them into periods when the pole is at a particular position. This is a slowly changing period, and I imagine we can group the observations into 1-month segments. We would then have essentially a plot of the moments of inertia, as a function of time. Then returning to our equation involving the Love number and the observed position of the reference pole, we could use this relationship and the determinations of F and G to compute a value for k. Since we postulate that we can determine F and G to 1 part in $10^{10}$, since F has a value of $10^{-8}$, we should be able to compute k to three figures. The value of k has been measured from the period of the Eulerian nutation, from the seismic observation, and from the tides to be about 0.27 with the uncertainty in about the second digit. Having these data now, and the time dependence of F and G, we could do other things; however, if we assume k from the values of seismology, we could then determine the values of $\ell$ and m, the coordinates of the pole. We could then use this as a method of measuring the motion of the pole. Presumably, data could be made available quickly enough from electronic systems so that the results of this analysis could be available perhaps within time measured in weeks rather than in months and years, as it is now.

On the other hand, if we return to the fundamental definition of F and G as mass integrals of the earth around this rotation axis, and if we are willing to accept the density data from seismology, we can then perform this integral and find out how much of the earth is participating in this motion. This is of immense interest today, because of the uncertainty in the relation of the core to the mantle, and because of the dynamics of the possible rotation of the core. These last few subjects are primarily of geophysical interest, but they demonstrate the benefit of the space program both in the study of the solid earth and in the study of the subject of celestial mechanics. It demonstrates the seemingly endless questions of geodynamics and celestial mechanics.

## BIBLIOGRAPHY

GAPOSCHKIN, E. M.
 1967. A dynamical solution for the tesseral harmonics of the geopotential
        and station coordinates using Baker-Nunn data. In Proceedings
        VII Intl. COSPAR Space Sci. Symp., North-Holland Publ. Co.
        (to be published).
JEFFREYS, H.
 1962. The Earth. 4th ed., University Press, Cambridge, England,
        438 pp.
LOVE, A. E. H.
 1909. The yielding of the earth to disturbing forces. Proc. Roy. Soc.,
        vol. 82A, pp. 73-88.
MUNK, W. H., AND HASSEN, S. M.
 1961. Atmospheric excitation of the earth's wobble. Geophys. Journ.,
        vol. 4, pp. 339-358.
MUNK, W. H., AND MAC DONALD, G. J. F.
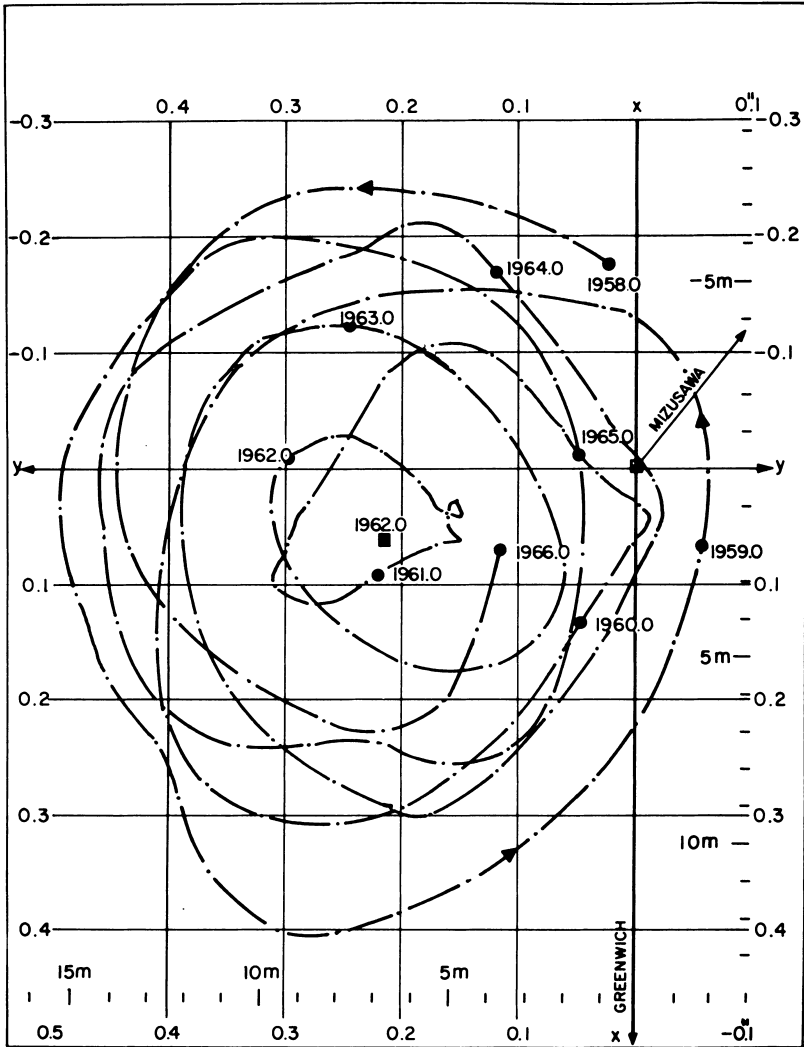 1960. The Rotation of the Earth. University Press, Cambridge,
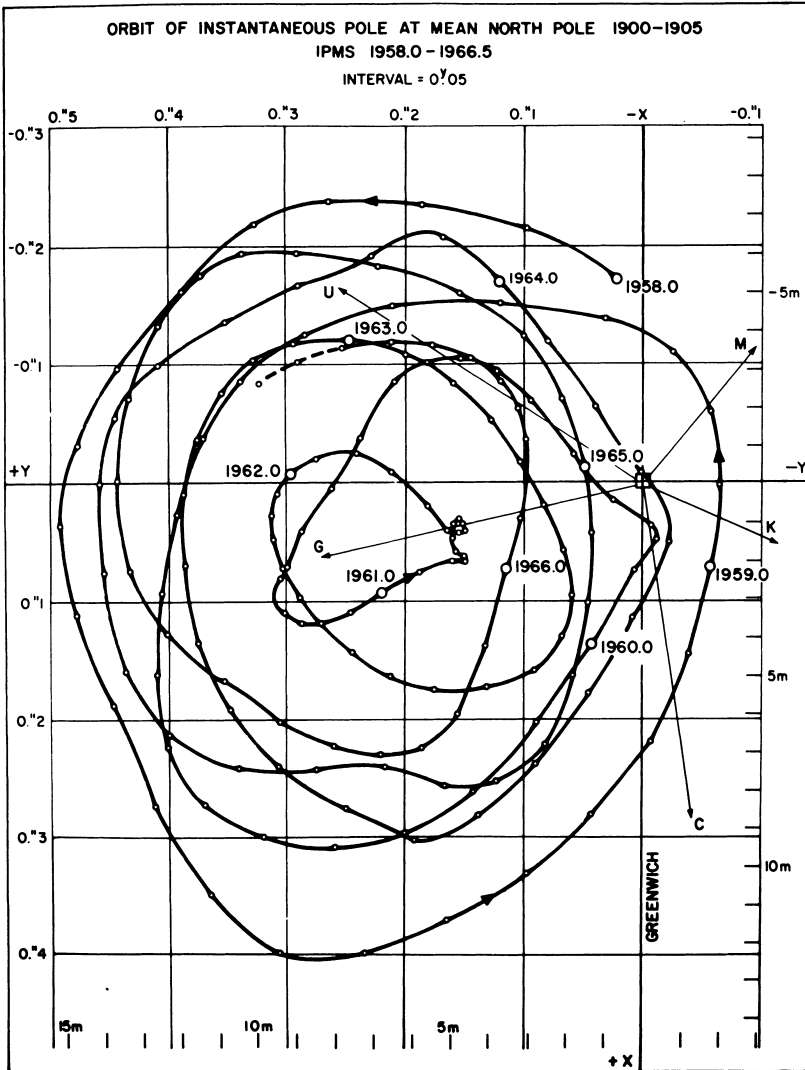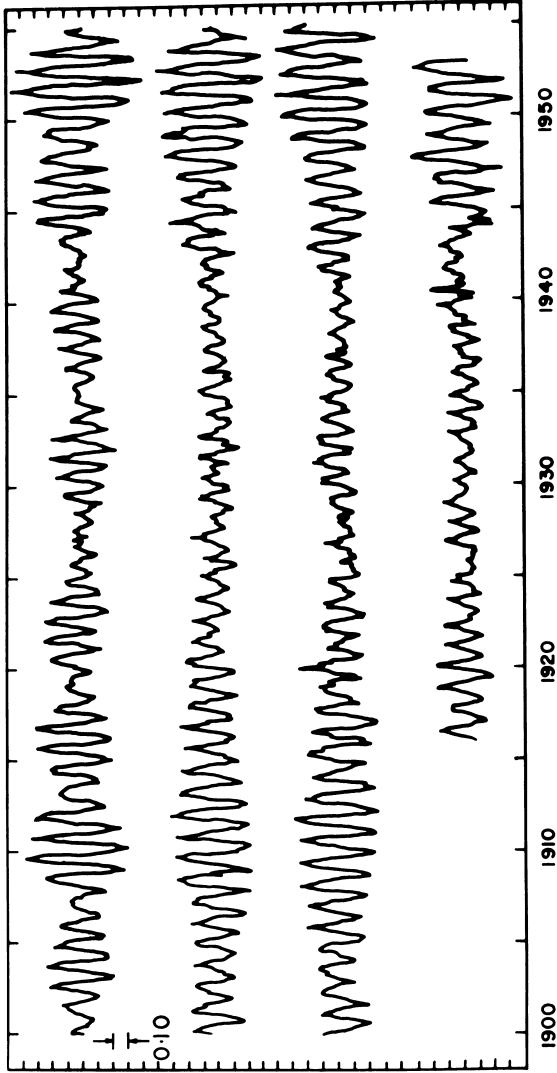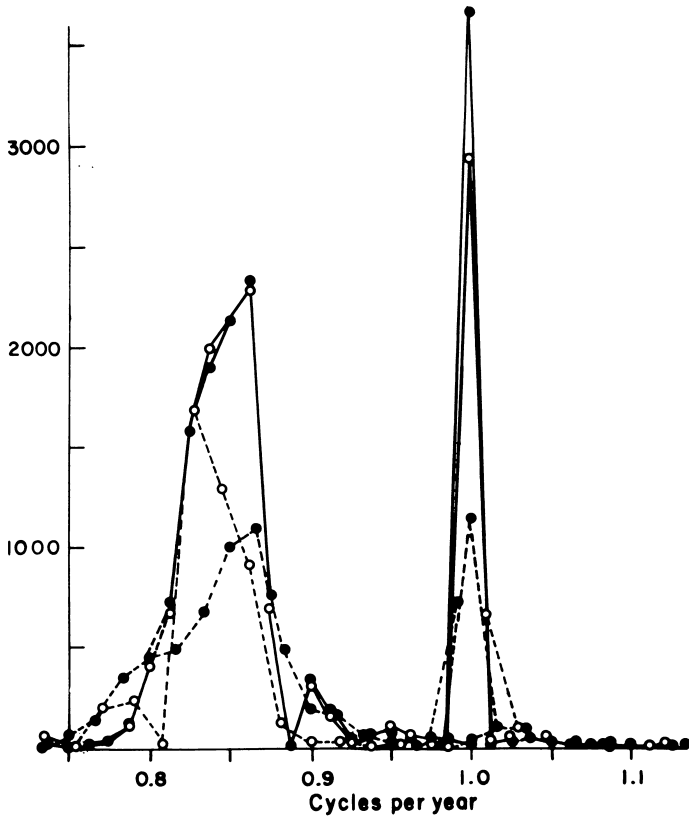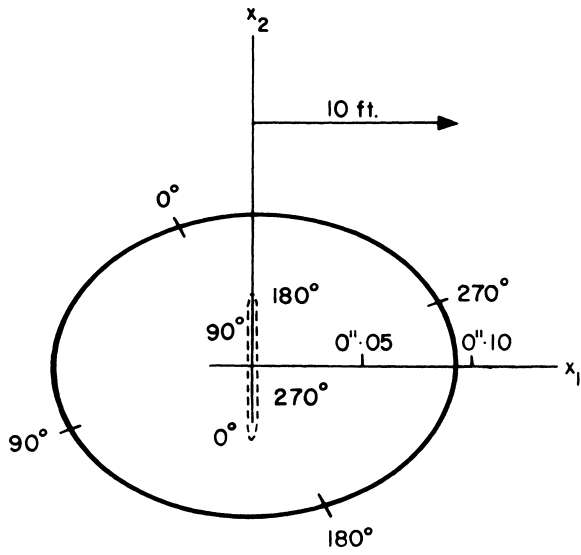        England, 323 pp.

Figure 1.

Figure 2.

Figure 3.

Figure 4.

Figure 5.

CENTRO INTERNAZIONALE MATEMATICO ESTIVO
(C. I. M. E. )

Irwin I. SHAPIRO

RADAR ASTRONOMY, GENERAL RELATIVITY, AND CELESTIAL MECHANICS

TABLE OF CONTENTS

# RADAR ASTRONOMY, GENERAL RELATIVITY, AND CELESTIAL MECHANICS

by

I. I. Shapiro (Cambridge, Mass.)

## I. INTRODUCTION

Radar astronomy is a young science. Yet, like most modern enterprises, it is already large -- certainly too large to be covered adequately in eight lectures. We will therefore concentrate on those aspects which have the most influence on general relativity and celestial mechanics.

The extension of radar measurements to interplanetary distances heralded a major advance in the accuracy achievable in estimating astronomical constants and planetary ephemerides and rotation rates. Nonetheless one is tempted to ask, for example, why it is of interest to determine the n+1 st significant figures in the descriptions of planetary orbits and constants. The answer is twofold: Theoretically, the more accurate the measurements the more stringent the test of the underlying physical theory; from a "practical" point of view, it was necessary, for example, to improve the determination of the astronomical unit of length in terms of a terrestrial unit in order to direct space probes to the planets. We shall amplify both points in subsequent lectures.

Perhaps the most startling of the discoveries made with radar concerns the rotation vectors of Mercury and Venus. Mercury's spin appears to be locked to its orbital motion such that the spin angular velocity is three-halves the orbital mean motion. Venus' spin is retrograde and, moreover, is apparently locked to the Earth with Venus rotating on its axis four times, as seen by an Earth observer, between successive inferior conjunctions. We shall discuss in detail the methods used to determine the rotation rates and the theories developed in an attempt to explain these unexpected spin-orbit resonances.

## II. CHARACTERISTICS OF RADAR

We begin with a brief description of the fundamentals of radar. The word itself is an acronym (radio detection and ranging) as well as a palindrome. The mai . advantages offered by radar stem from the control that can be exercised over the radiation. The experimenter can select (within technological constraints) the frequency, modulation, polarization, intensity, geometrical origin, time origin,

direction of propagation, etc. More formally, we may denote the electric vector of the transmitted wave in complex notation by

$$\vec{E}(t) \;=\; A(t) \exp \left[i(2\pi f_o t + \varphi_o)\right] \hat{p}(t) \quad , \tag{1}$$

where A(t) represents the so-called envelope or modulation of the transmitted wave and is slowly varying with respect to $f_o t$; $f_o$ denotes the carrier frequency and $\varphi_o$ the initial phase of the wave. The direction of polarization is indicated by the unit vector $\hat{p}(t)$. By studying the modifications introduced into the echo, one can estimate the motion of the target, its size, shape, and surface characteristics, as well as the properties of the intervening medium (e.g., the Earth's and planet's atmosphere and the solar corona or interplanetary medium).

The ability to detect on Earth the reflection of a radar wave from a target planet depends first of all on the power of the echo received. This power $P_r$ is expressed in the radar equation as

$$P_r \;=\; P_t \, \frac{G_t}{4\pi R^2} \, \frac{A_r}{4\pi R^2} \, \sigma \tag{2}$$

where

$$A_r \;\equiv\; \frac{\lambda^2}{4\pi} \, G_r \quad . \tag{3}$$

The various symbols are defined as follows: $P_t$ is the transmitted power; $G_t$ and $G_r$ the gains of the transmitting and receiving antennas, respectively; $A_r$ the effective area of the receiving antenna (roughly half the geometric area); R the distance to the target; $\sigma$ its radar cross section; and $\lambda$ the wavelength of the received radiation. This expression is valid only if the beamwidth of the transmitting antenna (in radians, approximately $\lambda/D$, where D is the antenna diameter) is larger than the angular size of the target. The derivation of the radar equation follows essentially by inspection. The received power is obviously directly proportional to the transmitted power and to the cross section of the target, suitably defined. The energy transmitted at each instant spreads in time over the surface of a sphere of ever increasing radius; hence the energy density will decrease with the surface area. The power incident on the target will thus be inversely proportional to $4\pi R^2$ -- the surface area of the sphere -- but modified by the anisotropy, or concentration, of the radiated power in the direction of the target. This latter factor is simply the antenna gain, $G_t$, which is related to the effective aperture,

or area, ∙ ⌐ he transmitting antenna as in Eq. (3). (The ratio of the effective to the geometric area is called the efficiency of the antenna.) Upon reflection, the echo signal power again spreads out; the fraction impinging on the receiving antenna will be proportional to $A_r/4\pi R^2$. The portion of the incident energy on the target that is backscattered towards the antenna is described by the radar cross section $\sigma$ which is normalized so that Eq. (2) is valid as written; i.e., so that $\sigma$ will be the geometric cross section if the target is a smooth, perfectly conducting sphere.

Ignoring questions of efficiency and assuming that the same antenna is used both to transmit and to receive, we may recast Eq. (2) as:

$$P_r \sim \left[ P_t \frac{D^4}{\lambda^2} \right] \frac{\sigma}{(4\pi R^2)^2} \quad . \tag{4}$$

Except for the fact that $\sigma$ may be a function of $\lambda$, the terms in square brackets represent quantities determined exclusively by the radar system whereas the other terms -- $\sigma/(4\pi R^2)^2$ -- are beyond the experimenter's control and depend only on the properties of the target.

To maximize the received power, the most desirable parameter to adjust is the antenna diameter whose influence grows with the fourth power. In view of the $\lambda^{-2}$ dependence, one wants the antenna to operate efficiently at as high a frequency (as small a $\lambda$) as possible. Of course, to utilize this capability, high-power transmitters at these frequencies must also be available although, as can be seen, the performance increases least rapidly with increases in transmitted power.

The quantity $\sigma/(4\pi R^2)^2$ is usually termed the <u>path loss</u> and is conventionally described in $db/m^2$ ("dee bee per square meter"). That is, one expresses the path loss L in units of square meters and calculates $10 \log_{10} L$ which thus yields the path loss in $db/m^2$. For instance, $L = 10^{-25} m^{-2}$ implies $L = -250 \, db/m^2$. The path loss is simply a measure of the weakening of the echo due to the distance and size of the target.

The values of path loss vary greatly for the different objects in the solar system. For the Moon, $L \approx -247 \, db/m^2$ whereas for Venus at inferior conjunction

(closest approach to the Earth) $L \approx -314$ db/m$^2$. Thus the echo from Venus is roughly $10^7$ times smaller than that from the Moon. For Pluto, $L \approx -400$ db/m$^2$ and represents, logarithmically speaking, a slightly bigger step from Venus at inferior conjunction than was the one from the Moon to Venus. Figure 1 shows path losses and two-way (round-trip) echo delays for the Sun, planets, and some satellites as well as for the asteroid Icarus which will make a close approach to the Earth in mid-June 1968.

The second most important factor in determining the detectability of the echo signal is the system noise temperature which includes external contributions (e. g. , sky background noise) as well as the noise of the receiver, waveguide losses, etc. One has here the classical problem of detecting a weak signal in the presence of additive, gaussian noise. The feasibility of detecting the signal depends not only on the received power and system noise temperature, but also on many other factors such as the bandwidth of the signal, the integration time, etc. We cannot probe these details here and shall merely point out qualitatively that to achieve detection the signal power must exceed that of the fluctuations of the noise which implies, of course, that the experimenter strives for as low a system noise temperature as possible. Signals as weak as $10^{-21}$ watts are now routinely detectable with some planetary radar facilities after several hour's integration.

The overall sensitivity of radar systems has undergone an explosive, but nonetheless steady, growth since the end of World War II. The ability to overcome path loss has increased, on average, by 5. 5 db each year -- that is, by almost a factor of four per year -- for the past 20 years. Continued advances in technology show that this growth rate can be sustained, with "reasonable" outlays of capital (a few tens of millions of dollars), for at least another decade. Whereas the Moon was first observed by radar only in 1946 and Venus in 1961, Pluto may well be detected by the early to middle 1970's.

At present the major instruments being used for lunar and planetary radar astronomy are the Haystack facility of the M. I. T. Lincoln Laboratory at Tyngsboro, Massachusetts; the giant radar (1000-ft diameter spherical antenna) of Cornell's Arecibo Ionospheric Observatory in Puerto Rico; and the radar systems of the C. I. T. Jet Propulsion Laboratory in Goldstone, California. The Jodrell

·Bank radar in England and the Crimean radar tracking facility in the U.S.S.R. have not recently been active in planetary radar work, presumably because of a lack of sufficient system sensitivity.

## III. TIME DELAYS AND DOPPLER SHIFTS

From the characteristics of the received echo signal, many properties of the target may be inferred as mentioned earlier. For the moment we shall concentrate simply on the round-trip time delay and on the Doppler shift of the echo, assuming that the medium between radar site and target is a vacuum. For a point target, the time delay, by definition, is the elapsed time, as measured at the radar site, between transmission of a signal and detection of the echo. We may express this delay in terms of the coordinates of the radar site and target with respect to the Sun. That is, we assume the existence of a solar ephemeris for each body in the solar system. (We defer until Section V both a more precise discussion of spatial and time coordinates consistent with the theory of general relativity and the influence of the latter on the expression for time delay.) Under the assumption that light propagates rectilinearly at a constant speed (not valid in general relativity), the delay $\tau$ is given by

$$\tau(t_3) \;=\; \frac{1}{c}\{\;|\;\vec{r}_2(t_2) - \vec{r}_1(t_1)\;| + |\;\vec{r}_2(t_2) - \vec{r}_3(t_3)\;|\;\} \quad , \tag{5}$$

where $t_1$ is the time of transmission, $t_2$ the time of reflection at the (point) target, and $t_3$ the time of echo reception. The position vectors $\vec{r}_1(t_1)$, $\vec{r}_2(t_2)$, and $\vec{r}_3(t_3)$ represent, respectively, the position with respect to the Sun of the radar site at $t_1$, of the target at $t_2$, and of the site at $t_3$. The speed of light is denoted by $c$. This expression for $\tau$ looks deceptively simple. In fact, given that an echo signal was received at $t = t_3$, one cannot obtain a closed form solution for the time $t_2$ at which the signal was reflected by the target -- even for circular orbits. But because $(v/c) << 1$, we can use an iterative scheme which yields sufficient accuracy for present practice after at most two iterations.

The derivation of the Doppler shift expression -- at least when general relativistic effects are neglected -- is more interesting. We consider first a derivation consistent with special relativity and employ a 4-vector notation: $A_\mu = (A_o, \vec{A})$ with scalar products given by

$$A_\mu B^\mu \;\equiv\; A_o B_o - \vec{A} \cdot \vec{B} \tag{6}$$

where $\vec{A} \cdot \vec{B}$ is the ordinary 3-vector (spatial) dot product. With the speed of light set equal to unity, the velocity 4-vector for a mass point may be written as

$$v^{\mu} = (1 - \dot{\vec{v}} \cdot \vec{v})^{-1/2} (1, \vec{v}) ; \qquad c = 1 , \tag{7}$$

where $\vec{v}$ is the ordinary velocity 3-vector and $\vec{v} \cdot \vec{v} \equiv v^2$. For the propagation of electromagnetic signals, we note that a wavefront is specified by a surface in an inertial space-time coordinate frame; its normal, the <u>wave number vector</u> $k_{\mu}$, in this frame, describes the direction of propagation and has a magnitude such that

$$k_{\mu} = f_s (1, \vec{e}) , \tag{8}$$

where $f_s$ is the frequency of the wave as measured in the inertial frame and $\vec{e}$ is the unit vector in the direction of spatial propagation in this same frame. The frequency of this wave measured by an observer whose velocity is $v^{\mu}$ with respect to this frame will be

$$f = k_{\mu} v^{\mu} = f_s (1 - v^2)^{-1/2} (1 - \vec{v} \cdot \vec{e}) . \tag{9}$$

Now let us assume that all times, distances, and velocities are expressed in the rest frame of the Sun and suppose that a monochromatic wave is transmitted from Earth at $t_1$, reflected from the planet at $t_2$, and received back on the Earth at $t_3$. If the frequency measured at $t_1$ by the Earth observer (4-vector velocity $v_1{}^{\mu}$ is $f_1$, then

$$f_1 \equiv k_{\mu} v_1{}^{\mu} = f_s (1 - v_1{}^2)^{-1/2} (1 - \vec{v}_1 \cdot \vec{e}_{12}) , \tag{10}$$

where $\vec{e}_{12}$ is the unit vector pointing from the position of the Earth observer at $t_1$ to the position occupied by the (point) target at $t_2$. Similarly the frequency measured at the target at $t_2$ by an observer at rest with respect to the target (4-vector velocity $v_2{}^{\mu}$) will be given by

$$f_2 = k_\mu v_2{}^\mu = f_s(1 - v_2{}^2)^{-1/2} (1 - \vec{v}_2 \cdot \vec{e}_{12}) \quad . \tag{11}$$

Since, as measured by the target observer, the reflected wave has the same frequency both before and after reflection we may also write $f_2$ as

$$f_2 = k_\mu' v_2{}^\mu = f_s' (1 - v_2{}^2)^{-1/2} (1 - \vec{v}_2 \cdot \vec{e}_{23}) \quad , \tag{12}$$

where

$$k_\mu' \equiv f_s' (1, \vec{e}_{23}) \tag{13}$$

represents the wave number 4-vector for the reflected wave with $f_s'$ being its frequency as measured in the solar rest frame and $\vec{e}_{23}$ being a unit vector in the direction from the spatial position of the target at $t_2$ to that of the Earth observer at $t_3$. Upon reception of the echo, the Earth observer measures

$$f_3 = k_\mu' v_3{}^\mu = f_s' (1 - v_3{}^2)^{-1/2} (1 - \vec{v}_3 \cdot \vec{e}_{23}) \quad . \tag{14}$$

We seek, of course, $\Delta f \equiv f_3 - f_1$ which is the Doppler shift. From Eqs. (11) and (12) we obtain

$$f_s' = f_s(1 - \vec{v}_2 \cdot \vec{e}_{12})(1 - \vec{v}_2 \cdot \vec{e}_{23})^{-1} \quad , \tag{15}$$

which leads through Eqs. (10) and (14) to

$$\Delta f = f_1 \left[ \left( \frac{1 - v_1^2}{1 - v_3^2} \right)^{1/2} \left( \frac{1 - \vec{v}_2 \cdot \vec{e}_{12}}{1 - \vec{v}_1 \cdot \vec{e}_{12}} \right) \left( \frac{1 - \vec{v}_3 \cdot \vec{e}_{23}}{1 - \vec{v}_2 \cdot \vec{e}_{23}} \right) - 1 \right]. \tag{16}$$

If we assume $\vec{v}_1 \approx \vec{v}_3$ and $\vec{e}_{23} \approx -\vec{e}_{12}$, we get the approximate relation

$$\Delta f \approx -2f_1(\vec{v}_2 - \vec{v}_1) \cdot \vec{e}_{12} [1 - (\vec{v}_2 - \vec{v}_1) \cdot \vec{e}_{12} + O(v^2)] \tag{17}$$

after simple algebraic manipulation.

It is instructive to present an alternative form for the Doppler shift whose derivation is based on only very fundamental principles and is applicable to electromagnetic propagation in any non-dispersive medium. Let all times and frequencies now be as measured by the Earth observer. Further, let $f_3(t_3)$ be the frequency received at $t_3$ and $\tau(t_3)$ be the two-way time delay associated with the wave "crest" received at $t_3$. The next crest will be received at approximately $t_3 + f_3^{-1}$ since the time interval between the reception of successive crests is approximately the reciprocal of the number of oscillations per unit time. (The expression is not necessarily exact because $f_3(t_3)$ is not necessarily equal to $f_3(t_3 + f_3^{-1})$, i.e., the frequency $f_3$ may be changing with time.) These two crests were, of course, transmitted at time $t_3 - \tau(t_3)$ and (approximately)at time $t_3 + f_3^{-1} - \tau(t_3 + f_3^{-1})$. These latter times are separated (approximately) by the inverse of the frequency of the transmitted wave. Thus,

$$f_1^{-1} \approx [t_3 + f_3^{-1} - \tau(t_3 + f_3^{-1})] - [t_3 - \tau(t_3)]$$

$$\approx f_3^{-1} + \tau(t_3) - \tau(t_3 + f_3^{-1}) \approx f_3^{-1} - \dot{\tau}(t_3) f_3^{-1} , \qquad (18)$$

where in the last step we replaced $\tau(t_3 + f_3^{-1})$ by the first two terms of its Taylor expansion, the dot superscript denoting differentiation with respect to time. Now the first part of Eq. (18) was approximate only because we treated the (finite) separation of adjacent crests. If we consider the fact that frequency is simply the time derivative of phase and deal with infinitesimal increments of phase, this first part can be made exact. By the same token, the last part involving only the first two terms of the Taylor expansion will by themselves be exact in the limit when $f_3^{-1}$ is replaced by the equivalent infinitesimal phase increment. The combined result of these two effects proves that the equation

$$f_1^{-1} = f_3^{-1} - \dot{\tau}(t_3) f_3^{-1} \qquad (19)$$

is, in fact, exact and not approximate. (The only reason infinitesimal phase increments were not introduced ab initio was to make clearer the basic idea behind the derivation.) Solving Eq. (19) for $\Delta f$ shows that

$$\Delta f \equiv f_3 - f_1 = -f \dot{\tau} , \tag{20}$$

where all times and frequencies are measured by the observer at reception. The corresponding equation based o time of transmission [e. g., $\tau(t)$ now will represent the time delay for a wave transmitted at time t] is easily shown to be

$$\Delta f = -f_1 \frac{\dot{\tau}}{(1 - \dot{\tau})} \tag{21}$$

We stress again that these equations [(20) and (21)] depend only on propagation in a nondispersive medium so that the phase delay (delay for a particular crest) and the group delay ($\tau$) are equal. In particular the result is independent of the theory of gravitation assumed -- provided only that it is nondispersive.

How may the form of the Doppler shift given, say, in Eq. (20) be related to that in Eq. (16)? Or, equivalently, how may $\dot{\tau}$ be expressed in terms of the $\vec{v}$'s and $\vec{e}$'s of the solar ephemeris? We shall not attempt an exact derivation but merely calculate $\dot{\tau}$ in terms of the $\vec{v}$'s and $\vec{e}$'s accurate to second order in v to demonstrate the equivalence between Eqs. (20) and (17). For $v << c$, we see that $t_3 - t_2 \approx t_2 - t_1 \approx \tau(t_3)/2$, which, in view of Eq. (5), leads to

$$\tau(t_3) = | \vec{r}_2(t_3 - \frac{\tau}{2}) - \vec{r}_1(t_3) | + | \vec{r}_2(t_3 - \frac{\tau}{2}) - \vec{r}_1(t_3 - \tau) | \tag{22}$$

when we again set c = 1. Differentiating Eq. (22) yields

$$\dot{\tau} \equiv \frac{d\tau(t_3)}{dt_3} \approx \frac{d}{dt_3}([\vec{r}_2(t_3 - \frac{\tau}{2}) - \vec{r}_1(t_3)] \cdot [\vec{r}_2(t_3 - \frac{\tau}{2}) - \vec{r}_1(t_3)])^{1/2}$$

$$+ \ldots$$

$$\approx |\vec{r}_2 - \vec{r}_1|^{-1}\{[\vec{r}_2 - \vec{r}_1] \cdot [\vec{v}_2 - \vec{v}_1 - \frac{\dot{\tau}}{2} \vec{v}_2]$$

$$+ [\vec{r}_2 - \vec{r}_1] \cdot [\vec{v}_2 - \vec{v}_1 - \frac{\dot{\tau}}{2} \vec{v}_2 + \dot{\tau} \vec{v}_1]\} , \tag{23}$$

since, for example, $\vec{v}_i \equiv \dot{\vec{r}}_i$ (i = 1, 2) and

$$\frac{d}{dt_3} \vec{r}_1(t_3 - \tau(t_3)) = \vec{v}_1(t_3 - \tau(t_3)) [1 - \dot{\tau}] . \tag{24}$$

Solving for $\dot{\tau}$ and using Eq. (20) we obtain, in agreement with Eq. (16),

$$\Delta f \; = \; -f_1 \, \dot{\tau}(t_3) \; \approx \; -2f_1 \, \vec{e}_{12} \cdot (\vec{v}_2 - \vec{v}_1) \, \{1 - \vec{e}_{12} \cdot (\vec{v}_2 - \vec{v}_1)\} \; , \qquad (25)$$

where, by definition,

$$\vec{e}_{12} \; = \; \frac{\vec{r}_2 - \vec{r}_1}{|\vec{r}_2 - \vec{r}_1|} \qquad . \qquad\qquad (26)$$

We also assumed in the above derivation that the change in the speed of the observer was negligible between $t_1$ and $t_3$ so that the rate of the observer's clock in this interval would be uniform with respect to that of the clock in the solar rest frame; under this condition the differences in the times measured in the two frames do not affect the result: the time derivative with respect to the observer's time of the delay measured by the observer equals the time derivative with respect to solar-rest-frame time of the solar-rest-frame delay.

Until now, we have considered only point targets. Obviously, planets do not fall into this category. The radius of a planet is typically at least six orders of magnitude greater than the wavelength of the radiation emitted by the radar. How then can we determine from where on the planet any particular part of the received echo signal was reflected? It turns out that we can use time-delay and Doppler-shift measurements to determine an almost unambiguous map of the target planet, i.e. an association of parts of the echo with their origins at physical locations on the planet. This method utilizes the principle of delay-Doppler mapping which we shall now describe. Assume the planet to be spherical and consider the transmission from the radar and towards the target of a monochromatic, very narrow pulse of energy. These two requirements on the transmission are not mutually contradictory from a practical point of view since we choose the pulse length $\Delta t$ to satisfy $\lambda \ll c \, \Delta t \ll \rho$. The first point on the target encountered by the pulse -- the subradar point S -- lies on the surface at the intersection of the line from the radar site to the planet's center. Hence the first part of a received echo signal must arise from reflections at or near the subradar point. The echo received slightly later arises from reflections from a ring or annulus on the planet's surface, suitably displaced from S and lying in a plane perpendicular to

the line from the radar site to the planet's center. By considering the echo as a function of delay, we therefore have a one-dimensional mapping of the surface: From the time of reception we know the ring from which the echo originated. If, for example, a right-hand circularly polarized pulse is transmitted and left-hand circularly polarized energy is received, then typically the received echo signal power falls rapidly as a function of delay. (Recall that upon reflection from a smooth sphere, the direction of polarization is reversed.) Thus the reflection is essentially specular, most of the backscattering power originating from the first few Fresnel zones around S. Backscattering from other parts of the planet arises primarily because of the roughness of the surface (undulations, discontinuities, rocks, etc.). The longer the probing wavelength, the smoother the planet appears (the radiation is, in effect, insensitive to roughness on a scale much smaller than $\lambda$) and the more rapidly the received power decreases with delay. As an example, we point out that the backscattered power in a given delay increment is about three orders of magnitude smaller 10 msec (two-way) back into Venus from S than at S for $\lambda \approx 70$ cm. The total two-way depth of Venus is about 40 msec since its radius is about 6050 km (only one hemisphere, of course, can contribute to the echo).

Now we consider the other coordinate: Doppler shift. The planet as viewed from the radar site appears to be rotating. This rotation consists of two components: an apparent rotation, due to the relative orbital motion of radar and target, and the inertial (or sidereal rotation) of the planet. Thus, even were the planet not rotating it would seem to rotate to an Earth observer who sees the planet sweep by in space; at different times different parts of the planet's surface would occupy the subradar point. Because of the rotation, each region on the surface of the planet will impart a particular Doppler shift to its echo, depending on the component of its relative velocity along the line to the radar site [see, e.g., Eq. (17)]. Since we are concerned now only with associating points on the surface with particular Doppler shifts, it is convenient to ignore the projection along the line of sight to the radar of the relative velocity between S and the radar site. That is, we need consider only the Doppler shifts relative to that for S. We denote the Doppler shift for a surface point P relative to that for S by $f_r(P)$:

$$f_r(P) \equiv \Delta f(P) - \Delta f(S) \quad . \tag{27}$$

Now S has the same component of velocity along the line of site to the radar as does the planet's center of mass. Thus considering only these relative Doppler shifts is equivalent to neglecting the contribution to the Doppler shift of the relative translational motion of radar site and center of mass and leaving only the contribution attributable to the total (apparent plus sidereal) rotation of the planet about its center of mass. The corresponding total angular velocity $\vec{\omega}$ will therefore be the vector sum of the apparent $\vec{\omega}_a$ and sidereal $\vec{\omega}_s$ contributions. If $\vec{\rho}(P)$ is the radius vector from the planet's center to P, then to lowest order in $\omega$,

$$f_r(P) = -2f_1[\vec{\omega} \times \vec{\rho}(P)] \cdot \vec{e}_{12} = 2f_1[\vec{\omega} \times \vec{\rho}(P)] \cdot \vec{e}_{23} \ . \tag{28}$$

Referring to Figure 2, in which $(\vec{R}/R) \equiv \vec{e}_{23} \approx -\vec{e}_{12}$, we see that the relative velocity $\vec{v}_r(P)$ can be expressed as

$$\vec{v}_r(P) \equiv \vec{\omega} \times \vec{\rho}(P)$$

$$= \rho\omega \sin \theta' \left\{ \sin \varphi' \ \frac{\vec{\omega} \times (\vec{\omega} \times \vec{e}_{23})}{\omega \ |\ \vec{\omega} \times \vec{e}_{23}\ |} + \cos \varphi' \ \frac{\vec{\omega} \times \vec{e}_{23}}{|\ \vec{\omega} \times \vec{e}_{23}|} \right\} \ . \tag{29}$$

Hence,

$$\vec{v}_r(P) \cdot \vec{e}_{23} = -\rho\omega \sin \psi \sin \theta' \sin \varphi'$$

$$= \rho\omega \sin \psi \sin \theta \sin \varphi \qquad , \tag{30}$$

where the last step follows from the preceding by simple geometric considerations (see Figure 2). Substituting into Eq. (28) yields

$$f_r(P) = 2f_1 \rho\omega \sin \psi \sin \theta \sin \varphi \ . \tag{31}$$

Again by simple geometry, we see that the distance of P from the $\vec{\omega} - \vec{R}$ plane is $\rho \sin \theta \sin \varphi$ which implies that all points with the same value of $f_r$ lie on planes parallel to the $\vec{\omega} - \vec{R}$ plane, i.e. lie on rings on the surface whose planes are perpendicular to those of the delay rings. The projection of the visible surface onto a plane normal to $\vec{e}_{23}$ is shown in Figure 3 where the locus of points of constant delay appear as rings and the locus of points of constant Doppler as "strips". By analyzing the echo in both delay and Doppler we can therefore construct a map of the planet's surface. But this map has unusual properties: it is double-valued

(the same delay and Doppler coordinates correspond in general to two points on the planet's surface as shown by the two small darkened regions in Figure 3) and the transformation from delay-Doppler to spherical coordinates is singular at the apparent equator as evidenced by the enlarged common area of the delay ring and Doppler strip at the point of tangency there.

The mathematical relation between the delay-Doppler and spherical coordinates is easily derived. In addition to Eq. (31) we make use of the relation between the (two-way) delay $\tau_r$ to P relative to the delay to S. From Figure 2, we see that

$$\cos \theta = 1 - (c \tau_r / 2 \rho) \quad . \tag{32}$$

The Jacobian J of the transformation can now be calculated in a straightforward manner yielding

$$J = \frac{\partial(\tau_r, f_r)}{\partial(\theta, \varphi)} = 4 f_1 \omega \rho^2 \sin \psi \sin^2 \theta \mid \cos \varphi \mid$$

$$= [\tau_r(\frac{4\rho}{c} - \tau_r)]^{1/2} [(f_1 \omega \sin \psi)^2 \tau_r(\frac{4\rho}{c} - \tau_r) - f_r^2]^{1/2} \quad . \tag{33}$$

The surface element of area is thus expressible as

$$\rho^2 \sin \theta \, d\theta d\varphi = \frac{\rho \, d\tau_r d f_r}{2[(f_1 \omega \sin \psi)^2 \tau_r(\frac{4\rho}{c} - \tau_r) - f_r^2]^{1/2}} \quad . \tag{34}$$

It is clear from the above that, for a given $\tau_r$, the possible values of $f_r$ are bounded by

$$|f_r| \leq f_1 \omega \sin \psi [\tau_r(\frac{4\rho}{c} - \tau_r)]^{1/2} \quad ; \tag{35}$$

that is, any given delay annulus intercepts Doppler strips whose values lie between the limits shown. Therefore the bandwidth B of the spectrum of the echo from a given annulus is

$$B(\tau_r) = 2 f_1 \omega \sin \psi [\tau_r(\frac{4\rho}{c} - \tau_r)]^{1/2} \quad . \tag{36}$$

If the backscattered power reflected from an annulus is independent of $\varphi$, then the corresponding power spectrum $P(f_r, \tau_r)$ has a characteristic shape:

$$P(f_r, \tau_r) = \begin{cases} P_o \left[ \left( \dfrac{B(\tau_r)}{2} \right)^2 - f_r^2 \right]^{-1/2} & ; \ |f_r| \leq B(\tau_r)/2 \\ 0 & ; \ |f_r| > B(\tau_r)/2 \ , \end{cases} \tag{37}$$

where $P_o$ is independent of frequency. The singularity of J causes P to become infinite at $f_r = \pm B/2$ and is useful for estimating bandwidths as will be described below in connection with the determination of planetary rotation rates.

The above idealized picture of delay-Doppler mapping is, of course, not realized in practice. The finite extent of the transmitted pulses, the necessity to integrate the echoes from successive pulses with proper account of relative phases, the filters used in the receiver system, and the backscattering law obeyed by the planet all combine in influencing the association of points on the planet's surface with particular parts of the echo signal. These topics, however, are too complicated and too far afield to delve into here. Our main point -- to show that, in principle at least, accurate measurements could be made of time delays and Doppler shifts of echoes from a particular part of a planetary target -- has nonetheless been achieved. Our treatment was more detailed than necessary simply to establish this point because we wished at the same time to lay the groundwork for our discussion in Section VI of the determination of planetary rotation rates from radar data.

A theoretical discussion of the accuracy with which measurements can be made of the time delay and Doppler shift associated with, say, the subradar point is also beyond our scope. On the other hand, the achieved results are crucial to it. For the planets Venus and Mercury, measurements made with the Haystack radar have errors in time delay of only 10 μsec even when these planets are farthest from Earth (near superior conjunction). When closest, for example, the errors in Earth-Venus time-delay measurements are no more than about 3 μsec. Near opposition this spring, delays to the subradar point on Mars were measured with comparable accuracy. The errors in Doppler-shift measurements range from several tenths of a Hertz to 2 Hz, depending on conditions, for the Haystack value of $f_1$ which is

about $8 \times 10^9$ Hz. To show clearly the implication of these accuracies, we consider the fractional errors:

$$\frac{\delta \tau}{\tau} \approx \frac{10^{-5} \text{ sec}}{2 \times 10^3 \text{sec}} \approx 5 \times 10^{-9} \left. \begin{array}{c} \\ \\ \\ \\ \\ \end{array} \right\} \qquad , \qquad (38)$$

$$\frac{\delta \Delta f}{\Delta f} \approx 6 \times 10^{-7}$$

where for $\Delta f$ we used a radial velocity of 20 km/sec. By contrast the claimed accuracy for determining the direction of a planet from Earth with respect to the stellar background is about $2 \times 10^{-6}$ radians. Although of a different type, the radar time-delay measurements are thus between two and three orders of magnitude more accurate than conventional optical ones. More important from the theoretical point of view, the relative errors are no longer large compared to $v^2/c^2$. Hence we can expect relativistic effects to be discernible thereby providing the opportunity to subject the basic theory of gravitation -- general relativity -- to additional experimental tests.

## IV. GENERAL RELATIVITY - THEORY

Before discussing the experimental tests in detail, we shall give a brief introduction to general relativity. This theory was given its definitive form by Einstein in 1916 and is unique in having achieved almost universal acceptance with very few experimental verifications. The reasons are not hard to find. Such was and is the reputation of Einstein that any theory bearing his imprimatur would almost automatically be accepted until experimental proof of its deficiencies could be given. But experimental tests of general relativity are not easy to come by: The differences between its predictions and those of the older theory of gravitation -- Newton's -- are generally minuscule. So small, in fact, are the differences that, even with modern technology, one most often requires the solar mass to disclose these -- a far grander scale of experiment than physicists have become accustomed to requiring. Only in the most sophisticated and sensitive experiments will the Earth's mass suffice.

We might first enquire why Einstein felt the need to devise a theory to supplant Newton's. The latter's laws of motion and gravity had proven most adequate

for more than two centuries before Einstein came upon the scene.   Even today
the other lecturers at this symposium are using them without any question.   Al-
though by the beginning of the 20th Century there was a well established, albeit
small, discrepancy between Newton's predictions and observations, Einstein's
motivation characteristically was one of fundamental principle.   He objected to
Newton's law of gravity because of its action-at-a-distance feature: The force
exerted on body A at a given (Newtonian) instant by body B depends only on the
position of body B at that instant -- no matter how greatly separated A and B
might be.   In other words the speed of propagation of the force was a sumed by
Newton to be infinite.   This troubled Einstein because, based on his theory of
special relativity, he expected that no signal (or force) could propagate faster
than the speed of light.

The primary principle used by Einstein to guide his thoughts in creating a
new theory was the so-called principle of equivalence.   In its weak form, it
states that the ratio of the gravitational to the inertial mass of a body is inde-
pendent of its composition.   In its strong form, the statement is that, locally,
a gravitational force is indistinguishable from an acceleration in inertial space.
That is, the effects of a gravitational field, locally, are completely equivalent
to an acceleration in inertial space.   Hence any body (of small enough mass and
physical extent) would follow the same trajectory if give  the same initial   con-
ditions -- regardless of its composition.   The trajectory is not influenced by the
body traversing it, provided the body has a sufficient y small mass.   In this
sense, the trajectory is an inherent property of the space in which the body is
moving.   This principle led Einstein to assume that the effect of a gravi ational
field could be adequately described by an appropriate geometry.   After some
study, he chose Riemannian geometry as being rich enough, with the metric to be
determined by the sources of the gravitational field.

It was clear to Einstein, as well as to others, that a field equation was needed
which would be a generalization of Poisson's equation in Newtonian theory:

$$\nabla^2 V \quad = \quad 4\pi d \quad , \tag{39}$$

with the generalization somehow incorporating a finite speed of propagation.

[In Eq. (39), of course, V denotes the gravitational potential and $\underline{d}$ the mass. density.] A generalization involving the scalar wave equation had been developed by Nordström in 1911, but had failed because of an unfortunate attribute: the predictions disagreed with experiment. (The wrong result was obtained for the advance of Mercury's perihelion.) Einstein tried a tensor equation, reasoning approximately as follows: Special relativity showed matter and energy to be equivalent and one might expect the mass density to be replaced in the generalized equation by the total energy density from all sources. But since the scalar theory fails and special relativity also shows energy and momentum to be connected in a 4-vector, Einstein proposed that the appropriate generalization of mass density would be the (symmetric) energy-momentum tensor $T_{\mu\nu}$. The generalization of Poisson's equation would therefore look like

$$G_{\mu\nu}(g_{\mu\nu}) = \kappa\, T_{\mu\nu} \quad , \tag{40}$$

where the left side depends on the metric tensor $g_{\mu\nu}$ -- the analog of the gravitational potential in Poisson's equation -- and $\kappa$ is a constant. (The Greek letter indices run from 1 to 4.) What should be the form of $G_{\mu\nu}$? From the principle of equivalence it follows that locally the effects of a gravitational field can be transformed away (i.e., they can be nullified with an opposing acceleration introduced by a suitable coordinate transformation), resulting in a locally "flat" coordinate system. The energy momentum tensor satisfies the well known (special relativistic) conservation law: its divergence vanishes. If we assume that the divergence of $T_{\mu\nu}$ vanishes identically in one coordinate system, then it must vanish in all. The left side of Eq. (40) must then also have a vanishing divergence. Requiring in addition that $G_{\mu\nu}$ be a function only of $g_{\mu\nu}$ and its derivatives -- no higher than the second (again in analogy to Poisson's equation) -- determines the generalized field equations uniquely:

$$R_{\mu\nu} - \frac{1}{2}\, g_{\mu\nu} R + \lambda g_{\mu\nu} = -\,\kappa\, T_{\mu\nu} \quad , \tag{41}$$

where $R_{\mu\nu}$ and $R$ are, respectively, the doubly and quadruply contracted Riemann curvature tensor. The undetermined constant $\lambda$, termed the cosmological constant, is usually set equal to zero; its history is interesting but again beyond

the intentions of this brief, heuristic development.  The final field equations
are

$$R_{\mu\nu} - \frac{1}{2} g_{\mu\nu} R \; = \; - \kappa T_{\mu\nu} \quad . \tag{42}$$

Because of symmetry these 16 equations reduce to 10, and because the vanish-
ing of the divergence of each side yields 4 conditions (Bianchi identities), the
number of equations is effectively reduced to 6.

What about the equations of motion?  In Newtonian theory we must postu-
late equations of motion separately from the field equations, viz.

$$m \; \ddot{\vec{r}} \; = \; - \nabla V(\vec{r}) \quad . \tag{43}$$

Einstein at first also proposed separate equations of motion;  in particular he
assumed that mass points would follow geodesics in the Riemannian space de-
termined by the metric tensor $g_{\mu\nu}$ found from the solution to the field equa-
tions.  These geodesic equations are

$$\left. \begin{array}{l} \dfrac{d^2 X^{\mu}}{ds^2} + \Gamma^{\mu}_{\rho\sigma} \dfrac{dX^{\rho}}{ds} \dfrac{dX^{\sigma}}{ds} \; = \; 0 \\[4ex] \Gamma^{\mu}_{\rho\sigma} \; = \; \dfrac{1}{2} g^{\mu s} (g_{\sigma s, \rho} + g_{\rho s, \sigma} - g_{\rho\sigma, s}) \\[4ex] g_{\rho\sigma, s} \; \equiv \; \dfrac{\partial g_{\rho\sigma}}{\partial X^s} \\[4ex] g_{\mu\nu} g^{\lambda\nu} = \; \delta^{\lambda}_{\mu} \end{array} \right\} \quad , \tag{44}$$

with the interval ds being given by

$$ds^2 \; = \; g_{\mu\nu} dX^{\mu} dX^{\nu} \tag{45}$$

which is assumed to vanish ($ds^2 = 0$) for light rays, determining thereby the
null geodesics.  (Summation over repeated indices is implied in all equations. )

In 1927 Einstein discovered an amazing fact: The equations of motion could be derived directly from the field equations and, as one no doubt expected, the results agreed with the previously assumed equations of motion. That the field equations should themselves determine the equations of motion could only happen with nonlinear field equations. For linear field equations, as in Newtonian theory, the applicability of the principle of superposition implies that these equations could not determine equations of motion. The general necessary and sufficient conditions for field equations to determine equations of motion have, however, not been established.

The solutions to Einstein's equations, as is well known, are notoriously hard to come by. The N-body problem has been considered only formally by making expansions in inverse powers of $c$ to obtain ordinary differential equations for each stage of the expansion. The constant $\kappa$ is established by correspondence of the lowest-order solution with the Newtonian result; it is simply related to the gravitational constant. What to choose for $T_{\mu\nu}$ is considered by some to be a moot question. In particular Einstein remained unhappy with the marriage he brought about between $T_{\mu\nu}$ and $G_{\mu\nu}$; his attempts at solutions were always referred to empty space where $T_{\mu\nu}$ vanishes. Various schools (e.g., Fock's in the U.S.S.R. and Infeld's in Poland) have developed which carry on bitter polemics in regard to the proper attack on, for example, the N-body problem. Yet all get the same result, at least up to the post-Newtonian terms which is as far as anyone has attempted a detailed development. (The post-Newtonian is the approximation which includes only the next order terms in $c^{-1}$ after the Newtonian ones.)

What problems can be solved in closed form? The $g_{\mu\nu}$ appropriate for a single mass point can be determined. Demanding that, outside the mass point, the metric tensor be spherically symmetric, static (i.e., independent of time), and "flat" at spatial infinity, one is led to the celebrated Schwarzschild exterior solution which can be written as

$$ds^2 = \left(1 - \frac{2r_o}{r}\right) dt^2 - \left(1 - \frac{2r_o}{r}\right)^{-1} dr^2 - r^2(d\theta^2 + \sin^2\theta d\varphi^2)$$

$$= \left(1 - \frac{2r_o}{r}\right) dt^2 - \left(\delta_{ij} + \frac{2r_o X_i X_j}{r^2(r-2r_o)}\right) dX^i dX^j \quad , \tag{46}$$

where the Latin indices run from 1 to 3. Here $X_i(i = 1 \rightarrow 3)$ are cartesian coordinates and $r_o \equiv (GM/c^2)$ is the gravitational radius of the mass point. For the Sun

$$r_o \approx 1.5 \text{ km} \quad . \tag{47}$$

Requiring that the metric be flat (i.e., Minkowskian as in special relativity) at spatial infinity was Einstein's method of incorporating Mach's principle which states that (somehow) the distant matter in the universe is responsible for the inertia of "laboratory" masses. Speaking approximately, one can say that far from the mass point being studied ("spatial infinity") the behavior of test particles should be the same as in the inertial system postulated in special relativity, this behavior being assured by the distant matter which is really distant -- much further than the "spatial infinity" at which the boundary conditions are being applied.

We may use the Schwarzschild solution in particular to study the trajectories of test particles and light rays outside the Sun which represents the mass point. But the planets are not test particles; their masses are appreciable. How is the solution to the equations of general relativity affected by their presence? DeSitter showed that because of the relatively slight mass of the planets, all problems with moving bodies in the solar system could for all practical purposes be treated by the _ad hoc_ addition of the strictly Newtonian perturbations caused by the planets and other masses to the effect of the Schwarzschild solution for the Sun.

## V. GENERAL RELATIVITY - EXPERIMENTAL TESTS

### A. Equivalence of Gravitational and Inertial Mass

Except for null experiments, all tests of general relativity to date refer strictly to predictions based on the Schwarzschild metric. The most famous

null experiment is probably the test of the equivalence of gravitational and in-
ertial mass. Early in the 20th Century, Eötvös in essence balanced the gravi-
tational attraction of the Earth and the centrifugal force due to its rotation to
show that the ratio of gravitational to inertial mass was independent of compo-
sition to about 1 part in $10^8$. (Incidentally, Newton himself had used a simple
pendulum apparatus to establish this result to 1 part in $10^3$.) More recently,
Roll, Krotkov, and Dicke devised a more sensitive experiment utilizing the or-
bital motion of the Earth and deduced that the equivalence held to a few parts in
$10^{11}$. These results imply, for example, that the gravitational and inertial
properties of electrons, neutrons, protons, and binding energies are all equivalent.

    B.  <u>Classical Tests</u>

    The three "classical" tests of general relativity -- all originally discussed
by Einstein -- are: (1) the red shift of the frequency of light; (2) the deflec-
tion of the path of light rays by matter; and (3) the advance of orbital perihelia.
The first experiment concerns the predicted effect of the gravitational potential
on the measured frequency of a light wave. When the light wave passes from a
region of high gravitational potential (say, from near the Sun) to a region of low
potential (say, near the Earth) its measured frequency should appear red-shifted.
Of course, a violet shift is predicted for the reversed path. Since the original
proposal concerned comparing the frequency of the emission from a transition
between two given energy levels of an atom in the Sun's atmosphere, as meas-
ured on Earth, with the frequency of the corresponding emission from an Earth-
based atom, this test is traditionally known as the red-shift test. Unfortunately
measurements of emissions from the Sun proved hard to deal with because of
the difficulty in accounting for the ordinary Doppler shift due to the motion of the
atoms undergoing the transitions. Use of more massive stars, in which the
Doppler shift will be relatively less important, is not a panacea because of the
difficulty in accurately estimating its mass and radius. In fact, the most accu-
rate red-shift experiment so far performed was done completely in a terrestrial
laboratory using the difference in the Earth's gravitational potential over a 25 m
path. Needless to say, extraordinarily accurate frequency measurements were
required to detect any effect at all, the fractional predicted change in frequency
being about $10^{-15}$. By utilizing the Mössbauer effect, which narrows the line
width of $\gamma$-ray emission from an $Fe^{57}$ crystal, Pound and Rebka were able to

detect the gravitational frequency shift and, ultimately, Pound and Snider were
able to verify the predictions to within a 1% estimated error. Even with the aid
of the Mössbauer effect this accomplishment was by no means trivial: the frac-
tional line-width of the natural emission was $10^{-12}$. Yet with very clever and
painstaking experimentation, it proved possible to measure frequencies to 1 part
in $10^{17}$ or to 1 part in $10^5$ of the natural line-width. Despite the impressive ex-
perimental tour-de-force, the theoretical significance of this "red-shift" result
is often minimized since the same result can be deduced directly from the (weak)
principle of equivalence and conservation of energy. The formal structure of
general relativity is not required. On the other hand, perhaps one should view
this experiment as establishing the principle of equivalence for photons -- at
least to a 1% accuracy. About an order of magnitude higher accuracy is antici-
pated within the next few years by means of a somewhat different experiment --
the comparison of the frequency of (or, more precisely, the time kept by) an or-
biting hydrogen maser with an identical device on the ground. This experiment,
suggested by many, is actively being developed by Ramsey and Kleppner.

   The second experiment involves the prediction that the path of starlight is
deflected towards the Sun. The total change in angle $\alpha$ for a ray emanating far
from the Sun, passing by it, and receding to infinity is given by

$$\alpha \approx \frac{4r_o}{b} \tag{48}$$

where b is the impact parameter of the ray or, to a sufficient accuracy, the
distance of closest approach of the ray path to the center of the Sun. Knowing
$r_o$ [Eq. (47)] and the solar radius ($\rho_s \approx 7 \times 10^5$ km), one easily calculates
that the maximum value of $\alpha$ -- for a ray just grazing the solar limb -- is about
$1.''75$. The bending is clearly not a large effect; moreover, the effect decreases
inversely with the distance of closest approach so that, for example, at only $1^o$
from the Sun's center, the total bending will be less than $0.''5$. Attempts to ob-
serve this deflection began during the total solar eclipse of 1919 when expeditions
were sent to several parts of the world that were in the eclipse's path. Only dur-
ing a total eclipse was it felt possible to detect the deflection because, during or-
dinary daytime conditions, the glare from the scattered light in the atmosphere
prevented the detection of stars whose light rays passed near the Sun. But, even

during a total eclipse, how could the deflection be observed? The classical procedure was to take photographs of the star field in the vicinity of the Sun and measure the radial distance of each from the position corresponding to the center of the Sun. These distances on the photographic plates are, of course, directly related to the angular separation of the rays as they arrive at Earth. By comparing such plate measurements with corresponding ones made on photographs taken, say, six months earlier or later when the same star field is visible near midnight on Earth, one can check to see whether the differences in radial separation distances on the two sets of plates are in accord with the predicted deflections. Although straightforward in principle, this procedure encounters a number of experimental difficulties. The results are affected significantly by slight misalighments of plates with respect to the telescope axis, by radial distortions of the plates, and, to a lesser extent, by the quite different, turbulent atmospheric conditions prevailing during a total eclipse as compared to the relatively quiescent midnight atmosphere. Because of the infrequent occurrence and short duration of suitable eclipses the total observing time since 1919 has only mounted to about 90 minutes! The results have not been very definitive. Different observers measuring the same plates and the same observers measuring different plates have all come up with different results. More often than not, the differences are substantially greater than the combined, quoted probable errors accompanying the results. One might reasonably enquire how different observers using the very same plates could obtain such different results. The answer lies mostly in the necessity to make somewhat subjective plate corrections to account for radial distortions. Although no star has been observed closer than two solar radii from the Sun's center (maximum deflection about $0\rlap{.}''9$), the results are always quoted as the equivalent deflection at the limb of the Sun, with the scaling being based on the assumption that the deflection varies inversely with $b$ as predicted by general relativity. Numerically, values ranging from $1\rlap{.}''6$ to $2\rlap{.}''4$ have been reported with most probable errors being in the range of a few tenths of an arc-second. It is therefore usually considered that the predictions of general relativity for this phenomenon have been verified only semi-quantitatively, i.e. to about $\pm 25\%$.

Could the predicted gravitational deflection of electromagnetic waves be detected with radar, say by measuring the direction of arrival of the echoes received

from observations of a planet passing behind the Sun? One obvious advantage of
using the radio part of the spectrum would be that waiting for eclipse conditions
is not necessary. The atmospheric constituents, being small compared to the
radio wavelengths, don't cause appreciable scattering; likewise the ionosphere
is little cause for concern at wavelengths sufficiently short compared to the criti-
cal value for reflection. But, examining the effective angular resolution of the
antenna (given approximately from diffraction theory by $\lambda/D$), we find that the best
available in planetary radars is about 4' of arc or about 150 times larger than the
maximum predicted bending. It would therefore seem impossible to detect the
deflection by radar -- even if a suitable target were available. However some
months ago it occurred to me that by using a pair of separated radar receivers
(but only a single transmitter) in an interferometric arrangement, far greater
angular resolution becomes feasible. In effect, one can achieve the same angu-
lar resolution as with a giant antenna whose diameter equalled the distance be-
tween the pair forming the interferometer. Naturally, the received echo will be
weaker in the latter case; in addition there are ambiguities inherent in the inter-
ferometer configuration which are fortunately resolvable through a series of ob-
servations. The angular direction of the target is determined by a comparison of
the phases of the echo signal recorded at the two sites which of course must be
connected by a phase-stable link. With a 10-km baseline and $f_1 \sim 10^{10}$ Hz, the
angular resolution achievable in principle is on the order of $0.''01$.

We may next enquire about the theoretical expression for the deflection in the
radar-planet case. A relatively simple, but somewhat cumbersome, derivation
shows that the difference $\eta$ between the angles of arrival at the observer of a ray
travelling along a straight path and of one travelling along the curved path pre-
dicted by general relativity is given by

$$\eta = \frac{2r_o}{r_e} \left[ \tan\left(\frac{\theta}{2}\right) + O\left(\frac{r_o}{r_e}\right) \right] , \tag{49}$$

where $r_e$ is the Sun-Earth observer distance and $\theta$ is the angle between the line
from the Sun to the planet and the line from the Sun to the Earth observer. This
result for $\eta$ is really quite surprising: it is independent of $r_p$, the Sun-planet dis-
tance. That is, the difference in angle of arrival between rays propagating

rectilinearly and rays following the curved, relativistic path is the same regard-
less of the position on a given radial line from which the electromagnetic signal
emanates. How may we understand this result? The further out on a given radius
that the planet lies, the further away from the Sun will be the point of closest ap-
proach of the ray path to the observer; hence the less will be the maximum rate of
deflection. But the further the planet lies from the origin, the longer will be the
total ray path and, consequently, the greater will be the time over which the rate
of deflection will be integrated. These two opposing tendencies just compensate.
If we take the limiting form of Eq. (49), appropriate for ray paths passing close
to the solar limb, we find

$$\eta \approx \frac{4r_o}{b} \left( \frac{r_p}{r_e + r_p} \right) \; ; \; b = \frac{r_e r_p \sin \theta}{[r_e^2 + r_p^2 - 2r_e r_p \cos \theta]^{1/2}} << r_e, r_p \, .$$

$$(50)$$

The maximum value of $\eta$ is about $0\overset{''}{.}75$ for Venus the target planet. If we con-
sider the limit of Eq. (50) in which $r_p \to \infty$, we obtain agreement with Eq. (48).
However, if we let $r_e$ and $r_p$ both approach infinity, with $r_e = r_p$, we find that
$\alpha = 2\eta$. That such a result should hold follows from the different definitions of
$\alpha$ and $\eta$. The former represents the total deflection of a ray (i.e., the angle be-
tween the final and initial directions of propagation), whereas the latter is the
difference in angle at the receiver between the final directions of propagation of
the rays following rectilinear and curved paths, respectively. A simple draw-
ing suffices to show that $\alpha = 2\eta$ in all cases for $r_e = r_p$ because of the symmetry
of the paths about the bisector of the lines that connect the Sun and planet (source)
and the Sun and Earth (observer).

One might logically wonder how the predicted deflection in the radar-planet
case could be determined operationally. After all, with radar, we do not observe
the apparent positions of the planet against the stellar background. How would
we know, for example, which angle of arrival supported the predictions of general
relativity? From other optical and radar observations of the planets, all made
when the paths of the relevant electromagnetic waves never pass near the Sun, the
orbits can (and have) been determined very accurately. Hence one then knows
theoretically what angle of arrival to expect from the echo of radar signals trans-
mitted towards the planet when near superior conjunction (i.e., when the ray path
passes near the Sun). In this manner the predictions of general relativity concern-
ing deflection can be tested. Of course, the orbit calculations perforce required

some theoretical basis and, so, what is really involved is a complicated test of self-consistency of the theory and all of the data.

Aside from problems of interpretation, would it be feasible in practice to measure angles of arrival with errors substantially less than the predicted angle $\eta$, at least near superior conjunction? An a priori answer to this question must be based on a careful analysis of a number of factors, viz. the minimum possible $\underline{b}$ considering the "blinding" effects of the Sun through the antenna sidelobes, the effects of the Earth's atmosphere and the solar corona, the problems of maintaining absolute phase stability between widely separated receivers, the angular size of the target planet, the required radar-system sensitivity, the precise determination of receiver locations, etc. After examining all of these possible impediments to a successful experiment, I concluded that it was most likely feasible at radar frequencies in the vicinity of $10^{10}$ Hz. The main uncertainty concerns the properties of the Earth's atmosphere. A preliminary estimate for the accuracy achievable gives a result between 10 and 15% of the predicted relativistic deflection.

With the development of highly precise atomic clocks, it will soon be feasible to maintain phase stability between widely separated radio receivers for usefully long periods without recalibration. That is, each site would be slaved to its own atomic-clock standard which periodically would be recalibrated with the others to prevent excessive phase drift. It should prove possible to perform the calibration without the necessity of a direct radio link between the separate sites. To describe the techniques in detail would again take us too far afield.

In addition to the radar-planet possibility for measuring the gravitational bending of the path of electromagnetic waves, passive radio interferometry might be used for the same purpose: If point sources in the sky are found that are (1) periodically occulted or nearly occulted by the Sun, and (2) strong emitters of high frequency ($\approx 10^{10}$ Hz) radio waves, then their apparent direction in the sky could be monitored in the vicinity of solar occultation to measure the gravitational deflection of the path of the radio waves. The accuracy achievable might be on the order of 1% of the predicted bending. The improvement over the anticipated accuracy from radar-planet measurements is primarily due to the point source nature of the postulated target for the radio experiment. As in the radar proposal, the high frequency is desirable to minimize the effects of the solar corona which have

a frequency dependence varying as the inverse square. In principle, the coronal effects could be measured simultaneously by employing two frequencies. However if these were chosen from regions of the spectrum much below $10^{10}$ Hz, say by a factor of 5 or 10, then turbulence in the corona might well prove an insuperable barrier to precise measurements. With the use of atomic clocks, it might eventually prove feasible to utilize intercontinental separations of the radio receivers.

Although not strictly relevant to our main discussion, I shall mention several other interesting measurements that might be made with long-baseline radar or radio interferometers. The Earth's rotation (length of the day and polar wandering) could be monitored with perhaps an order of magnitude greater precision, and the separation of any two points on the Earth's surface could be determined with an error of maybe only a few centimeters! This latter possibility implies that the tidal motions of the Earth's crust could be measured and that intercontinental drifts of the order of a few centimeters per year would be detectable after several years. These possible applications must of course be considered speculative at the moment, but I strongly suspect they will be realities within a decade.

Why can't such techniques be used at optical wavelengths? The basic obstacle is the Earth's atmosphere which, at optical wavelengths, produces such a large number of rotations of the phase vector that it has, at least until now, proved infeasible in practice to use long-baseline interferometers to improve the accuracy with which the angular distance between two separated point sources (stars) could be measured. To put the matter quantitatively, the effective increase in the optical path attributable to the atmosphere is about 10m which corresponds at $\lambda = 5000 \overset{o}{A}$ to $2 \times 10^7$ wavelengths, or phase vector rotations. Even if the uncertainty in the optical path through the atmosphere were as low as 0.1%, there would still remain an ambiguous range of $2 \times 10^4$ phase rotations caused by the atmosphere. Use of a sufficiently wide bandwidth and long enough integration time to resolve the ambiguity would probably be prevented both by uncertainty in the wavelength dependence of the atmospheric index of refraction and rapidly varying, small-scale turbulence. One might immediately think of performing a satellite interferometer experiment outside the effects of the Earth's atmosphere where, in addition, one eliminates the problem of scattered sunlight. However, carrying out such an experiment in orbit is immensely difficult. The pointing requirements have not yet been achieved, let alone the orbital precision. So far as I know, there are no groups actively working on a design for this type of experiment.

A different, ground-based approach is, however, actively being developed by
Hill and Zanoni. They have constructed a very sophisticated optical system
where uncompensated differential refraction, multiple scattering, flexure, and
temperature effects are about an order of magnitude smaller than heretofore
achieved. By adding photoelectric scanning techniques, coherent detection elec-
tronic circuitry, and a laser interferometer, they hope to measure with high ac-
curacy the apparent position of stars near the Sun under normal daylight condi-
tions. The Sun's diameter will provide the basic unit of length. First, of course,
they must measure the extent of the visual deviations from circular symmetry of
the Sun. Results of measuring the gravitational deflection with errors of only 1
to 2% of the predicted effect are expected to be obtained within the next few years.
It remains to be seen whether any unanticipated effects will appreciably degrade
the expected accuracy.

The third classical test of general relativity involves the prediction that the
angular position of the perihelion of the orbit of Mercury will advance about 43"
of heliocentric arc per century more than is predicted by Newtonian theory. Such
a discrepancy between observations and predictions based on Newtonian precepts
was actually noticed decades before Einstein's birth. Leverrier in the 1850's
made a detailed study of the accumulated observations of Mercury's transits
across the solar disc. Mercury's orbit is inclined to the Earth's by about $7^\circ$;
thus only when both Mercury and Earth simultaneously lie on or near the same
"side" of the line of intersection of their orbital planes can Mercury be seen from
Earth as a silhouette against the Sun. These transits can take place in either
May or November, when Earth passes the line of intersection of the two orbits,
and actually occur about 13 times each century. Observations of a transit are
traditionally concerned solely with the determination of the times of occurrence
of four events: the instant of apparent tangency of the discs of Mercury and the
Sun with Mercury "outside" the solar disc ("first external contact"); the instant
of apparent tangency with Mercury "inside" the solar disc ("first internal con-
tact"); and the two corresponding instants as Mercury passes through the other
side of the solar limb. Such data have been obtained by astronomers at every
transit since 1677. There are two difficulties with the results. Firstly, the
clock used, until the last decade or so, was based on the Earth's rotation which

was assumed to be uniform. We now know this assumption to be false and its
effect on the data is substantial; a reliable estimate of the errors so introduced
into the data is difficult to make, even if other observations such as of those of
the Moon are considered. Secondly, and more importantly, the visual event to
be associated with the well-defined mathematical concept of apparent tangency
is by no means obvious. Through the years attempts were made by observers
of the transit phenomenon to develop unambiguous verbal descriptions that would
lead to uniformity in the data of different observers of the same transit and of
different transits. Nonetheless, Newcomb at the end of the 19th century found
no particular a posteriori correlation between the verbal descriptions and the
data. It appeared as if the verbal references to the different visual stages did
not, in fact, serve to identify a physical stage. Thus, Newcomb estimated the
interval of ambiguity surrounding a contact to be at least 10 sec, irrespective
of any detailed descriptions of the observed phenomena. (There is also the pos-
sibility that secular drifts away from data uniformity might result from improve-
ments in equipment with time. Such changes, however, would be expected to
cancel on average for the symmetrically placed contacts.)

Leverrier had access to most of the transit data gathered from 1677 to
1848. What could be done with it? Leverrier's original purpose was to study
the long-term secular perturbations of Mercury's orbit and, for example, to
estimate the mass of Venus. Of course, all of the elements of the orbits of
Mercury and Earth and their secular variations could by no means be deter-
mined from the transit data. All of the observations are made at just two dif-
ferent points of the orbits. Therefore, only certain linear combinations of the
secular variations were well determined by the data. When Leverrier found a
relatively large discrepancy here between theory and observation, he felt the
necessity to assign it to one particular element. His reasoning, which was
mostly logical and yet to a certain extent fortuitous, led to the conclusion that
Mercury's perihelion position was moving at a rate of 36" per century greater
than expected. The expected advance was approximately 5500" per century.
Somewhat over 5000" are due simply to the precession of the Earth's axis of
rotation, primarily under the action of the lunar and solar gravitational torques
exerted on the Earth's equatorial bulge. Since the coordinate systems conven-
tionally used by astronomers involve the intersection of the Earth's orbital
plane with its equatorial plane, these 5000 plus seconds of arc are convention-
ally attributed to the motion of planetary perihelia whereas they are perhaps

more properly considered as part of the motion of the astronomical coordinate system with respect to the "fixed" stars. The remainder of the expected secular advance, slightly under 500" per century, are attributable to planetary perturbations, the largest being 267"/century due to Venus and 130"/century due to Jupiter. Except for Venus, whose mass was quite uncertain, the other planetary contributions were known with more than sufficient accuracy. The precession was considered reliable to within a few seconds of arc per century, and despite the uncertainty in Venus' mass, the "observed" advance (or the equivalent combination of secular changes) seemed beyond question to represent a bona fide discrepancy with theory. Leverrier was unable to find a satisfactory explanation. Contemporary physicists for the most part did not concern themselves with the problem. The matter lay dormant for about three decades until Newcomb reanalyzed all of the transit data through 1881. He confirmed Leverrier's result, but found an even larger discrepancy: 42".95 per century in the perihelion motion. The agreement with the then unknown relativistic calculation is remarkable but not of great significance; it is clear from Newcomb's work that the uncertainty in this result is at least of the order of 10%. Nonetheless, the effect was real and had to be dealt with. It was further confirmed by Newcomb's monumental analysis in the 1890's of the meridian-circle observations of the four inner planets which disclosed about a 40" per century discrepancy for the perihelion advance of Mercury.

What explanations were put forward to resolve the discrepancy? A number were tried -- a gross error in the then accepted mass of Venus, an as yet undiscovered planet inside the orbit of Mercury, the zodiacal dust cloud, and a far larger than expected solar oblateness -- but all proved to have "side effects" which rendered them inconsistent with other astronomical data. Some people even suggested that the inverse square law of gravitational interaction was not quite exact and several ad hoc modifications were proposed. With the 1916 publication of general relativity the matter seemed elegantly resolved. Of course, new impetus was thereby given to checking the observational facts and a number of astronomers tried to reproduce Newcomb's results. Because of misinterpretations and errors of various sorts, the conclusions reached for the "anomalous"

advance of the perihelion varied from about 30" to 50" per century. The last,
and most comprehensive, treatment since Newcomb's was carried out by
Clemence who in 1943 concluded that the difference between observation and
Newtonian theory was 43".0/100 yr. He judged the probable error in the result
to be 1"/100 yr, with the major contributor being the effect of the uncertainty
in Venus' mass on the Newtonian prediction of the advance. Subsequent work by
Duncombe in the 1950's on Venus' motion and, more recently, the results of the
Mariner-Venus flyby, have essentially eliminated this source of uncertainty,
yielding an estimated error of 0".4/100 yr or, equivalently, of 1% of the predicted
additional advance.

The story does not end here. A few years ago Dicke revived the solar oblate-
ness idea with a slight twist. He reasoned that although the entire 43"/100 yr
non-Newtonian advance could not be satisfactorily explained by a solar gravitational
oblateness, 10 to 15% could be so explained without introducing any concomitant
disagreements with other orbital data. But the known mass, size, and rotation rate
of the Sun are apparently inconsistent with such a relatively large oblateness (dif-
ference between equatorial and polar moments of inertia of about 1 part in $10^5$).
To circumvent this obstacle, Dicke proposed that all is not what it seems and that
the exterior of the Sun is effectively decoupled rotationally from the interior which
is rotating far more rapidly with a period of between 1 and 2 days. (Such an hypo-
thesis, if valid, has the rather pleasing implication that the apparent great dispar-
ity between the distributions of angular momentum among planets and Sun is not
really so great!) What led Dicke to make this proposal? He, and others before
him, felt that general relativity did not incorporate Mach's principle in a funda-
mental enough manner and proposed a modification in which a scalar field was
added to the tensor gravitational field, the former providing a direct coupling with
the "distant" matter. This theory, developed by Brans and Dicke, is quite simi-
lar to an earlier one of Jordan and contains an arbitrary parameter $\underline{s}$ that repre-
sents, in essence, the fraction of the total gravitational field to be ascribed to the
scalar interaction. From other arguments of a cosmological nature, Dicke con-
cluded that $\underline{s}$ would be of the order of 5 to 10%. Since the Brans-Dicke theory pre-
dicts a non-Newtonian perihelion advance differing from that of general relativity
by a factor of $(1 - \frac{4}{3} s)$, and since the observations seemed in agreement with gen-
eral relativity, Dicke proposed that $(4/3)s$ of the advance might really be due to an

unexpectedly large solar gravitational quadrupole moment. Clearly an inde-
pendent measure of this quantity was needed. As implied above, the available
orbital data were not precise enough to determine the solar quadrupole mom-
ent. Although its secular effects fall off with the inverse seven-halves power
of the orbital major axis and the general relativistic contribution with the in-
verse five-halves power, the influences on the other planets are in either case
too small to be reliably detected, no less distinguished. On the other hand,
the visual oblateness of the Sun might be accurately measurable. If, further,
conditions were such that the visible surface was essentially in hydrostatic
equilibrium, then the visible surface would coincide with an equipotential sur-
face. By separating out the centrifugal contribution, a simple relation would
be obtained between the gravitational and visual oblateness. (This argument,
resurrected by Dicke, goes back at least as far as Newcomb.)

Dicke and Goldenberg, with essential design help from Hill, recently car-
ried out a very sophisticated measurement of the visual oblateness. An opti-
cal system was carefully developed to create an image of the Sun free from
systematic distortion of its shape. Then, between the image and a photodetector,
they placed a spinning opaque disc with diameter slightly larger than that of
the solar image. The disc contained two "notches" at opposite ends of a diam-
eter with the angular extent of the two notches being different but their radial
depth (inner diameter) the same. The purpose of the inequality in angular ex-
tent was to detect any offset of the center of the disc from the center of the Sun.
A feedback loop from the photodetector, which recorded the solar light passing
through both notches, served to change the disc position to eliminate the first
harmonic of the signal (i. e. , to eliminate the center offset). Since the orienta-
tion of the rapidly spinning disc was known accurately as a function of time, the
difference in the photocurrent for any two selected orientations could be moni-
tored. Two pairs were actually used: the $45^\circ$ diagonals and the horizontal-
vertical directions. From such measurements both the orientation and the ob-
lateness of the (assumed) elliptical cross section of the solar image can be de-
termined. To distinguish different brightnesses from a difference in shape,
different magnifications of the solar image relative to the notched disc were
employed. (The assumption here is that the radial spacing of contours of equal
brightness will, at least near the surface, be independent of the orientation of
the radial line.) A crucial test of the validity of the results is to determine
whether the measured orientation with respect to the local vertical, of, say, the
minor axis of the visual solar disc, rotates with time in the same manner as
the geometry of the situation dictates for the Sun's polar axis. The Princeton

data passed this test and yielded a difference of equatorial and polar diam-
eters of 0."096 + 0."013. Given the assumption of hydrostatic equilibrium, this
difference translates into a gravitational oblateness sufficient to account for
about 8% of the non-Newtonian advance of Mercury's perihelion.

We are thus faced with a startling possibility: If the Princeton measure-
ments are correct, if the interpretation in terms of a gravitational oblateness
of the Sun is sound, and if, finally, the determination of the non-Newtonian
perihelion advance of Mercury is as accurate as advertised, then general rela-
tivity is not an acceptable theory of gravitation. There are no arbitrary par-
ameters in general relativity and hence none that could be altered to save the
theory. One would be forced to conclude that the 1916 prediction for the ad-
vance was right for the wrong reasons. Let us examine each of the "ifs" in
turn.

Although this measurement of the visual oblateness of the Sun was carried
out with extreme care, it should be realized that the observations extended
over only a few months and were all made through only one filter (red). Fur-
ther, the results obtained from the horizontal-vertical directions were affected
by systematic errors which are not yet fully understood. (These data were eli-
minated from the final reduction.) The presence of unsuspected syste-
matic errors in the diagonal component data cannot be considered beyond doubt.
The extreme delicacy of the experiment, coupled with the relatively limited ex-
perience gained so far with this technique, imply that further measurements
certainly should be made to verify the Princeton results.

One might well enquire at this point about previous measurements of the
Sun's visual oblateness. Were these consistent with a 0."1 difference between
equatorial and polar diameters? As far as I am aware, no earlier measure-
ments yielded a significant difference. The highest claimed accuracy
accompanied a sustained series of observations carried out in Germany by Schur
and Ambronn at the end of the last century. From heliometer measurements ex-
tending over about an 11 year solar cycle, their independent results for the equa-
torial minus polar diameter were: 0."007 ± 0."015 (Schur) and - 0."002 ± 0."009 (Am-
bronn). The basis of this type of measurement is an objective, separated into
two equal parts with each producing its own image. The two halves can be moved

so as to bring the two images into tangency, thus allowing a given diameter to be measured accurately. Again, although the German observers were very meticulous and had long experience with their instrument, the observations depended directly on personal judgement and the possibility of there being systematic errors substantially larger than the quoted ones is difficult to eliminate.

The interpretation of the visual oblateness in terms of an oblateness in the gravitational potential is not nearly so simple as might at first be thought. Very little is really known about the physical characteristics of the photosphere and interpretations may be strongly model dependent. Several incompatible explanations have already been published concerning this conversion of a 0!'1 visual oblateness into a gravitational one. Moreover, the Dicke interpretation seems to require the Sun to possess a rapidly rotating inner part. That the outer part could be effectively decoupled is by no means clear. Of course, even less is firmly known about the interior of the Sun than about its surface and so proposed models are not too greatly restrained by observations. A number of possibilities have been examined by fluid dynamicists who conclude that a large differential rotation between the outer and inner parts of the Sun could not be maintained. Although no support has yet emerged for Dicke's proposal, all agree that the case against his two-part model can by no means be considered beyond doubt. Most agree, as well, that without the rapidly spinning inner part, the sun's gravitational oblateness could not contribute appreciably to the advance of Mercury's perihelion.

Finally, we return to the question of the reliability of the determination of the non-Newtonian perihelion advance. The generally accepted value, as mentioned, is 43!'0 ±0!'4 per century. But a difference of 0!'4 in the position of perihelion corresponds to a maximum difference in Mercury's position (given the same orbital period) of only about 0!'08 since its orbital eccentricity is approximately 0.2. In other words, if one were to consider two planets initially lying in a straight line with the Sun and travelling along identical orbits, except for a difference of 0!'4 in their arguments of perihelion, the maximum difference in heliocentric angular position of the two planets would be only about 0!'08 for e = 0.2. For Mercury, viewed from the Earth, the maximum angular difference would be reduced from 0!'08 to at most 0!'06 in virtue of the Earth-Mercury separation at inferior conjunction being four-thirds the average distance of Mercury from the Sun. (Of course except for transits, Mercury is never observed optically near inferior conjunction. )

On the other hand, individual series of meridian-circle observations of Mercury show spreads in residuals of between ± 1" and ± 2" of arc and observations of individual transits of Mercury yield contact times usually varying over more than 30 seconds (equivalent to a spread greater than about 4" in geocentric arc). I find it quite conceivable that the systematic errors present in these observations, or in their reductions, are large enough to prevent meaningful deductions from being made about centennial variations in observed angular positions of the order of 0.''06. For example, the Newcomb constant of general precession -- a key element in the comparison of the advance of Mercury's perihelion with predictions -- is widely recognized to be in error by about 1" of arc, which is alone two and a half times the accepted error in the perihelion motion comparison. Because of these and other doubts about the limits of reliability of the optical observations, my colleagues Ash and Smith and I have undertaken a re-analysis of these data. The project is a tedious one involving, for example, the conversion to machine-readable form of over 100,000 optical observations of the planets. We hope to obtain at least preliminary results within a year.

Having completed our discussion of the logical chain attached to the Princeton measurements, we may enquire how radar can be used to help resolve the controversy. Ideally, we wish to distinguish between the effects on planetary orbits that might be attributable to general relativity and those that would be caused by a solar gravitational quadrupole moment. The secular motion of the perihelion induced by the former varies as the inverse five-halves power of the planet's semimajor axis whereas the corresponding motion induced by the latter decreases with the inverse seven-halves power of $a$. The difficulty in basing a discrimination on this difference is clear: Although the fractional predicted differences are large when comparing results for planet's with widely differing $a$'s, the absolute effect decreases very rapidly with increasing $a$ in both cases. Nevertheless, if radar measurements of Earth-Mercury, Earth-Venus, and Earth-Mars time-delays are continued for about a decade with each having an uncertainty no greater than 10 μsec, it should be possible either to determine the Sun's gravitational quadrupole moment or to place a useful upper bound on it. Preliminary calculations indicate that a solar quadrupole moment sufficient to cause a 2"/100 yr advance in Mercury's perihelion could be detected with these measurements. The effects on inclination and ascending node (measured with respect to the ecliptic)

of the quadrupole moment are very difficult to detect with time-delay measure-
ments since first-order changes in these elements introduce only second-order
changes in the delays: the corresponding changes in orbital position occur in a
direction normal to the line-of-sight.

Artificial planets with radio transponders could conceivably allow far
greater precision to be obtained in estimating the quadrupole moment. Aside
from the economic and political difficulties of such an enterprise, one must be
careful to insure that the desired results are not vitiated by substantial and un-
predictable effects of gas leakage from attached rocket motors and of sunlight
pressure. The latter will introduce complications to the extent that the solar
constant changes and that the orientation and reflection properties of the artifi-
cial planet are unknown. In principle these difficulties can be eliminated by en-
closing an inert "planet" with a shield that senses the position of the planet elec-
tromagnetically and, by means of rockets, adjusts its motion to follow the purely
gravitational behavior of its partner. The pressure of the sensing radiation can
obviously be reduced to negligible proportions. In the end, the limit on accuracy
may be set by the "noise" introduced via the gravitational perturbations of the
myriad of uncharted asteroids. In any event, the engineering realization of a
shielded artificial planet belongs to the distant future.

The thoughtful reader might at this point begin to mull over a simple ques-
tion: Since interplanetary radio and radar observations make no reference to
the "fixed" stars, how can a change in perihelion position be detected? The cor-
respondingly simple answer (although not complete) is that the perihelion motion
of each observed planet can be determined with respect to the orbital position of
Earth. We have, in essence, a closed dynamical system and through interplan-
etary time-delay measurements we can monitor the relative spin and orbital mo-
tions of each of the components. The overall orientation of the system with re-
spect to the "fixed" stars is indeterminate and irrelevant. It is true, of course,
that at present by no means all of the planets are accessible to radar observa-
tions. But the future looks bright.

To summarize, we reiterate that radar and radio observations offer the pro-
mise -- with fulfillment expected within a decade -- of a resolution to an accepta-
ble accuracy of the controversy raised by the specter of a possibly large solar
gravitational oblateness.

C. Increase in Interplanetary Time Delays

Several years ago it became evident that another experimental test of general relativity was technically feasible. The experiment, which I suggested, was designed to verify the prediction that the speed of propagation of a light ray decreases as the ray passes through a region of increasing gravitational potential. The test could be performed, for example, by measuring the round-trip time delays of radar signals reflected by Venus or Mercury as either passes on the other side of the Sun from the Earth -- the superior conjunction alignment. The slowing down of the propagation speed by solar gravity would be manifested by an increase in the round-trip delay as the ray paths pass closer to the solar limb. These increases, of course, would be superposed on the expected delay attributable to the separation between the Earth and the target planet. How may we estimate this delay effect quantitatively? For the present illustrative purpose we may consider the two planets to remain stationary during the round-trip travel time of the radar signal [$(v/c) \approx 10^{-4}$]. Of course, making such an approximation in an actual experiment would be disastrous. Under this approximation, the round-trip coordinate-time delay is given by

$$t = 2 \int_{\vec{r}_e}^{\vec{r}_p} \frac{d \, path}{v} \tag{51}$$

where $\vec{r}_e$ and $\vec{r}_p$ are the positions of the Earth and target planet, respectively; and where $v$ is the speed of propagation of the light or radar signal. The integration is to be carried out over the actual (curved) spatial path as shown in Figure 4. We can, however, replace this path by a rectilinear one without introducing any substantial error. This replacement has two effects: (1) the path length is decreased, and (2) the gravitational potential is greater along the new path. Both consequences are easily shown to be of second order in $r_o$ and so are negligible. For example, we note from Figure 4 that the difference in path lengths is proportional to $1 - \cos \eta$, i.e. to $\eta^2$, where $\eta$ itself (the deflection) is proportional to $r_o$. Application of Fermat's principle establishes the same result. To first order in $r_o$, we may then replace Eq. (51) by

$$t = 2 \int_{x_e}^{x_p} \frac{dx}{v} \quad , \tag{52}$$

where

$$v = \frac{dx}{dt} = c \left[ 1 - \frac{r_o}{r} \left( 1 + \frac{x^2}{r^2} \right) + 0 \left( \frac{r_o^2}{r^2} \right) \right] \quad . \tag{53}$$

Eq. (53) follows directly from Figure 4 and Eq. (46) since for light rays $ds^2 = 0$. Here $\underline{c}$ (previously set equal to unity) represents the speed of light propagation far from the Sun. No loss of generality is involved by taking the x-axis parallel to the rectilinear path between the Earth and the planet.

Integrating Eq. (52) yields

$$t = \frac{2}{c} (x_e - x_p) + \frac{2r_o}{c} \left[ 2 \ln \left( \frac{x_e + r_e}{x_p + r_p} \right) - \left( \frac{x_e}{r_e} - \frac{x_p}{r_p} \right) \right] . \tag{54}$$

What we seek, though, is not the coordinate-time delay but rather the expression for the proper-time delay $\tau$ which is the time measured by the Earth observer using, for example, an atomic clock. The interval of proper time is given by

$$\tau = \frac{1}{c} \int_{s_1}^{s_2} ds \tag{55}$$

where $s_1$ is the (four-dimensional) position of the Earth at transmission and $s_2$ is the corresponding position at echo reception. For the Earth $ds^2$ is given by Eq. (45) and since we assumed that the Earth remains stationary in space between transmission and reception of the radar signal, we find to first order in $r_o$:

$$\tau = \frac{2}{c} (x_e - x_p) + \frac{2r_o}{c} \left\{ 2 \ln \left( \frac{x_e + r_e}{x_p + r_p} \right) - \left[ \frac{2x_e - x_p}{r_e} - \frac{x_p}{r_p} \right] \right\} \quad .$$

$$\equiv \frac{2}{c} (x_e - x_p) + \Delta \tau \quad . \quad , \tag{56}$$

where $\Delta\tau$ represents the "excess" delay attributable to the slowing down of the propagation speed by solar gravity. Near superior conjunction, where the impact parameter $\underline{b}$ of the ray (see Figure 4) is small, we find

$$\Delta\tau \approx \frac{4r_o}{c}\left[\ln\left(\frac{4r_e r_p}{b^2}\right) - \frac{3r_e + r_p}{2r_e}\right]; \quad b << r_e, r_p, \tag{57}$$

whereas near inferior conjunction

$$\Delta\tau \approx \frac{4r_o}{c}\left[\ln\left(\frac{r_e}{r_p}\right) - \left(\frac{r_e - r_p}{2r_e}\right)\right]; \quad b << r_e, r_p. \tag{58}$$

The numerical behavior of $\Delta\tau$ as a function of the Sun-Earth-planet angle is shown in Figures 5 and 6 for Mercury and Venus, respectively. We see that the maximum "excess" delay is about 200 μsec out of a total round-trip time of nearly $2 \times 10^3$ sec--a maximum fractional increment of less than 1 part in $10^7$. In effect, this predicted extra delay would appear as a "bump" or localized swelling in the orbit of the target planet (or in the Earth's!) each time superior conjunction is passed.

What difficulties can be expected in performing this experiment? First of all, we must know what to expect in the absence of the extra delay; in other words we must know the planetary orbits. Although, to be sure, centuries of optical observations have enabled these orbits to be determined to $\underline{n}$ significant figures of accuracy, for our purposes $\underline{n}$ is not big enough. Even if it were, we would still be faced with an important question: What coordinates of general relativity correspond to the Newtonian coordinates? The numerical positions of the planets were obtained by reducing the observations in accord with Newtonian theory. In general relativity, there exists a rich choice of coordinate systems (standard Schwarzschild, isotropic, harmonic, etc.) for which the coordinate values of a given observable would in general be different. Into which of these different systems should we take over the Newtonian coordinate values? The answer is simply that the question is the wrong one to ask. Coordinate values themselves have no operational significance in general relativity. The only philosophically defensible procedure is simply to take over the observations and reprocess them consistently in terms of the theoretical framework

being tested which in this case is general relativity. In other words, the orbital elements and the other unknown parameters of the dynamical system are theory-dependent as well as measurement-dependent numbers. This time-delay test of general relativity can therefore be viewed as a gigantic exercise in the statistical theory of parameter estimation or hypothesis testing. We may take all of the relevant observations and from them determine the maximum likelihood estimates of the unknown parameters. If the subsequent comparison of the observations with the corresponding theoretical predictions (based on the previously determined parameter estimates) shows agreement to within the measurement errors, we can conclude that these data are consistent with the theory. Otherwise the theory may have to be discarded, at least in part. Clearly the data must redundantly span the parameter space or the test will be manifestly shallow. Instead of making such a massive self-consistency check, practical considerations allow us to effectively separate out a test of the predicted "excess" time delays. The parameters associated with the orbital motions can be determined solely from observations made far from superior conjunctions. The time-delay observations near superior conjunction can then be compared directly against the theoretical predictions. Of course, "near" and "far" are relative terms and therein lies the arbitrariness of this procedure.

What other problems interfere with the simple-minded performance of this experiment? A moment's reflection reveals that the plasma constituting the solar corona will also slow the propagation of radio waves. How can this effect on the time delays be distinguished? There seem to be at least three options open to the experimenter. All stem from the fact that $\Delta \tau_{p\ell}$, the delay increment attributable to plasma, varies with the inverse square of the frequency of the electromagnetic radiation employed:

$$\Delta \tau_{p\ell} \sim \frac{r}{f^2} \tag{59}$$

First, since general relativity is a dispersionless theory -- all radiation is predicted to have the same speed of propagation regardless of its frequency -- one can at least in principle separate out the plasma effect by measuring simultaneously both the group and the phase delay of the radar signals. Second, measurements of group delay simultaneously at two frequencies suffice to determine the plasma induced delay increment. Third, use of a single, sufficiently high frequency radar signal can ensure that $\Delta \tau_{p\ell}$ remains smaller than the measurement

uncertainty. For practical reasons only, this last approach is the one being used. By employing a frequency of $10^{10}$ Hz, for example, the plasma increment to the delay can be expected to remain less than 1% of the predicted gravitational increment for impact parameters greater than about 3 solar radii.

In relation to the second method mentioned for circumventing plasma problems, we note that the rate of change of phase delay (i.e., the Doppler shift) will have a contribution from the general relativistic effect on time delay. Since this delay increment changes with time because of the relative motion of planets and Sun, there will be a corresponding increment $\Delta f_{gr}$ to the Doppler shift [see Eq. (20)]. Near superior conjunction, this increment can be approximated by

$$\Delta f_{gr} \approx \mp \frac{8r_o f}{c} \left| \frac{r_e v_p - r_p v_e}{(r_e + r_p) b} \right| \tag{60}$$

where

$$v_{e,p} \equiv \left| \frac{d\vec{r}_{e,p}}{dt} \right| . \tag{61}$$

The upper sign holds before superior conjunction when the delay increment is increasing; the lower signal applies after the conjunction when the increment is decreasing. The magnitude of $\Delta f_{gr}$ decreases rapidly with an increase in impact parameter so that for $b$ equal to three solar radii, $\Delta f_{gr}$ is only about 4 Hz for Mercury and 1 Hz for Venus at $f \approx 10^{10}$ Hz. These effects are too small to be reliably detected at present.

Finally, we come to the important question: What are the experimental results from the time-delay test? Unfortunately, none of significance yet. But considerable progress has already been made in refining the necessary estimates of the orbits of the inner planets and associated parameters, such as planetary masses and radii. We will describe briefly the theoretical procedures used and then the results so far obtained by Ash, Shapiro, and Smith. Our choice of units is for the most part in accordance with standard astronomical practice. We set the mass of the Sun equal to unity and take the day as the unit of time where 1 day, by our definition, equals 86,400 seconds of atomic time. This second, in turn, is defined by a large, specified number of oscillations of the electromagnetic radiation emitted by the Cesium atom in making a particular transition

under specified conditions. In particular our time measurements are synchronized with the atomic time maintained at the U.S. Naval Observatory. The unit of length, by definition one astronomical unit (a.u.), is determined by giving the gaussian constant $\underline{k}$ its standard value: 0.01720209895.

The parameters characterizing the rotation of the Earth (precession and nutation matrices) are taken from standard astronomical sources with universal time, which measures the angular orientation of the Earth, being obtained from U.S. Naval Observatory bulletins. The theoretical positions and velocities of the observed planets, in the chosen harmonic coordinates, are obtained by numerical integration with the concomitant errors kept below 1 part in $10^{10}$ over the intervals of interest. A variable step size, predictor-corrector algorithm is used with initial estimates, obtained from prior work, for the initial conditions and necessary parameters. The theoretical values for the radar (time-delay and Doppler-shift) and optical (meridian-circle) observations are then calculated and compared with the observations. New, maximum likelihood estimates are obtained for all the parameters in the usual iterative manner. The estimated standard errors and the correlation, or moment, matrix of the parameter estimates are also calculated. All computations are, of course, carried out with a high-speed digital computer.

Using U.S. Naval Observatory optical data of the Sun, Mercury, and Venus, spanning the period from 1950 to 1965, and all of the planetary radar data obtained at Lincoln Laboratory and at the Arecibo Ionospheric Observatory, through mid-1966, we have obtained the results illustrated in Figures 7 through 12 and Tables 1 through 4. The first of these figures shows the time dependence of the osculating elements of Mercury's orbit and illustrates the perturbations attributable to the other planets. One can see, for example, that the fluctuations in the osculating semi-major axis are as great as $2 \times 10^{-5}$ a.u. or about 3000 km which is some 100 times larger than the maximum predicted effective increase in path due to the general relativistic decrease in speed of propagation of the radar signals. Figures 8 to 10 exhibit samples of the post-fit residuals obtained for the optical data. The spread is seen to be generally on the order of 1" of arc. In the Venus residuals a systematic error is clearly in evidence. From the correlation of the sign of the residual with the illuminated part of the planet, we may infer that the corrections made for planetary phase were probably inaccurate.

Figures 11 and 12 show samples of the time-delay residuals for Earth-
Mercury and Earth-Venus observations, respectively. For the Earth-Mercury
measurements, the observed minus computed (O-C) residuals are shown for
three sets of calculations. In the first, the computations were based on the
ephemerides produced at the Jet Propulsion Laboratory to match Newcomb's
orbits for the inner planets. (The values for the planetary radii and a.u. were,
however, taken from the solutions based on radar data.) The residuals have
large systematic oscillations with amplitudes up to 4 msec. The lower two sets
of residuals show the comparisons when the orbits are fitted to the optical and
radar data. The two sets differ in that the calculations leading to the upper
residuals were carried out in conformity with the theory of general relativity
whereas the lower residuals were based on solely Newtonian calculations. The
superiority of the former is especially apparent near perihelion and is presum-
ably attributable to the additional advance of the perihelion predicted by general
relativity. From Figure 12, one can see that for the Earth-Venus measure-
ments, the orbits fitted to the radar data yielded a solution far superior to that
obtained from Newcomb's orbits. (For the Earth and Venus, relativistic effects
are too small to be observable and the Newtonian and general relativity fits are
virtually indistinguishable.)

A total of 26 parameters were adjusted in making these weighted-least-
squares fits to the radar and optical data: 18 initial conditions for the orbits of
Mercury, Venus, and the Earth-Moon barycenter; the masses of Mercury,
Venus, the Earth, the Moon, and Mars; the equatorial radii of Mercury and
Venus (assumed independent of longitude); and the a.u. In Table 1 the osculat-
ing orbital elements are given for both fits with the concomitant formal stand-
ard errors being shown in Table 2. The values for the planetary masses and
for the planetary radii and astronomical unit are contained, respectively, in
Tables 3 and 4. The radius of Venus that we obtained is seen to be substantially
smaller than previous estimates; the other parameter determinations are in
reasonable agreement with the results of other analyses. In regard to the orbital
initial conditions, the power of the radar data is shown dramatically in Table 2.
Although the time span of the radar data is several fold less than that of the op-
tical observations used, the results obtainable from the radar data for eccentricity,

semi-major axis, and initial mean anomaly are at least an order of magnitude
more precise. Of course, the radar data are relatively insensitive to changes
normal to the line-of-sight; hence the values found for $i$ and $\Omega$ are far less pre-
cise than those deduced from the optical measurements alone. Combination of
both types of data yields a significant reduction in the uncertainty of the estim-
ate of $\omega$.

Having determined the parameters for a given theory from the observations,
we may of course predict the results of future observations. In addition, we may
estimate the expected accuracy of these predictions in a straightforward manner
from the formal errors in the parameter estimates. (The basic measurement
errors were assumed to be independent and gaussianly distributed with zero
means.) We applied this technique to radar data taken after mid-1966 and found
the actual prediction errors to be quasi-periodic with an amplitude approximately
five times greater than expected. Presumably the cause of these discrepancies
lies in as yet undiscovered systematic errors -- either in the measurements, the
physical model, or the computer program! For this reason especially one should
only be surprised if future analyses do not yield any parameter estimates differ-
ing from ours by more than the standard errors.

We return now to the discussion of the time-delay test proper. The radar
data analyzed above represented observations through the fall of 1966. We have
not been idle since. The first attempt at performing the test was made as Venus
passed through superior conjunction in November 1966. No superior conjunction
observations had been made earlier because the radar system at Haystack had to
be improved in order for useful data to be obtained at superior conjunction. The
transmitter power was increased from 100 kw to about 250 kw (the designed im-
provement envisioned a 500 kw klystron transmitter system but that figure has not
yet been achieved). The receiver system was also redesigned to yield an overall
system noise temperature of about $55^{\circ}K$ as opposed to the earlier value of $140^{\circ}K$:
The maser receiver (a helium-cooled ruby enclosed by a 4500-gauss supercon-
ducting magnet) contributes $8^{\circ}K$; the hardware losses from the internal flange of
the maser to the horn of the antenna contribute $20^{\circ}K$, the horn itself $7^{\circ}K$, and the
radome and atmosphere about $20^{\circ}K$. (Heavy clouds in summer can contribute an-
other 10 to $20^{\circ}K$ and, when the radome is wet, the system noise temperature can
rise to about $100^{\circ}K$.)

Only a few, inconclusive data points were obtained from the Venus observations in November. The next opportunity, the superior conjunction of Mercury in January 1967, was less favorable since Mercury was low in the sky (thus curtailing the daily observing time) and did not pass closer to the Sun than two degrees. Again the results were somewhat on the meager side. A massive effort, however, was mounted for the favorable May 1967 superior conjunction of Mercury. (The planet actually passed behind the Sun on the 11th.) The crew of engineers operating Haystack during this period started preparations each morning at 4:30 a.m. and continued with the observations and data reduction until 7 p.m. each evening. This schedule was maintained almost every day from 1 May to 23 May. The most important contributions in designing the radar system and directing its operation were made by Pettengill, Stone, Price, Ingalls, and Brockelman of Lincoln Laboratory. We expect to be able to publish the results of these observations in the near future. Because of the rapid variation with time of the predicted "excess" delays, this test should not be seriously affected by the possible presence of (slowly-varying) systematic errors such as might be responsible for the prediction errors described earlier.

Before leaving the subject of experimental general relativity, I wish to describe briefly one more of our results and another experimental test of general relativity that is actively being developed at Stanford University.

D.  Possible Change in Gravitational Constant

In addition to estimating the planetary orbital elements, masses, and radii from the radar and optical data, we also introduced another parameter to represent a possible secular change of the gravitational constant. There have been a number of theoretical speculations since Dirac's first such conjecture in 1937 that the gravitational constant G might in fact be decreasing with time. But even apart from theoretical arguments, it is of interest to test for a possible time dependence of G since such a dependence could have extremely important consequences for the history of the Earth, the evolution of stars, and for cosmological models in general. We therefore introduced an ad hoc parameterized form for G:

$$G = G_o - \dot{G}_o(t - t_o) \tag{62}$$

and estimated the constant $G_o$, as well as its uncertainty, from the data. Our preliminary results indicate that $\dot{G}_o$ is less than $10^{-9} G_o$ per year. That is, if G is time dependent, it does not change by more than one part in $10^9$ per year. With the accumulation of three more years of Earth-Venus and Earth-Mercury radar data at the currently achievable level of precision, we should be able to detect reliably any change in G larger than about 5 parts in $10^{11}$ per year. Since the main consequence of $\dot{G}_o$ is to produce a change in the relative longitudes of the planets, the effect will increase approximately with the square of the observing interval.

### E. Gyroscope Precession

The so-called gyroscope test of general relativity involves the predicted behavior of an angular momentum vector in a gravitational field. According to Newtonian theory, a perfectly spherical spinning gyroscope freely falling in a gravitational field will continue to point in the same direction in inertial space, i.e. with respect to the "fixed" stars. But, according to general relativity, a freely falling angular momentum vector (gyroscope) can precess with respect to the fixed stars. This consequence of relativity theory was recognized in 1916 by de Sitter who calculated the rate of precession of the Earth's spin axis caused by its being in orbit about the Sun. This effect, known as the geodesic precession, amounts to slightly less than 2" of arc per century for the Earth. Somewhat later, in 1918, Lense and Thirring calculated the additional effect due to the spinning of the Sun. These theoretical results have remained unverified for 50 years. But the outlook for an experimental test is no longer so dim. Soon after the launching of the first artificial satellite, Schiff suggested that these two predicted effects on the orientation of an angular momentum vector might be detectable using a gyroscope in orbit about the Earth. The calculations show that if the angular momentum vector lies in the orbital plane, this vector will precess in the plane at a rate of about 7"/yr when in an 800-km altitude orbit. The contribution to the gyroscope's precession from the Earth's spin is about two orders of magnitude smaller. To separate the two effects, a polar orbit can be used. With the gyroscope's spin axis normal to the orbital plane, the precession will be due solely to the Earth's rotation and is predicted to be about 0".05/yr. Thus, with two gyroscopes inside a polar orbiting satellite -- one with its spin axis in the orbital plane, the

other with the axis perpendicular to the plane -- the two effects can be separated and studied simultaneously.

How can such small precessions be detected reliably in the presence of the presumably far larger effects of the ordinary torques acting on an orbiting gyroscope? Obviously this is not an easy experimental task. However, the Stanford experimenters, principally Fairbank and Everitt, are undaunted and are developing a very sophisticated apparatus with which they hope to measure these relativistic precessions with an error of only 0.001"/yr. The essence of their approach involves the use of a superconducting thin film to coat a nearly perfect (quartz) sphere. The problem of determining the direction of the angular momentum vector of a perfectly spherical gyroscope is non-trivial: Suppose we were to create such a gyroscope and put a small scratch on it from which to determine the direction of the spin axis. This scratch, to be visible, would likely create sufficient asymmetry for the resultant gravitational torque to cause a precession larger than the relativistic effect being sought! Calculations show that the gyroscope must be spherical to about 1 part in $10^6$ and homogeneous in density to 1 part in $10^5$ to insure that $(\Delta I/I) < 10^{-5}$, where $\underline{I}$ is the moment of inertia and $\Delta I$ the maximum difference in the principal moments.

The ingenious use of superconductors, suggested by Fairbank, to surmount this problem of the observer influencing the observed, can be described briefly as follows: a spherical quartz ball is coated uniformly with a thin superconducting film, spun up by gas jets, and electrostatically suspended by three mutually perpendicular electrostatic fields, the whole apparatus being maintained at liquid helium temperatures. A magnetic moment will be created in the spinning superconducting film (London moment) caused in essence by the drift of electrons with respect to the lattice ions in the superconductor. This magnetic moment is very precisely aligned with the angular momentum vector; even if the body axis drifts with respect to the spin axis by one radian, no harm is done. Measuring accurately the direction of the London moment is also non-trivial. Superconducting circuitry and parametric amplifiers have been specially invented for this purpose.

The maintenance of superconducting temperatures around an orbiting gyroscope for periods of the order of one year requires that enormous amounts of

boiled-off liquid helium must be disposed of during the course of time. In
fact, the amount turns out to be far greater than that needed for the satellite's
attitude control system. One is thus confronted with a problem opposite to the
conventional difficulty in designing attitude control systems: How can large
amounts of gas be emitted in a controlled manner? This problem is still under
study.

The comparison of the orientation of the gyroscope with the "fixed" stars
presents a further problem. The proposed solution envisions a rigid connec-
tion between a telescope and the gyroscope mount with its read-out system.
The telescope optics and associated electronics are designed to take maximum
advantage of the good seeing conditions by "splitting" the diffraction patterns
of the reference stars.

Although there are many other problems connected with the gyroscope test,
all are being studied and it is hoped by the experimenters that some of the com-
ponents might be tested in orbit by 1970. Actual experimental results are not
expected before the early to middle '70's.

Let us summarize the main assumptions involved in the gyroscope experi-
ment: (i) The angular momentum vector is defined by the London moment;
(ii) The reference stars used in the experiment are, with respect to their
transverse motions, "fixed" in an inertial frame; and (iii) general relativity
gives a correct description of the precession of angular momentum vectors.
If the observations agree with predictions, one has established the self-consis-
tency of the assumptions. If results disagree, there are various possible con-
clusions. For example, the reference stars may not correspond to an inertial
frame. How could we test this possibility? We could put two or more sets of
gyroscopes in orbit simultaneously at different altitudes and observe their pre-
cession with respect to the same set of stars. In this way, we could compare
the precessional motion of one set of gyroscopes with respect to another inde-
pendently of whether or not the stars form an inertial system. In other words,
the precession predicted by general relativity is function of the orbital altitude;
hence we can test the predictions simply by comparing relative precessions of
two or more sets of gyroscopes. The stars form a convenient intermediary
reference, but since the comparison of each gyroscope set with these stars is
made at the same time, any transverse motion of the stars will not "ruin" the
results. If the results are consistent with general relativity, then we would

have a method for finding or forming a stellar reference system that has the desired inertial property. Or, put in a slightly different manner, we would have a method (albeit very expensive!) for determining the proper motions of stars. If, as the experimenters hope, the gyroscope pointing can be monitored with a precision of 0.''001 for periods of more than a year, then this method for the determination of an inertial stellar reference system might eventually prove useful.

The last possible conclusion -- that general relativity is wrong -- would only be reached after an exhaustive study had ruled out all other causes of any observed discrepancies.

## VI. AXIAL ROTATIONS OF MERCURY AND VENUS

The final topic to be discussed is the rotational motion of the planets Mercury and Venus. We will first describe the radar techniques that have been used to measure the rotation, or spin, vectors of these planets and then the present status of the theories propounded to explain these anomalous axial rotations.

### A. Radar Determination of Planetary Spins

In our analysis of delay-Doppler mapping, we derived for a homogeneous planet the expression for power vs. frequency for the energy returned from a particular ring on the planet [Eq. (37)]. This equation and its predecessors show that the frequency extent of this spectrum, the bandwidth, is given by

$$B(\tau_r) = 2f_1 |\vec{\omega} \times \vec{e}_{12}| [\tau_r(\frac{4\rho}{c} - \tau_r)]^{1/2} , \qquad (62)$$

where

$$\vec{\omega} = \vec{\omega}_a + \vec{\omega}_s ,$$

with $\vec{\omega}_a$ being the apparent angular velocity attributable to the relative orbital motion of the radar site and target planet, and with $\vec{\omega}_s$ representing the intrinsic or sidereal angular velocity of the planet. Our object is to determine $\vec{\omega}_s$ from measurements of $B(\tau_r)$. Since the power density reaches a maximum at each limb of the spectrum [see Eq. (37)], the signal-to-noise ratio at these

two points is in general higher than elsewhere. Thus the bandwidth, which is simply the difference of the limb frequencies, can be determined easily and accurately from the observed spectrum corresponding to echoes from a given ring. This fact, of course, governed our selection of the bandwidth data as the basis for estimating $\vec{\omega}_s$.

Since $B(\tau_r)$ is proportional to the projection of $\vec{\omega}$ on the plane perpendicular to $\vec{e}_{12}$ (unit vector along the line-of-sight from radar to target), we will have, as time goes on, different vectors $\vec{\omega}$ (as $\vec{\omega}_a$ changes) projected on different planes (as $\vec{e}_{12}$ changes). Therefore from a time series of measurements of $B(\tau_r)$, for one or a set of values of $\tau_r$, all three scalar parameters specifying $\vec{\omega}_s$ can be estimated in a weighted-least-squares sense. The accuracy achievable by this method is limited essentially by the accuracy with which B can be determined. Continuation of these measurements over an arbitrarily long time interval will not yield arbitrarily accurate results for the average rotation rate $\omega_s$. On the other hand if we could recognize and follow a feature in the radar "picture" of the planet, then by continued observation of the movement of the feature, the average $\omega_s$ could be estimated with an ever increasing accuracy. Planets, like the Earth, are not homogeneous and have anomalous scattering regions. These latter form the basis of this feature method of determining planetary spin vectors. As an example, we show in Figure 13 a typical spectrum of the total radar echo from Venus obtained at Haystack. Several features are clearly discernible. By following the position of the feature in the spectrum, one can estimate $\vec{\omega}_s$. Of course, for a given feature there are now five unknowns: the three scalars characterizing $\vec{\omega}_s$ plus the two describing the planetary latitude and longitude of the feature. In principle, the longer the feature is followed the more accurate will be the resultant estimate of $\omega_s$. In practice, a difficulty arises. It is not always obvious that a given distortion in the spectrum on one day corresponds to that on the spectrum for another day: The appearance on a radar map of a given physical feature may depend importantly on the orientation from which it is viewed. Nonetheless this method has been used very successfully on Venus, partly because of special circumstances. At every close approach between Earth and Venus, the latter presents essentially the same aspect to the radar and hence features have the same appearance at these times when the signal-to-noise ratios in the spectra are highest.

The reader might wonder why only the spectral location of a feature has been discussed. The reason is mostly historical; the first applications of the feature method employed only spectral data. One could equally well follow the positions of features in delay; in fact, with sufficiently high signal-to-noise ratios, one can use delay-Doppler maps and follow the two-dimensional positions of features. For each additional feature, we add only two unknowns: the three associated with $\vec{\omega}_s$ are common to all.

What results have been obtained? We discuss Mercury first and begin with the very interesting history of prior optical determinations of its rotation rate. Mercury has always been notoriously difficult to observe telescopically because of its small size and proximity to the Sun. No distinguishing features were seen until the early 1800's when Schröter, a German astronomer, thought he observed mountains about 60 km high. From a series of Schröter's observations of such features, Bessel deduced a rotation period very close to 24 hours. Some were delighted with this result: Mercury and its fellow inner planets Earth and Mars had nearly identical sidereal spin periods. Many were skeptical. In the 1880's, the famous Italian astronomer, Schiaparelli, maintained a long and close watch on Mercury. He concluded that Mercury was rotating slowly and, in fact, that its spin period was 88 days -- exactly equal to its orbital period. Many were aesthetically delighted: Mercury like the Moon always presented the same face to its primary. Some remained skeptical. Nonetheless, every observer after Schiaparelli -- without exception -- supported his conclusion. Less than a decade ago, a leading specialist in planetary observations claimed that the accumulated series of drawings and photographs of Mercury allowed one to conclude that the spin and orbital periods matched to better than 1 part in $10^4$. Only when it became possible to apply the completely objective radar bandwidth method to determine the spin vector was this myth finally exploded. From a series of bandwidth measurements made at Arecibo during two successive inferior conjunctions in 1965 (see Figures 14 and 15), Dyce, Pettengill and I were able to show that Mercury's sidereal spin period was only $59 \pm 3$ days with an axis inclined by less than $25^\circ$ to its orbital plane. This represented only a slight improvement over the original determination of $59 \pm 5$ days made by Pettengill and Dyce on the basis of data from the first conjunction only. Since the direction of Mercury's spin is prograde, one can easily show that

the average "day" on Mercury is not eternally long, but rather 176 ± 9 days. Because of the large orbital eccentricity the day is rather peculiar: the direction of motion of the Sun across the sky as seen from Mercury goes through a double reversal near perihelion.

The myths about Venus' rotation were not nearly so widely held. Because of its thick cloud cover, which is essentially opaque to visible radiation, few observers claimed to have seen the presumed solid surface. However, observations of Doppler shifts in the spectral lines of atmospheric components (mainly $CO_2$), led to the general belief that the rotation was slow, perhaps synchronous like the Moon and, as was thought, Mercury. There were nonetheless recent claims by Kuiper that the rotation was direct with a period of about 15 days. The first radar observations of Venus at Lincoln Laboratory and at JPL in 1961 established that the rotation was slow with a period of about 200 days. The direction of the rotation was not definitely established although Carpenter and Smith independently pointed out the possibility of retrograde rotation. In 1962, with greater available radar system sensitivity, the rotation direction was definitely established by Carpenter and Goldstein as being backwards with the rotation period being about 250 ± 40 days. Thus, on Venus the Sun rises in the West and sets in the East. In the past few years the spin vector has been further refined. A value of 244 ± 2 days was obtained from an analysis of the 1964 Arecibo bandwidth data and is illustrated in Figure 16. The curve of bandwidth vs. time is seen to be a minimum at inferior conjunction. Why? Since $\vec{\omega}_a$ [see Eq. (63)] is approximately inversely proportional to the interplanetary distance, the contribution of $\vec{\omega}_a$ to $\vec{\omega}$ is a maximum at the close approach of the two planets. But, because of the retrograde sidereal rotation, $\vec{\omega}_a$ opposes $\vec{\omega}_s$ at inferior conjunction which leads to the minimum in bandwidth.

The most recent solution I have obtained by combining bandwidth and feature data shows the period to be 243.09 ± 0.18 days and the axis direction to be within a few degrees of normal to Venus' plane. The solar day on Venus is therefore 117 terrestrial days, i.e. less than its sidereal period because of the retrograde direction of rotation. The significance of the sidereal value will become obvious later.

B.  Theoretical Analysis of Spin-Orbit Resonances

    The reader may have wondered how it was possible for so many ob-
servers of Mercury to have failed throughout 85 years to notice the difference
between an 88-day period and a 59-day period.  Since, as Colombo was the
first to realize, 59 is almost exactly 2/3 of 88, it occurred to us that perhaps
all optical observations of features on Mercury were made at approximately
even multiples of 88 days and hence at integral multiples of 59 days.  This
possibility can be ruled out rather quickly.  However, if the rotation rate had
always been determined from series of observations that, in total, spanned
only about a decade and if each set of such observations was made at about the
same time of year (more precisely, at the favorable elongation that occurs
every third synodic period), then the 88-day result would be understandable:
The synodic period is 116 days and hence   3   synodic periods is nearly equal
to 4 orbital periods which, in turn, is nearly equal to 6 spin periods of 59 days
each.  A close examination of all available data on Mercury's features showed
that not all of the rotation determinations could be explained by this numerical
relationship between the synodic and orbital periods.  There were definitely
some misidentifications of features -- as well as the glaring and universal er-
ror of failing to realize that an essentially diophantine equation may have more
than one solution.

    How can this "new" 59-day rotation period of Mercury be explained?  Is it
to be interpreted as a final, stable spin state or as an intermediate, transient
state?  The first proposal, put forward by Peale and Gold, argued that the
present rotation state was stable.  The authors reached this conclusion by
studying the effect on the spin of the tidal torque exerted by the Sun.  The mag-
nitude of the tidal bulge raised by the Sun is inversely proportional to the cube
of its distance from Mercury.  Similarly, the torque exerted on the bulge var-
ies with the inverse cube of the distance leading to the well-known result that
the tidal torque varies with the inverse sixth power of the Mercury-Sun dis-
tance.  Of course, we must realize that if the axis of symmetry of the tidal
bulge coincided with the Mercury-Sun line, then the tidal torque would vanish.
But the cyclically varying relative attraction of the Sun on the various parts of
Mercury, due to the differences in spin and orbital angular velocities, causes
oscillations in the planet  which are dissipative.  These energy losses cause

the tidal bulge axis to lag behind the motion of the Sun-Mercury line and are responsible for a net tidal torque which opposes the spin of Mercury as seen from the Sun. If Mercury were initially spinning rapidly and if its orbit were circular, then given sufficient time the solar torque would slow down Mercury's spin until its rotation and orbital periods were equal: eventually Mercury would be locked with the same face always pointing towards the Sun. But Mercury's orbit is decidedly noncircular (e ≈ 0.2). What will happen with such a large eccentricity? The orbital angular velocity of Mercury will be a maximum at perihelion, corresponding to an instantaneous period of about 56 days, and a minimum at aphelion, corresponding to an instantaneous period of about 132 days. As soon as Mercury's spin slows sufficiently for its rotation period to exceed 56 days, the qualitative picture changes. Near perihelion, where the orbital angular velocity will exceed the spin angular velocity, Mercury will appear to a solar observer to reverse its direction of rotation. As a consequence, the lagging tidal bulge will now lie on the other side of the Sun-Mercury line. The direction of the tidal torque will therefore be reversed near perihelion. Since the tidal interaction is most intense there, it is easy to see that the spin period would not have to increase much above the 56-day value before the solar tidal torque vanished when averaged over an orbital period. In other words, for this model of the tidal torque, the sign of the torque will reverse twice during an orbital period when the spin period lies between 56 and 132 days. Because of the larger torque magnitude near perihelion, the average torque will be zero for a spin period not much in excess of 56 days. Thus, Peale and Gold concluded that the 59-day period was probably stable against further tidal evolution. In this analysis the actual stable period could not be predicted since it was presumably dependent on the precise details of the tidal energy losses.

As we mentioned earlier, 59 days (or, more accurately, 58.65 days) is two-thirds of Mercury's orbital period. Is this relationship merely a coincidence or might there be more to the spin problem than is contained in the tidal-torque model? Colombo suggested that perhaps Mercury's spin period is <u>exactly</u> two-thirds the orbital value with the resonance lock being made possible by a permanent axial asymmetry of the planet's moment of inertia ellipsoid -- independent of the changing asymmetry caused by the solar tidal forces. If Mercury possessed such a deviation from axial symmetry then, even if its rotation axis were normal it its orbital plane, the Sun would exert a torque on this permanent deformation which

would also affect Mercury's spin. A successful model of the spin evolution must therefore consider the actions both of this torque and of the tidal torque.

Colombo and I developed a theoretical model to study the general proper-ties of these spin-orbit resonances, particularly the stability conditions in the presence of tidal torques. Before describing the model in detail, we must ex-amine a simple question: Can the torque exerted by the Sun on a permanent deformation exceed that exerted on the tidal bulge? If the answer is negative then the permanent deformation will not play a substantial role in the evolu-tion of the spin state. Now, of course, we don't know either the deformation or the tidal characteristics of Mercury and so we cannot calculate either torque. But we can make reasonable assumptions. If the energy lost in a tidal oscilla-tion is about 1% of the maximum stored elastic energy due to the tidal deforma-tion (i. e., if the Q of Mercury is about 100 which is the approximate value for the Earth) and if Mercury's permanent equatorial asymmetry is like the Moon's (i. e., if the fractional difference in principal equatorial moments of inertia is equal to $2 \times 10^{-4}$), then a simple computation shows that

$$\frac{T_p}{T_t} \approx 10^5 \ , \tag{64}$$

where $T_p$ denotes a typical value of the solar torque exerted on the permanent deformation and $T_t$ a typical value of the tidal torque. We therefore cannot ig-nore $T_p$: It is very unlikely that our Earth-Moon analogy could be in error by more than one or two orders of magnitude.

Although the radar data have so far shown only that the axis of rotation of Mercury is inclined by less than $25^\circ$ to the orbit normal, we shall consider only a two-dimensional model in which the spin axis remains normal to the or-bital plane. (It is easily shown from Euler's equations that small deviations from perpendicularity will not alter our conclusions.) We shall assume also that Mercury's orbit is a Keplerian ellipse -- the implications of this restric-tion will be pointed out later. Letting $\theta$ be the angle between the orbital semi-major axis and Mercury's principal axis of minimum moment of inertia, we may write the equation of motion for the spin as

$$C \ddot{\theta} \ = \ \vec{T}_p + \vec{T}_t \ , \tag{65}$$

where C denotes the largest principal moment of inertia, assumed to be about the polar axis. (The dot notation signifies differentiation with respect to time.) The torque $T_p$ has the well-known form

$$\vec{T}_p = -\frac{3}{2} GM_\odot (B-A) \frac{\sin 2(\theta - f)}{r^3} \hat{k} , \qquad (66)$$

where $M_\odot$ is the solar mass, $A < B$ are Mercury's principal equatorial moments of inertia, $f$ is the true anomaly of Mercury's orbit, $r$ the Sun-Mercury distance, and $\hat{k}$ a unit vector normal to the orbital plane and positive in a northerly direction. The expression for the tidal torque is less certain. We assumed the lag angle of the axis of symmetry of the tidal bulge to be constant in magnitude with its sign depending only on the relative values of the spin and orbital angular velocities. Remembering the inverse sixth power radial dependence, we then obtain

$$\vec{T}_t = \frac{\eta}{r^6} \text{ sign} (\dot{\theta} - \dot{f}) \hat{k} , \qquad (67)$$

where $\eta$ is a constant proportional to the lag angle and inversely proportional to Q for small lag angles (high Q). The equation of motion can be rewritten as

$$\ddot{\theta} = -\frac{\beta' \sin 2(\theta - f)}{r^3} - \frac{\alpha' \text{ sign} (\dot{\theta} - \dot{f})}{r^6} , \qquad (68)$$

where $\alpha'$ is proportional to $Q^{-1}$ and $\beta'$ to $(B-A)/C$. How can one solve this second-order highly nonlinear differential equation? Obviously no closed-form solution exists and one must resort to approximation methods. All must take advantage of the fact that the disturbing torques are small: the spin state is not changed appreciably during an orbital revolution. We will limit ourselves here to the examination of only one approach. First, we introduce new variables, taking the mean anomaly

$$M = nt \qquad (69)$$

as the independent variable ($n$ denotes the mean motion). Then using

$$r = \frac{a(1 - e^2)}{1 + e \cos f} , \qquad (70)$$

where $\underline{a}$ is the orbital semimajor axis and $\underline{e}$ the eccentricity, and choosing units in which both $\underline{a}$ and the orbital period $P_o$ are unity, we find

$$\frac{d^2\theta}{dM^2} = -\beta \left(\frac{1 + e \cos f}{1 - e^2}\right)^3 \sin 2(\theta - f) - \alpha \left(\frac{1 + e \cos f}{1 - e^2}\right)^6 \cdot$$

$$\cdot \text{sign} \left(\frac{d\theta}{dM} - \frac{df}{dM}\right) , \tag{71}$$

where

$$\alpha \equiv \frac{\alpha'}{4\pi^2} ; \quad \beta \equiv \frac{\beta'}{4\pi^2} . \tag{72}$$

Since both $\alpha$ and $\beta$ are very small, we shall solve Eq. (71) only to first order in $\alpha$ and in $\beta$. We write the solution in the form

$$\theta = \theta_o' + \omega_o' M + \alpha \theta_{1\alpha}(M) + \beta \theta_{1\beta}(M) - \mathit{l} \pi , \tag{73}$$

where the integer $\mathit{l}$ is inserted to restrict $\theta$ to the interval $(-\frac{\pi}{2} \le \theta \le \frac{\pi}{2})$ since dynamically we cannot distinguish a rotation of $180^o$ when only second moments are being considered. (That is, the inertia ellipsoid is symmetrical with respect to reflection through a plane determined by any two principal moments.) We impose the obvious boundary conditions

$$\theta(0) = \theta_o' ; \quad \frac{d\theta}{dM} = \omega_o' \quad \left(= \frac{1}{2\pi} \left. \frac{d\theta}{dt}\right|_{t=o}\right) , \tag{74}$$

and

$$\theta_{1\alpha, \beta}(0) = \left. \frac{d\theta_{1\alpha, \beta}}{dM}\right|_{M=0} = \omega_{1\alpha, \beta}(0) = 0 , \tag{75}$$

Since there is no closed form expression for f(M), we use a Fourier series to obtain the formal first order solution. In particular, we make use of

$$-\left(\frac{1 + e \cos f}{1 - e^2}\right)^3 \sin 2(\theta_o' + \omega_o' M - f) \equiv \sum_{j=-\infty}^{\infty} P_j(e) \sin \left(\left[j - 2\omega_o'\right]M - 2\theta_o'\right)$$

$$-\left(\frac{1 + e \cos f}{1 - e^2}\right)^6 \text{sign} \left(\omega_o' - \frac{df}{dM}\right) \equiv \sum_{j=o}^{\infty} T_j(e, \omega_o') \cos jM , \tag{7}$$

where only the positive values of $j$ are required in the second series since the left side is symmetric with respect to the transformation $M \to -M$. The explicit expressions for $T_j$ and $P_j$ for arbitrary $j$ are quite complicated. Of the $T_j$'s, only $T_o$ will be of interest to us: it is proportional to the tidal torque averaged over an orbital period:

$$\overline{T}_t \sim T_o = -(1 - e^2)^{-9/2} \begin{cases} -g(e); & \omega < (1-e)^{1/2}/(1 + e)^{3/2} \\ \\ h(e); & (1-e)^{1/2}/(1+e)^{3/2} < \omega < (1+e)^{1/2}/(1-e)^{3/2} \\ \\ g(e); & \omega > (1+e)^{1/2}/(1-e)^{3/2} \end{cases}$$

$$\tag{77}$$

where

$$g(e) = 1 + 3e^2 + \frac{3}{8} e^4 , \tag{78}$$

$$h(e) = \frac{1}{2\pi} \left[ 2(\pi - 2f_c) - 16e \sin f_c + 6e^2(\pi - 2f_o - \sin 2f_c) \right.$$

$$- \frac{16}{3} e^3 (3 \sin f_c - \sin^3 f_c)$$

$$\left. + \frac{1}{8} e^4 \{6(\pi - 2f_c) - 8 \sin 2f_c - \sin 4f_c\} \right] , \tag{79}$$

and

$$\cos f_c = e^{-1} \{(1-e^2)^{3/4} \omega^{1/2} - 1\}; \quad 0 \le f_c \le \pi . \tag{80}$$

The true anomaly $f_c$ denotes the orbital position at which the tidal torque reverses sign. For $\omega$ corresponding to spin periods $P_s$ greater than 132 days or less than 56 days, $T_o$ will have a constant magnitude but a positive or negative sign, respectively. No reversal of the sign of the tidal torque occurs during an orbital period in either of these latter cases. As samples of the $P_j$ values, we write

$$P_j = \begin{cases} \frac{17}{2} e^2 + 0(e^4) & ; \quad j = 4 \\ \\ \frac{7}{2} e + 0(e^3) & ; \quad j = 3 \\ \\ 1 - \frac{11}{2} e^2 + 0(e^4) & ; \quad j = 2 \\ \\ -\frac{1}{2} e + 0(e^3) & ; \quad j = 1 \\ \\ 0 & ; \quad j = 0 . \end{cases} \tag{81}$$

After one orbital revolution the solution given formally in Eq. (73) is:

$$\theta_{1\alpha}(2\pi) = 2\pi^2 T_o$$

$$\omega_{1\alpha}(2\pi) = 2\pi\, T_o$$

$$\theta_{1\beta}(2\pi) = \sum_{j=-\infty}^{\infty} P_j(e) \left\{ \frac{2\pi \cos 2\theta'_o}{(j - 2\omega'_o)} + \frac{\sin(2\theta'_o + 4\pi\omega'_o) - \sin 2\theta'_o}{(j - 2\omega'_o)^2} \right\}$$

$$\omega_{1\beta}(2\pi) = -\sum_{j=-\infty}^{\infty} P_j(e) \left\{ \frac{\cos(2\theta'_o + 4\pi\omega'_o) - \cos 2\theta'_o}{(j - 2\omega'_o)} \right\} \quad . \tag{82}$$

The essential point to notice is the appearance of resonances at $\omega'_o = (j/2)$ ; $j = \ldots -1, 0, 1, 2, \ldots$ The strength of the resonance depends on $P_j(e)$. As can be seen from Eq. (81), except for accidental degeneracies, the lowest power of $\underline{e}$ present in the expression for $P_j(e)$ increases linearly as $\underline{j}$ deviates more (in either direction) from $j = 2$, the synchronous value. For $e = 0.2$, the resonance is very strong for $j = 3$ which corresponds to $P_s = \frac{2}{3} P_o$. (The orbital period $P_o$ should not be confused with $P_j(e)$ for $j = 0$.)

To study the resonances we look for periodic first-order solutions to Eq. (71). That is, we seek solutions for which

$$\theta(2\pi) = \theta(0) \equiv \theta_p$$

$$\omega(2\pi) = \omega(0) \equiv \omega_p \quad . \tag{83}$$

For an investigation near the kth resonance, we assume that

$$\left| \omega'_o - 2k \right| << 1 \quad , \tag{84}$$

and retain only first-order terms in the difference. To this accuracy we then find that the conditions for a periodic solution are satisfied for

$$\sin 2\theta_p = \frac{\alpha\, T_o(e, \frac{k}{2})}{\beta\, P_k(e)} \quad , \tag{85}$$

and

$$\omega_p - \frac{k}{2} = -\beta Q_k(e) \cos 2\theta_p \quad , \tag{86}$$

where

$$Q_k(e) = \sum_{\substack{j = -\infty \\ j \neq k}}^{\infty} \frac{P_j(e)}{(j-k)} \quad . \tag{87}$$

For given values of $\underline{k}$ and $\underline{e}$, there are two real solutions for $\theta_p$ provided that $(\alpha/\beta)$ is small enough. In fact, the right side of Eq. (85) is essentially the ratio of the average tidal torque to the average torque exerted on the permanent deformation and, from the discussion following Eq. (64), we expect this ratio to be very small indeed. Of these two periodic solutions, the one with $|\theta_p| < \pi/4$ has the principal axis of minimum moment of inertia inclined at perihelion by less than $45°$ to the radius vector from the Sun.

One might wonder why $\omega_p$ is not exactly equal to k/2 for the resonance solution. The answer is simply that although the $\underline{average}$ spin angular velocity exactly equals k/2 for the resonance solution, the $\underline{instantaneous}$ spin angular velocity at perihelion deviates from k/2 by the amount given in Eq. (86).

The important remaining question about these solutions concerns their stability. This can be investigated in the usual manner by examining the behavior of the system in the neighborhood of the periodic solution. For $\alpha, \beta \ll 1$, it is relatively easy to show that one of the two periodic solutions is unstable and the other asymptotically stable if and only if

$$\frac{\partial T_o(e,\omega)}{\partial \omega}\bigg|_{\omega = \frac{k}{2}} < 0 \quad , \tag{88}$$

i.e., if and only if $\overline{T}_t(\omega)$ has a negative slope at $\omega = k/2$.

Having developed these general characteristics of the solution to Eq. (71), let us return to a specific examination of the implications for Mercury. Its present rotation rate puts Mercury near (or at) the k = 3 resonance. The stable solution here corresponds to $\theta_p \approx 0$: the "long" equatorial axis, the one corresponding to the minimum moment of inertia, is inclined only slightly to the Mercury-Sun line at perhelion. In Figure 17 we show several different orientations of Mercury as it orbits the Sun under the assumption that the planet is

locked in the k = 3 resonance. (The small value of $\theta_p$ at perihelion appears
to vanish because of the scale of the figure.) At successive passages through
perihelion, points a and b alternate pointing towards the Sun. These two are
the only points on Mercury's surface which, on the basis of this model, ever
have the Sun directly overhead at perihelion. We also see from the figure
that Mercury's spin angular velocity matches the orbital angular velocity very
closely for a wide arc around perihelion which accounts for the strength of
this resonance: · The slight inclination of the principal axis of minimum mom-
ent of inertia to the Sun-Mercury line is maintained approximately throughout
this region and so the corresponding contributions to the average torque are
all of the same sign.

For contrast, we show in Figure 18 the orientation of Mercury at several
orbital positions on the assumption that it is in the k = 2 (synchronous) reso-
nance. There is no close match of spin and orbital angular velocities at peri-
helion; hence, for comparable orientations at perihelion, the average torque
due to the permanent deformation, is greater for the k = 3 resonance. Fur-
thermore, for k = 3 Inequality (88) is satisfied showing that there is an asymp-
totically stable spin state for this resonance (see Figure 19).

One crucial question has yet to be answered: What is the likelihood, or
probability, that Mercury would actually be captured into the k = 3 spin-orbit
resonance state during the course of its evolution? To understand better a
discussion of this question, it is useful to describe first the behavior of the
spin state in the phase plane. Since there is only one degree of freedom --
$\theta$ -- the phase plane is two-dimensional. Without loss of generality, we may
restrict the discussion to the phase strip $[-(\pi/2) \leq \theta \leq (\pi/2)]$ since for our
model the dynamical system is invariant under a $180^\circ$ rotation of the planet.
It is also convenient to consider only a stroboscopic view of the spin state.
Since the change in spin angular velocity $\omega$ during one orbit is very small and
the change in $\theta$ almost uniform, it is unnecessary to follow the continuous evo-
lution of the spin state. We merely indicate its value at the time of each peri-
helion passage. With this procedure, the evolution of the spin state can be
characterized by a set of points in the phase plane. Since these points will in
general be quite close together, we shall draw continuous lines through them
for the sake of simplicity.

In Figures 20 to 22, we illustrate for several different physical conditions the possible evolutionary behavior of the spin state in the phase plane near a resonance. For Figure 20, the tidal torque was assumed to vanish. The stable periodic solution is at $\theta_p = 0$ (point $C_2$) and the unstable one is at $\theta_p = \pi/2$ (point $C_1$). In the absence of a tidal torque the spin does not slow down secularly and we do not have asymptotic stability. Each trajectory is traced and retraced with time. For a given set of initial conditions, of course, only one trajectory will be followed. The shape of each trajectory shown in the figure is determined in essence by the conversion between kinetic and potential spin energy, with the angular velocity considered relative to the resonance value. For $\theta$ near zero at perihelion, the potential energy is a minimum and the magnitude of the difference between the spin angular velocity and its resonance value is a maximum. With a constant tidal torque present tending to slow down the spin, we have the phase-plane motion illustrated in Figure 21. Note the gradual decrease in spin angular velocity along a trajectory relative to the behavior along the corresponding trajectory in Figure 20. In view of Eqs. (84) and (85), the positions in the phase plane of the stable and unstable periodic solutions move to the left and right, respectively, by equal amounts. Since the tidal torque is constant in this case, the trajectories starting from, and arriving at, $C_1$ intersect the $\theta$ axis (drawn at the resonance value of the spin angular velocity) at the same point I. Here the condition (88) for asymptotic stability is not satisfied: If the spin state initially lies outside of the trajectories connecting $C_1$ and I, subsequent evolution will not take it inside and, in particular, the spin state will never reach the periodic solution $C_2$. For a nonconstant tidal torque satisfying (88) in the vicinity of the resonance spin state, we have a situation like that shown in Figure 22. Because the tidal torque is weaker below the $\theta$-axis than above it, the point I separates into two as shown, i. e. the differential tidal torque drives $I_1$ to the left of $I_2$. Any initial spin state lying in the shaded zones will evolve so as to eventually "lock into" the stable periodic solution at $C_2$. However, for an initial spin state between the shaded zones, the subsequent evolution will <u>not</u> lead to the stable periodic solution; the planet will not be captured at $C_2$ despite the existence of this asymptotically stable resonance state. Whether or not capture occurs will depend in this example on the initial conditions. There exists a set for which eventual capture is ensured and a complementary set for which the resonance "barrier" will be penetrated.

A clearer understanding of the behavior of this type of dynamical system in the vicinity of a resonance can be obtained by considering the analogy with a pendulum suspended from a friction bearing. This analogous dynamical system, pointed out by    Counselman, illustrates the major aspects of the capture process and we shall describe it in some detail. Ignoring the friction in the bearing upon which the pendulum swings, we obtain the well-known equation of motion

$$\ddot{\theta} \;=\; -\,\omega_o^2 \sin\theta \quad, \tag{89}$$

where

$$\omega_o^2 \;=\; \frac{g}{\ell} \tag{90}$$

with $g$ being the acceleration of gravity and $\ell$ the length of the pendulum whose attachment rod we assume to be rigid but massless. Suppose now that the rod is attached rigidly to a circular bearing which, in turn, is in frictional contact with a rotating shaft. (Such a pendulum device can be constructed easily by connecting a rod to the shaft of an ordinary table fan, as was done by Counselman.) What torque will the rotating shaft exert on the pendulum rod? For simplicity we assume the torque $T_s$ to be proportional to the difference between the angular velocity $\omega_s'$ of the rotating shaft and the angular velocity $\dot{\theta}$ of the pendulum. That the torque should depend on the difference is reasonable: Were the two angular velocities equal, there would be no relative motion and, hence, we would expect no frictional torque. For this model, the equation of motion becomes

$$\ddot{\theta} \;=\; -\,\omega_o^2 + \gamma\,(\omega_s' - \dot{\theta}) \quad, \tag{91}$$

where $\gamma$ is simply a coupling constant, and where the positive direction of rotation of the pendulum is opposite to direction of rotation of the shaft [i.e., $\omega_s'$ in Eq. (91) is negative]. Replacing $\omega_s'$ by its absolute value $\omega_s$ yields

$$\ddot{\theta} \;=\; -\,\omega_o^2 \sin\theta - \gamma\,(\omega_s + \dot{\theta}) \quad, \tag{92}$$

and shows that the torque exerted by the shaft is greatest when pendulum and shaft rotate in opposite directions. To make the model similar to the planetary

situation, we assume $\gamma \omega_s << \omega_o^2$ which is the analog of $\alpha << \beta$. We also assume that $\omega_s$ is constant, the shaft's motion being primarily controlled by a powerful external agency (say, the fan's motor).

If the pendulum is moving in the direction of increasing $\theta$ with sufficient angular velocity, it will go "over the top", i.e. the pendulum will rotate. The shaft's torque will then oppose the rotation and tend to slow it down. Suppose that, after a time, the pendulum has slowed sufficiently so that it just manages to go over the top, i.e. so that at the top the angular speed is essentially zero. The total energy E at this position will then be solely potential:

$$E(\theta = \pi) \approx -\omega_o^2 \cos \pi = \omega_o^2 \quad , \tag{93}$$

where the value for E follows from the first integral of Eq. (92) with the last term ignored. On the next swing, with the torque due to the shaft still opposing the motion, the pendulum will <u>not</u> return all the way to the top but will stop and reverse direction at an orientation

$$\theta = \pi - \Delta \theta \quad , \tag{94}$$

where $\Delta\theta$ is small and represents the energy lost by the pendulum due to the frictional contact with the shaft:

$$E(\theta = \pi) - E(\theta = \pi - \Delta\theta) \approx \omega_o^2 (1 - \cos \Delta\theta) \approx \omega_o^2 \frac{(\Delta\theta)^2}{2}$$

$$\approx \int_{-\pi}^{\pi - \Delta\theta} T_s \, d\theta = \gamma \left[ (2\pi - \Delta\theta)\omega_s + \int_{-\pi}^{\pi - \Delta\theta} \dot{\theta} d\theta \right] .$$

$$\tag{95}$$

To evaluate the integral we invoke our assumption that the friction torque does not affect the motion appreciably during any given orbit and we replace $\dot{\theta}$ by its value for $\gamma = 0$:

$$\dot{\theta} = \omega_o [2(1 + \cos \theta)]^{1/2} \quad . \tag{96}$$

Hence, we easily find the desired expression for $\Delta\theta$ in terms of $\gamma$, $\omega_s$, and $\omega_o$:

$$\Delta\theta \approx \frac{2[\gamma(\pi\omega_s + 2\omega_o)]^{1/2}}{\omega_o} \quad . \tag{97}$$

Now suppose that the pendulum slowed down until, when at the top of its swing, its kinetic energy was $\delta E$. Then, provided that

$$\delta E < \left| \int_{-\pi}^{\pi} T_s \, d\theta \right| \quad , \tag{98}$$

the pendulum during the ensuing swing will stop and reverse its direction of motion at an angle $\theta = \pi - (\Delta\theta - \delta\theta)$ where

$$\delta\theta \approx \frac{\delta E}{2\omega_o[\gamma(\pi\omega_s - 2\omega_o)]^{1/2}} \quad . \tag{99}$$

When Inequality (98) is an equality, the pendulum just makes the return to the top (i.e., $\delta\theta = \Delta\theta$). We have therefore established that friction causes the initially rotating pendulum to reverse direction for the first time at an angle $\theta$ satisfying $\pi - \Delta\theta \leq \theta < \pi$. What happens next? With the direction of swing reversed, the torque $T_s$ now acts in the same direction in which the pendulum is moving and feeds energy to the pendulum. There will be some critical angle $\delta\theta_c$ such that if the pendulum first stops at $\theta = \pi - \delta\theta_c$, then $T_s$ will feed in just enough energy to allow the pendulum to again reach the top of its swing but this time while travelling in the opposite direction. Using the same energy arguments as before we find

$$\delta\theta_c \approx \frac{2[\gamma(\pi\omega_s - 2\omega_o)]^{1/2}}{\omega_o} \quad , \tag{100}$$

which is identical with the expression for $\Delta\theta$ except that $2\omega_o$ is replaced by $-2\omega_o$ in the numerator. This crucial difference occurs because the $\gamma\dot\theta$ part of $T_s$ always opposes the motion of the pendulum whereas the $\gamma\omega_s$ part either opposes or aids depending on the direction of pendulum motion.

The following conclusions should now be apparent: If the pendulum first stops at a value of $\theta$ satisfying $\pi < \theta \leq \pi - \delta\theta_c$, then in its reverse swing $T_s$ will add enough energy for the pendulum to go "over the top". The pendulum will then continue to receive energy from $T_s$ and will rotate faster and faster in a direction opposite to the original one. This condition is the analog to the penetration of the resonance "barrier". If, on the other hand, the pendulum first stops at a value of $\theta$ satisfying $\pi - \delta\theta_c < \theta \leq \pi - \Delta\theta$, then $T_s$ cannot add enough energy on the reverse swing for the pendulum to reach the top. Instead, the pendulum will again stop, but at a value of $\theta$ satisfying $-\pi < \theta \lesssim -(\pi-\Delta\theta + \delta\theta_c)$. The motion will again reverse with the maximum value of $\theta$ achievable being smaller than for the previous oscillation. We have, therefore, the conditions for damped oscillations and the pendulum will eventually come to rest at a (negative) value of $\theta$ where the gravity torque is just balanced by $T_s$. This final configuration is the analog of capture into the spin-orbit resonance.

Let us re-emphasize the reason for there being damped oscillations. On the clockwise part of the swing, the pendulum gains energy proportional to $2\gamma(\pi\omega_s - 2\omega_0)$ but on the counterclockwise return trip, it loses energy proportional to $-2\gamma(\pi\omega_s + 2\omega_0)$, where the proportionality constant is essentially the same for both parts of a given oscillation. The $2\gamma\pi\omega_s$ contributions cancel: the constant part of $T_s$ does not contribute to the damping. The $4\gamma\omega_0$ losses however, always accumulate causing the damping and representing that part of $T_s$ which allows the analog to Inequality (88) to be satisfied.

The question of whether capture or penetration will occur rests with the initial conditions. We may speak of the probability of capture if we assign a priori probabilities to the possible initial conditions. For example, if we assume all values of $\delta E$ [see Eq. (98)] to be equally probable, then the capture probability for our model will be

$$P_c = 1 - \frac{\delta\theta_c}{\Delta\theta} \underset{\omega_0 << \omega_s}{\approx} \frac{2\omega_0}{\pi\omega_s} \quad , \tag{101}$$

and the penetration probability $p_p$ will be

$$p_p = 1 - \frac{\delta\theta_c}{\Delta\theta} = \left(\frac{\pi\omega_s - 2\omega_o}{\pi\omega_s + 2\omega_o}\right)^{1/2}$$

$$\underset{\omega_o << \omega_s}{\approx} \quad 1 - \frac{2\omega_o}{\pi\omega_s} \tag{102}$$

Having completed our discussion of the pendulum analogy, we return to our discussion of Mercury's axial rotation. From Figure 19, we see that Inequaltiy (88) is satisfied in the vicinity of the k = 3 resonance. Hence, for our model of the tidal torque and with the given orbital conditions, the capture of Mercury into the three-halves spin-orbit resonance state is possible. For the higher-order resonances, there is no possibility of capture under these assumptions because Inequality (88) is not satisfied for k > 3. The first calculations of the probability of Mercury's being captured into a resonance spin state were carried out by Goldreich and Peale. They assumed an a priori distribution of initial conditions that corresponded in essence to the assumption that all values of δE [see Eq. (98) and accompanying discussion] were equally probable. Using reasonable values for the Q of Mercury and several different models for the tidal torque, Goldreich and Peale showed that for e = 0.2 a substantial probability of capture into the k = 3 resonance exists only for [(B-A)/C] ~ $10^{-4}$. In other words, a fairly substantial axial asymmetry of the inertia ellipsoid is required for the probability of capture to be appreciable under these circumstances. These calculations were all based on the assumption that Mercury's orbit remains constant during the capture process. But the eccentricity -- the most important element -- is known to undergo oscillations caused by the perturbations of the other planets. The period of these oscillations seems to be neither very long nor very short relative to the time for capture to take place, given that the spin state is already in the vicinity of the resonance. To obtain more reliable results, these eccentricity variations must be considered.

In summary, we can say that considerable theoretical progress has been made in understanding the process of Mercury's capture into a spin-orbit resonance -- if indeed it is so captured -- but that some aspects remain to be investigated. To connect this discussion with the earlier ones, we remark

that the spin-orbit resonance coupling also contributes to the advance of Mercury's perihelion. When planetary observations become capable of distinguishing advances as small as 0.01 per century, this coupling effect must be considered in testing general relativity to this level of accuracy. Conversely, given that general relativity is correct, such measurements provide another means for the estimation of the differences in Mercury's principal moments of inertia.

The spin of Venus is even more interesting than Mercury's. The sidereal period of the retrograde rotation of Venus, as was pointed out earlier, is very nearly 243.1 days. This period is remarkable because it implies that Venus presents essentially the same face to the Earth at every inferior conjunction. Put more starkly, the Earth -- not the Sun -- appears to control the spin of Venus. The resonance between Venus' spin and the relative orbital motions of the Earth and Venus is precise for a period of 243.16; the measured value is in agreement to about 1 part in $10^3$ which seems too close to be simply a coincidence. How can we explain Venus' capture into such a resonance spin state? We may develop a simple model as for Mercury but with the addition of the torque $T_p^{(E)}$ exerted by the Earth on the permanent axial asymmetry of Venus. There are, of course, other torques; for example, the torque exerted by the Earth on the tidal bulge raised by the Sun. However, simple order-of-magnitude estimates suffice to show that each of these tidal-torque combinations is negligible except for the Sun's torque on the solar-induced bulge.

Unfortunately, there is not sufficient time remaining to develop the mathematical analysis of the evolution of Venus' spin state. Only a brief, qualitative discussion is possible. First we note that for Venus' spin to be held in a resonance "lock" by the Earth, the average $T_p^{(E)}$ must exceed the solar tidal torque. Otherwise, the Earth's torque could not possibly prevent penetration of the resonance. If we assume that Venus' primordial spin rate was rapid (period $\approx$ 1 day) and that the solar tidal torque has remained approximately constant over Venus' lifetime ($\approx 5 \times 10^9$ yr), then for $\overline{T}_p^{(E)}$ to exceed the tidal torque the fractional difference in Venus' equatorial moments of inertia would have to be about $10^{-3}$ -- a very large and difficult to understand value. By assuming either that Venus did not have a rapid primordial spin or that the tidal

torque was much larger in the past, we can relax this requirement to $[(B-A)/C] \gtrsim 10^{-4}$. Any smaller value would require the present Q of Venus to be "unreasonably" high. One might also wonder whether the solar torque exerted on this permanent asymmetry would not overwhelm the Earth's torque since, on an instantaneous basis, the former is at least $10^5$ times larger. In the close vicinity of an Earth resonance, however, the average solar torque vanishes, whereas the average Earth torque does not. A further problem arises when we consider the question of capture probabilities and asymptotic stability. For the tidal torque model considered in the Mercury analysis, the average value will be independent of $\omega$ for retrograde orbits. Hence, Inequality (88) cannot be satisfied and capture is not possible. To circumvent this difficulty, Bellomo, Colombo, and I introduced a generalized viscous tidal torque model for which Inequality (88) would be satisfied. More recently, Goldreich and Peale have proposed that Venus might have a core. By considering a two-degree of freedom problem (one orientation variable each for core and mantle) with a simple form for the coupling between the two parts, they calculated probabilities for capture into Earth-controlled spin resonance states. If the coupling is such that the response time of the core to changes in the motion of the mantle is about $3 \times 10^4$ years, then the maximum capture probabilities are obtained. These are still quite small -- less than 0.1, but certainly not negligible.

The final, and perhaps most important, question to be answered is: How did Venus' spin state evolve to the neighborhood of the 243-day-period resonance value? If Venus had a primordial, direct rotation like almost all the other planets, then, because of its small orbital eccentricity, Venus should have been captured into the k = 2 (synchronous) resonance with the Sun. It is hard to imagine any continuously-acting torque that would have prevented capture into this resonance and yet allowed capture into the Earth resonance. One could, however, invoke an impulsive torque by imagining that after Venus evolved to the synchronous resonance it collided with about a 200-km-diameter asteroid which produced a retrograde spin of magnitude close to the Earth resonance value. Such a cataclysmic event has the additional pleasing feature of providing a possible explanation for the rather high value of (B-A)/C required for Venus' spin to be controlled by the Earth.

Since we know so little about the formation of the planets, we can also postulate that Venus' spin was originally retrograde. Although the core-mantle model developed by Goldreich and Peale leads to a relatively small probability of capture into any given Earth resonance, the probability that Venus' spin will be captured into some Earth resonance state as its rotation slows becomes more substantial. (Because Venus' orbital eccentricity seems to remain very low, the likelihood of permanent capture into, say, the $k = -2$ Sun resonance is small.) Despite this substantial progress, further study -- both theoretical and experimental -- is required before the spin of Venus can be considered well understood.

Since radar astronomy has revivified the study of the spin and orbital motions of the planets, we can expect in the next few years to obtain a deeper understanding of the evolution of planetary spins as well as more definitive tests of the underlying theory of gravitation.

BIBLIOGRAPHY

M. E. Ash, I. I. Shapiro, and W. B. Smith, Astron. J. 72, 338 (1967).

E. Bellomo, G. Colombo, and I. I. Shapiro, in Mantles of the Earth and Terrestrial Planets (ed. S. Runcorn, Interscience, London, 1967), p. 193.

F. W. Bessel, Berl. Astron. Jahrbuch, p. 253 (1813).

C. Brans and R. H. Dicke, Phys. Rev. 124, 925 (1961).

R. L. Carpenter, JPL Res. Summary 36-14, p. 56 (1962).

_____, Astron. J. 69, 2 (1964).

_____, Astron, J. 71, 142 (1966).

G. Colombo, Nature 208, 575 (1965).

G. Colombo and I. I. Shapiro, Astrophys. J. 145, 296 (1966).

C. C. Counselman, Ph. D. Thesis, M. I. T., in preparation.

W. DeSitter, Mon. Not. Roy. Astron. Soc. 77, 155 (1917).

R. H. Dicke and H. M. Goldenberg, Phys. Rev. Lett. 18, 313 (1967).

P. A. M. Dirac, Proc. Roy. Soc. (London) A165, 199 (1938).

R. B. Dyce, G. H. Pettengill, and I. I. Shapiro, Astron. J. 72, 351 (1967).

A. Einstein, Ann. Phys. 49, 769 (1916).

R. V. Eötvös, D. Pekar, and E. Fekete, Ann. Physik. 68, 11 (1922).

J. V. Evans et al., Astron. J. 71, 904 (1966).

P. Goldreich and S. J. Peale, Astron. J. 71, 425 (1966).

_____, Astron. J. 72, 662 (1967).

R. M. Goldstein, Astron. J. 69, 12 (1964).

J. Lense and H. Thirring, Phys. Z. 19, 156 (1918).

W. E. McGovern, S. H. Gross, and S. I. Rasool, Nature 208, 375 (1965).

G. Nordström, Phys. Z. 13, 1126 (1912).

S. J. Peale and T. Gold, Nature 206, 1240 (1965).

G. H. Pettengill and R. B. Dyce, Nature 206, 1240 (1965).

G. H. Pettengill and I. I. Shapiro, in Ann. Rev. Astron. Astrophys. (ed. L. Goldberg, Ann. Rev., Palo Alto, 1965), Vol. 3.

R. V. Pound and G. A. Rebka, Phys. Rev. Lett. 4, 397 (1960).

R. V. Pound and J. L. Snider, Phys. Rev. Lett. 13, 539 (1964).

P. G. Roll, R. Krotkov, and R. H. Dicke, Ann. Phys. (N.Y.) 26, 442 (1964).

G. V. Schiaparelli, Astron. Nach. 123, 241 (1889).

L. I. Schiff, Proc. Nat. Acad. Sci. 46, 871 (1960).

I. I. Shapiro, Phys. Rev. Lett. 13, 789 (1964).

_____, Icarus 4, 549 (1965).

_____, Phys. Rev. 141, 1219 (1966).

_____, Science 157, 423 (1967.

_____, Science 157, 806 (1967).

_____, Astron. J. 72, 1309 (1967).

I. I. Shapiro, M. E. Ash, and M. J. Tausner, Phys. Rev. Lett 17, 933 (1966).

W. B. Smith, Astron. J. 68, 15 (1963).

## TABLE I

Osculating elliptic orbital elements at JED 2439340.5
referred to mean equinox and equator of 1950.0

| Orbital element | Mercury | | Venus | | Earth-moon barycenter | |
|---|---|---|---|---|---|---|
| | Newtonian fit | Relativity fit | Newtonian fit | Relativity fit | Newtonian fit | Relativity fit |
| Semimajor axis a (a.u.) | 0.3870984504 | 0.3870984149 | 0.7233299028 | 0.7233298596 | 1.0000038381 | 1.0000037947 |
| Eccentricity e | 0.2056252912 | 0.2056252616 | 0.0067893379 | 0.0067893036 | 0.0166821339 | 0.0166821414 |
| Inclination i (deg) | 28.6032096 | 28.6032349 | 24.4665221 | 24.4665152 | 23.4435784 | 23.4435741 |
| Right ascension of ascending node $\Omega$ (deg) | 10.8602455 | 10.8601939 | 7.9788463 | 7.9788228 | 0.0004331 | 0.0004358 |
| Argument of perihelion $\omega$ (deg) | 66.9331816 | 66.9333922 | 123.3335336 | 123.3336740 | 102.1909322 | 102.1909785 |
| Initial mean anomaly $l_0$ (deg) | 269.8932013 | 269.8930787 | 297.4518134 | 297.4517291 | 208.7069508 | 208.7069368 |

## TABLE II

### Formal standard errors of orbital element estimates

| | Mercury | | | Venus | | | Earth-moon barycenter | | |
|---|---|---|---|---|---|---|---|---|---|
| | Radar and optical | Radar | Optical | Radar and optical | Radar | Optical | Radar and optical | Radar | Optical |
| a(a.u.) | $1.0 \times 10^{-9}$ | $1.2 \times 10^{-9}$ | $1.4 \times 10^{-8}$ | $2.1 \times 10^{-9}$ | $3.1 \times 10^{-9}$ | $3.5 \times 10^{-8}$ | $7.4 \times 10^{-9}$ | $9.1 \times 10^{-9}$ | $1.4 \times 10^{-7}$ |
| e | $2.2 \times 10^{-8}$ | $2.3 \times 10^{-8}$ | $7.3 \times 10^{-7}$ | $1.7 \times 10^{-8}$ | $1.8 \times 10^{-8}$ | $3.4 \times 10^{-7}$ | $1.9 \times 10^{-8}$ | $2.1 \times 10^{-8}$ | $1.8 \times 10^{-7}$ |
| i | 0."04 | 6."5 | 0."16 | 0."02 | 6."5 | 0."05 | 0."02 | 6."5 | 0."02 |
| Ω | 0."10 | 27."0 | 0."31 | 0."08 | 31."7 | 0."12 | 0."07 | 33."1 | 0."07 |
| ω | 0."10 | 27."4 | 0."83 | 0."40 | 31."7 | 12."2 | 0."14 | 33."1 | 2."4 |
| $\ell_0$ | 0."03 | 0."03 | 0."86 | 0."40 | 0."43 | 12."2 | 0."12 | 0."13 | 2."4 |

# TABLE III

## Inverse masses of inner planets ($M_\odot = 1$) and Earth-Moon mass ratio

| | Present value | Newtonian fit | General relativity fit | Formal standard errors | | |
|---|---|---|---|---|---|---|
| | | | | Radar and optical | Radar | Optical |
| Mercury | 6110 000 ± 40 000 | 6 029 000 | 6 021 000 | 53 000 | 55 000 | 1 020 000 |
| Venus | 408 539 ± 12 | 408 450 | 408 250 | 120 | 124 | 2 200 |
| Earth and Moon | 328 906 ± 6 | 328 950 | 328 900 | 60 | 73 | 270 |
| Mars | 3 110 000 ± 7 700 | 3 106 700 | 3 111 200 | 9 000 | 11 600 | 27 000 |
| Earth-Moon mass ratio | 81.30 ± .01 | 81.3024 | 81.3030 | 0.005 | 0.005 | 0.34 |

TABLE IV

Astronomical unit and planetary radii

| | Present value | Newtonian fit | General relativity fit | Formal standard errors | | |
| | | | | Radar and optical | Radar | Optical |
|---|---|---|---|---|---|---|
| a.u. (light-sec) | 499.005 | 499.004785 | 499.004786 | $5.1 \times 10^{-6}$ | $7.4 \times 10^{-6}$ | ---- |
| Mercury radius (km) | 2420 | 2440.0 | 2434.0 | 2.2 | 2.4 | ---- |
| Venus radius (km) | 6100 | 6055.5 | 6055.8 | 1.2 | 1.4 | ---- |

## FIGURE CAPTIONS

Figure 1:  Radar path losses for solar-system targets as a function of echo delay. Minimum values associated with most and least favorable dates are shown for Mars and Mercury. Dotted bars and brackets indicate targets whose radar cross sections are unknown or highly variable; the values plotted were obtained under the assumption of 10% reflectivity and are given for reference only. The point for the minor planet refers to the close approach to earth on the date shown. For comparison, note that the path loss for the moon is $247$ db/m$^2$,

Figure 2:  Spherical-coordinate systems relating an arbitrary point P on the surface of the planet to the direction $\bar{R}$ of the radar site and that of the apparent angular velocity vector $\bar{\omega}$.

Figure 3:  Illustration of the principle of delay-Doppler mapping with the visible hemisphere of the planet projected on a plane normal to the radar-planet line.

Figure 4:  Geometric path of a radar pulse travelling between the Earth and a planet.

Figure 5:  Effect of general relativity on Earth-Mercury time delays.

Figure 6:  Effect of general relativity on Earth-Venus time delays.

Figure 7:  Time dependence of the deviation from initial values of the general relativistic osculating orbital elements of Mercury. (The values of semimajor axis and eccentricity are not zero at the initial time solely due to a slight translation error in placing the horizontal scale. )

Figure 8:  Sample of the residuals (observed minus computed: O-C) of the Earth-Sun optical observations obtained from a solution consistent with general relativity. (Note that time increases from right to left. )

Figure 9:  Same as Figure 8, except that the data refer to Earth-Venus observations. Note the systematic errors in the right-ascension residuals which indicate that the limb-to-center corrections for Venus need modification.

Figure 10: Same as Figure 8, except that the data refer to Earth-Mercury optical observations.

Figure 11: Residuals (O-C) of Earth-Mercury time-delay measurements obtained by comparison with (1) the JPL ephemerides; (2) a solution consistent with general relativity; and (3) a solution consistent with Newtonian theory. The JPL ephemerides are essentially Newcomb's orbits.

Figure 12: Same as Figure 11, except that the residuals refer to Earth-Venus time-delay measurements. The Newtonian solution, being not appreciably different from the relativistic one, has been omitted.
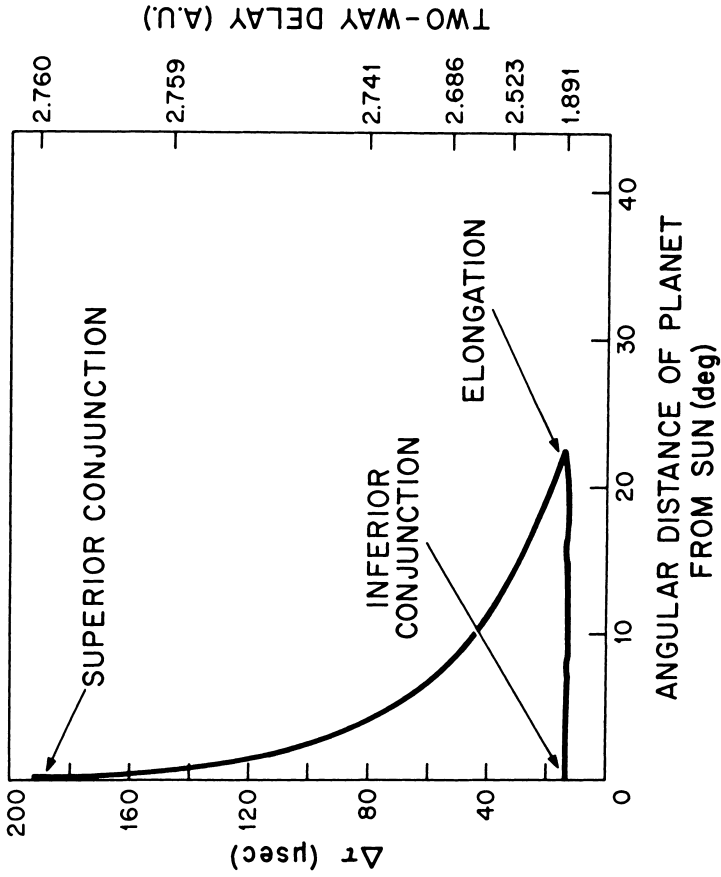
Figure 13:    Spectrum of 7750 MHz radar waves reflected from Venus and received at Haystack. Features correspond to those observed at other frequencies.

Figure 14:    Spectra of radar echoes from Mercury as a function of delay relative to the subradar point. Amplitudes have been scaled inversely by the factors listed in the right-hand column. The transmitted pulses were 100 $\mu$sec long; thus the spectra are essentially independent of each other. Error bars representing plus and minus one standard deviation of the measurement noise are shown near the left edge of each spectrum. The arrows indicate the expected positions of the edges of the spectra assuming a sidereal rotation period of 59 days with the rotation axis directed normal to Mercury's plane. These spectra were obtained at 430 MHz by the Arecibo Ionospheric Observatory.

Figure 15:    Determination of the rotation period of Mercury from observations taken at the Arecibo Ionospheric Observatory. Radar data similar to those shown in Fig. 14 were used to infer the bandwidths corresponding to reflections from the edges of the apparent disc. In the least-squares solution (upper dashed curves) the rotation axis was constrained to lie perpendicular to Mercury's orbital plane because of the limited accuracy of the available data. Removing this constraint changes the estimated rotation period by less than the quoted error of $\pm$ 3 days.

Figure 16:    Same as Figure 15, except that (1) the data refer to observations of Venus made during 1964, and (2) the rotation axis and the rotation period were determined simultaneously from the least-mean-square fit to the data.

Figure 17:    Expected orientation of Mercury's axis of minimum moment of inertia as a function of orbital position for the k = 3 resonance spin state $[P_s = (2/3)P_o \approx 58.65$ days$]$.

Figure 18:    Same as Figure 17, except that the orientations shown refer to the k = 2 (synchronous) resonance.

Figure 19:    Curves of the average values of $T_t$ and $T_p$ vs. $\omega(\equiv\dot{\theta}/n)$ for e = 0.2. Two curves are shown for $\overline{T}_p$, one with $\theta_o = 45^o$ and the other with $\theta_o = -45^o$. The ordinate scale is arbitrary and, in particular, the ordinates of $\overline{T}_t$ and $\overline{T}_p$ are not drawn to the same scale. The width of the resonances shown for $\overline{T}_p$ are meant only to be suggestive.

Figure 20:    Illustration of the phase-plane behavior of the spin state at successive perihelion passages in the absence of a tidal torque. For convenience, continuous curves are shown instead of points.

Figure 21:    Same as Figure 20, except that here a constant tidal torque is assumed to be present.

Figure 22:    Same as Figure 21, except that here the average tidal torque is assumed to have a negative derivative with respect to the spin angular velocity.
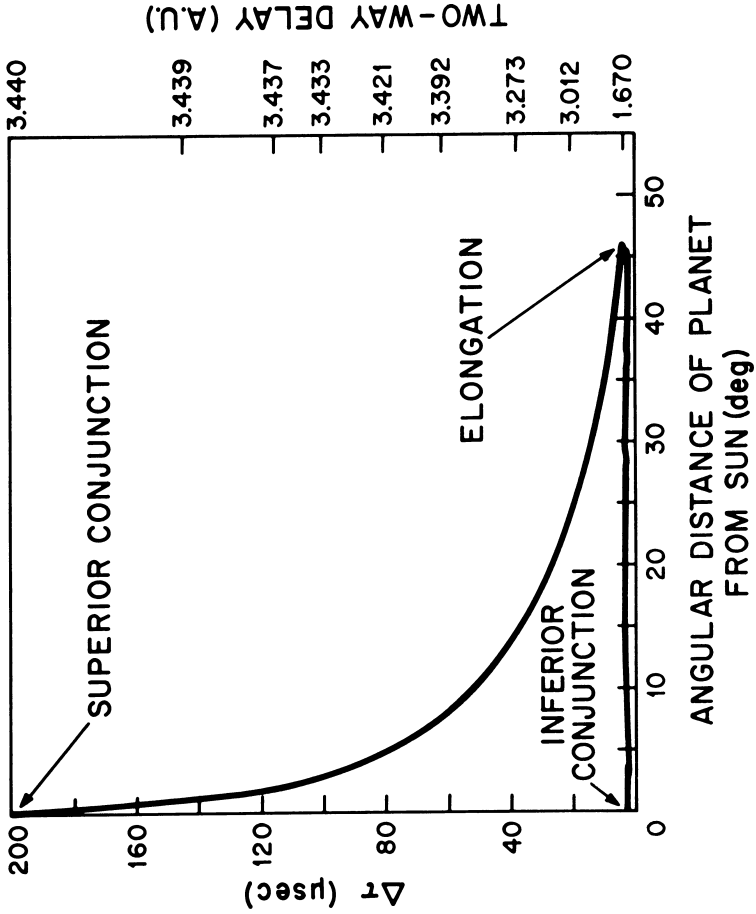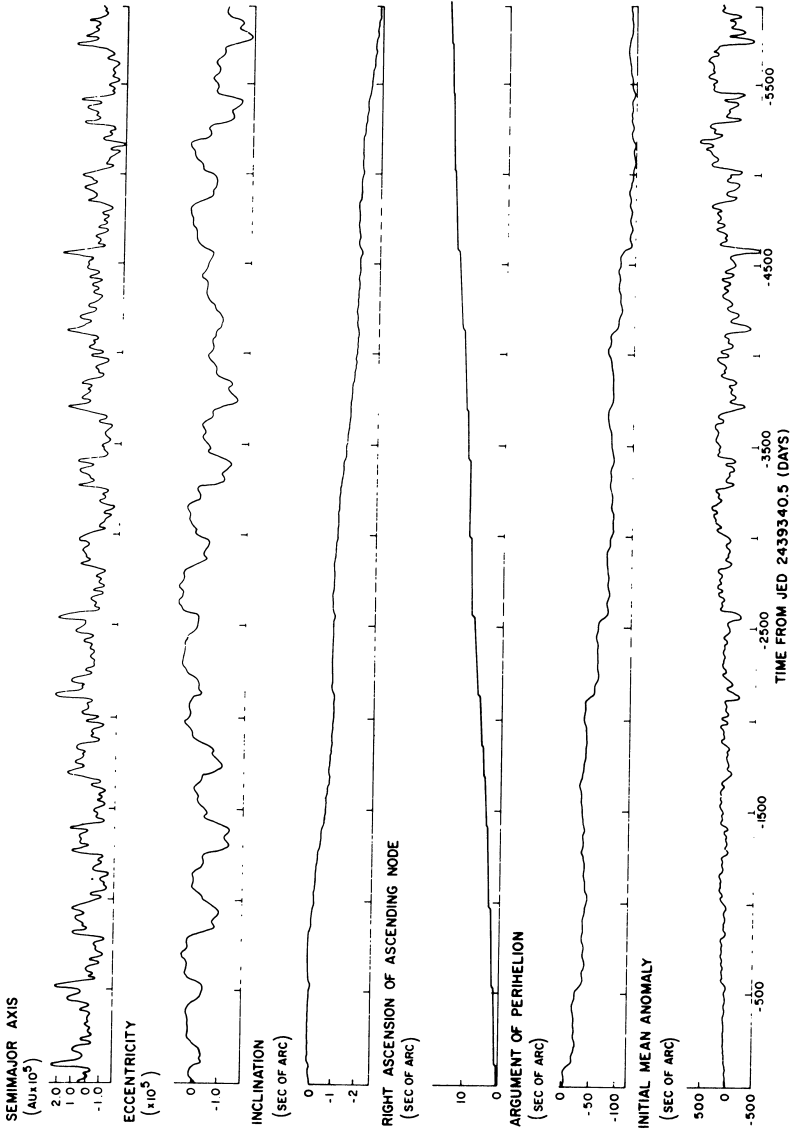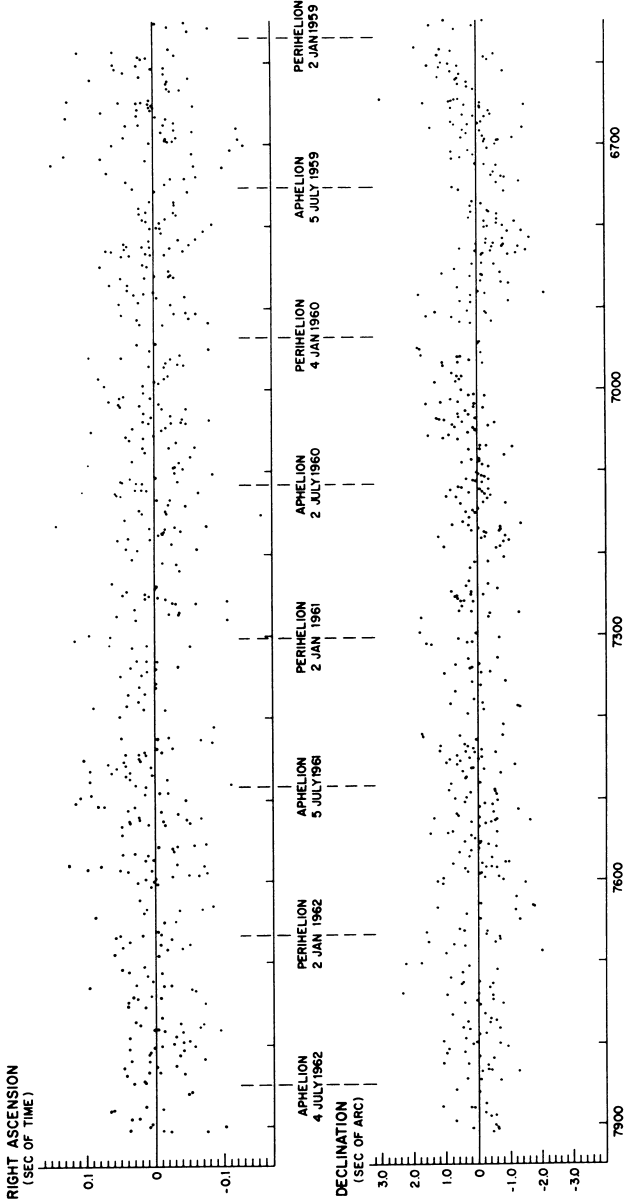
RELATIVISTIC EFFECT
ON EARTH-MERCURY TIME DELAYS

RELATIVISTIC EFFECT
ON EARTH-VENUS TIME DELAYS

SEMIMAJOR AXIS
(Au x 10⁻⁵)

ECCENTRICITY
(x 10⁻⁵)

INCLINATION
(SEC OF ARC)

RIGHT ASCENSION OF ASCENDING NODE
(SEC OF ARC)

ARGUMENT OF PERIHELION
(SEC OF ARC)

INITIAL MEAN ANOMALY
(SEC OF ARC)

TIME FROM JED 2439340.5 (DAYS)

VARIATION OF MERCURY OSCULATING ELLIPTIC ORBITAL ELEMENTS (GENERAL RELATIVITY)

RIGHT ASCENSION
(SEC OF TIME)

DECLINATION
(SEC OF ARC)

JULIAN EPHEMERIS DATE -2430000.5

U.S.NAVAL OBSERVATORY SUN RESIDUALS (GENERAL RELATIVITY FIT)

RIGHT ASCENSION
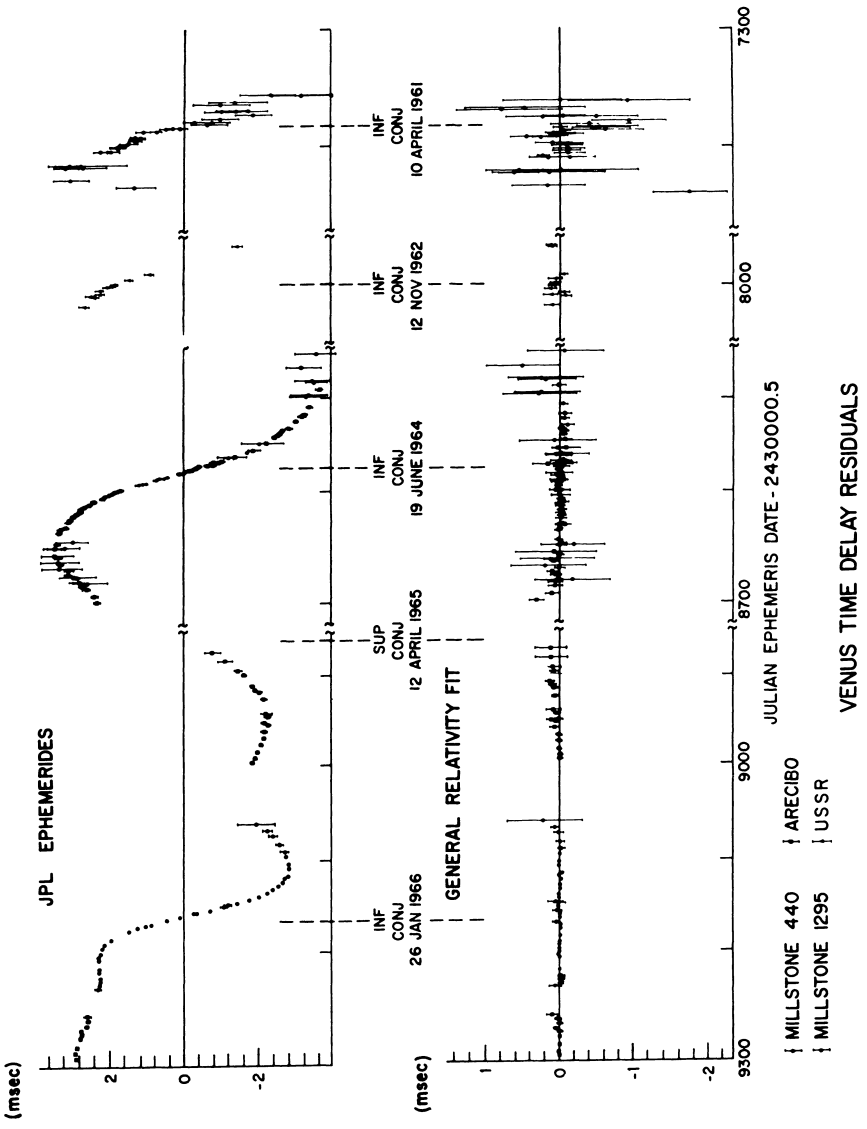(SEC OF TIME)

DECLINATION
(SEC OF ARC)

JULIAN EPHEMERIS DATE - 2430000.5

U.S. NAVAL OBSERVATORY VENUS RESIDUALS (GENERAL RELATIVITY FIT)

SUP
CONJ
27 JAN 1962

INF
CONJ
IO APR 1961

SUP
CONJ
22 JUNE 1960

INF
CONJ
7 AUG 1959

RIGHT ASCENSION
(SEC OF TIME)

DECLINATION
(SEC OF ARC)

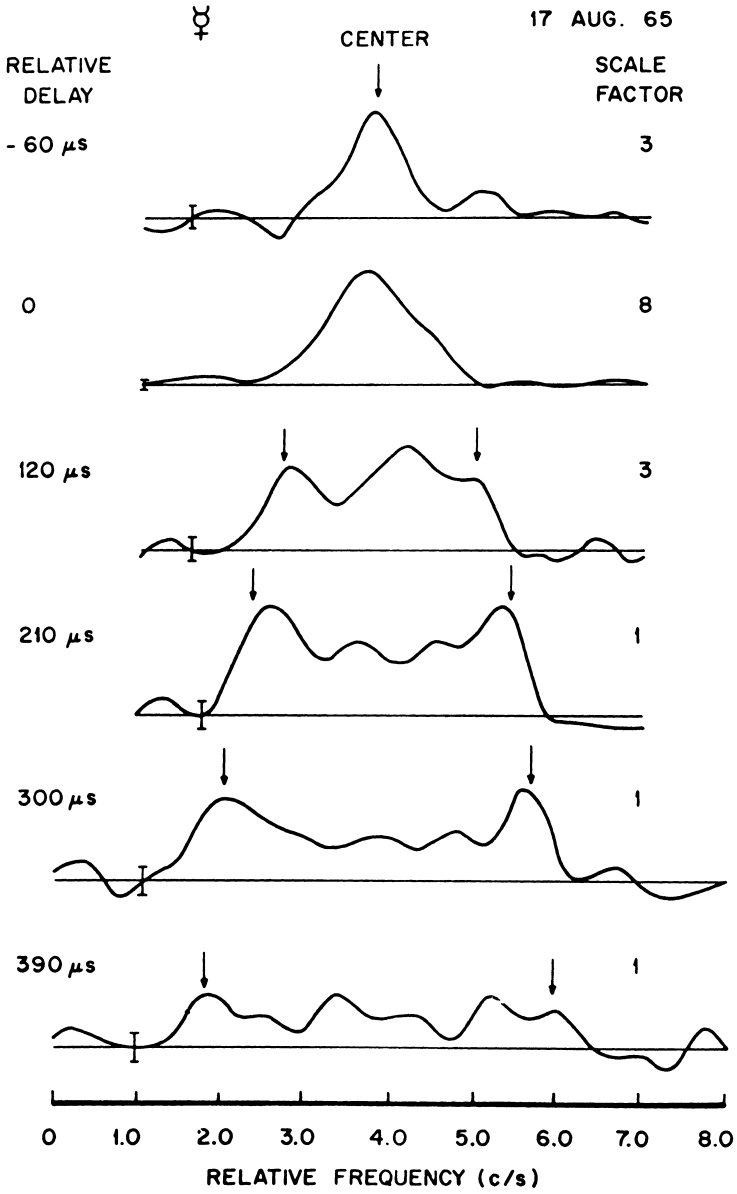JULIAN EPHEMERIS DATE -2430000.5
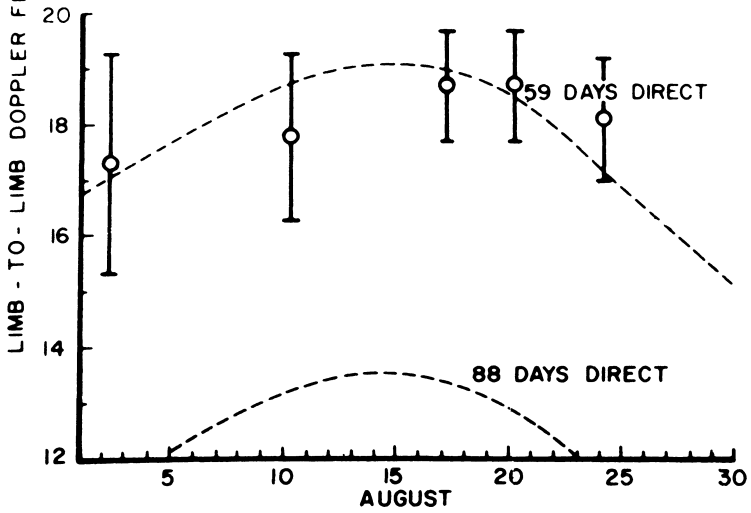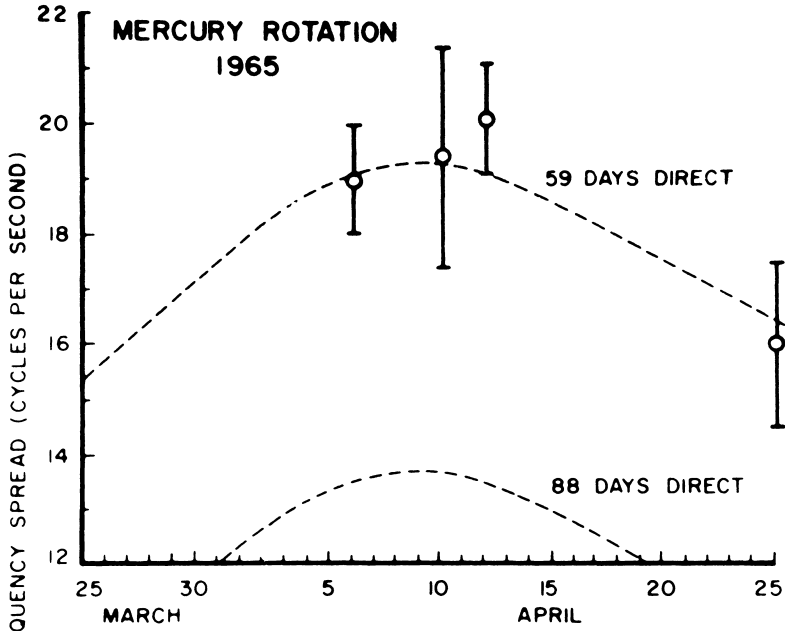
U.S. NAVAL OBSERVATORY MERCURY RESIDUALS (GENERAL RELATIVITY FIT)

MERCURY TIME DELAY RESIDUALS

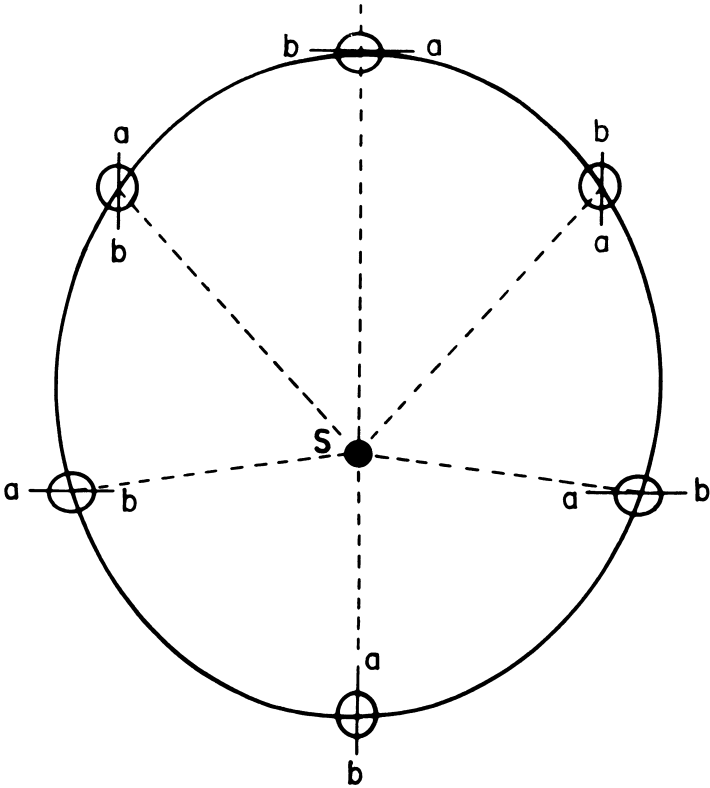VENUS TIME DELAY RESIDUALS

VENUS SPECTRUM
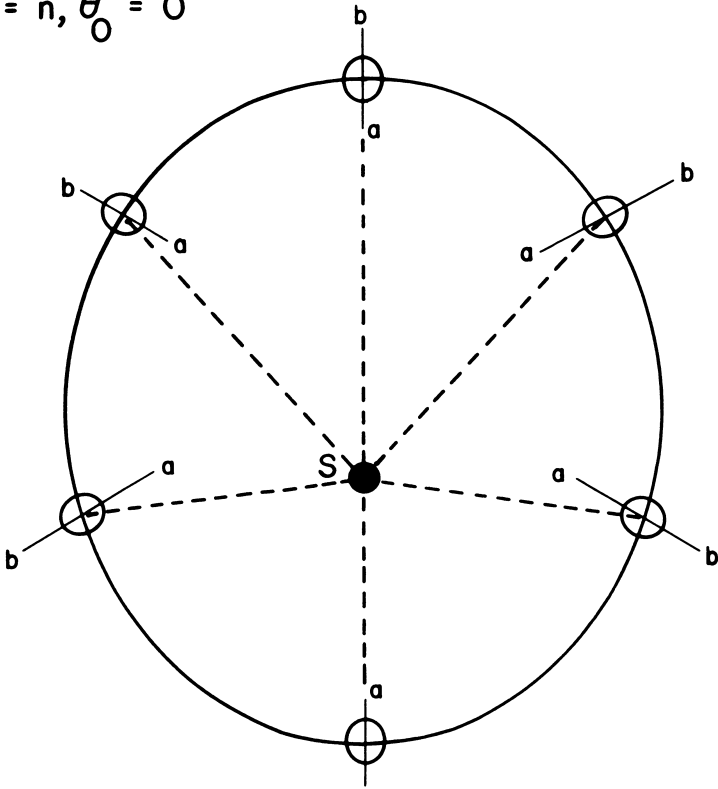λ = 3.8 cm
9 FEB 1966

$$\omega_0 = \frac{3}{2}n, \theta_0 = 0$$
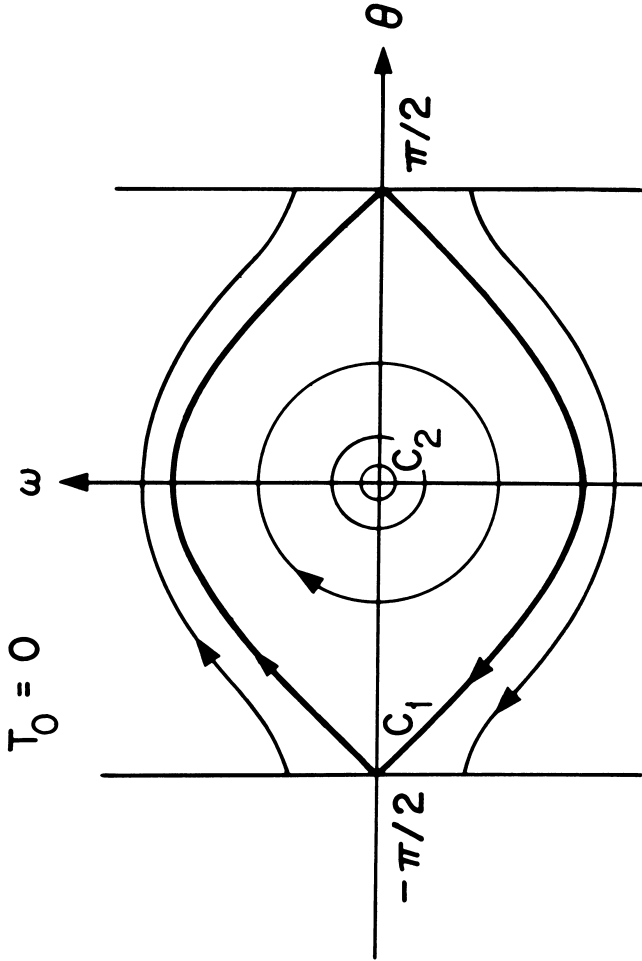
ORIENTATION OF MERCURY'S AXIS OF MINIMUM
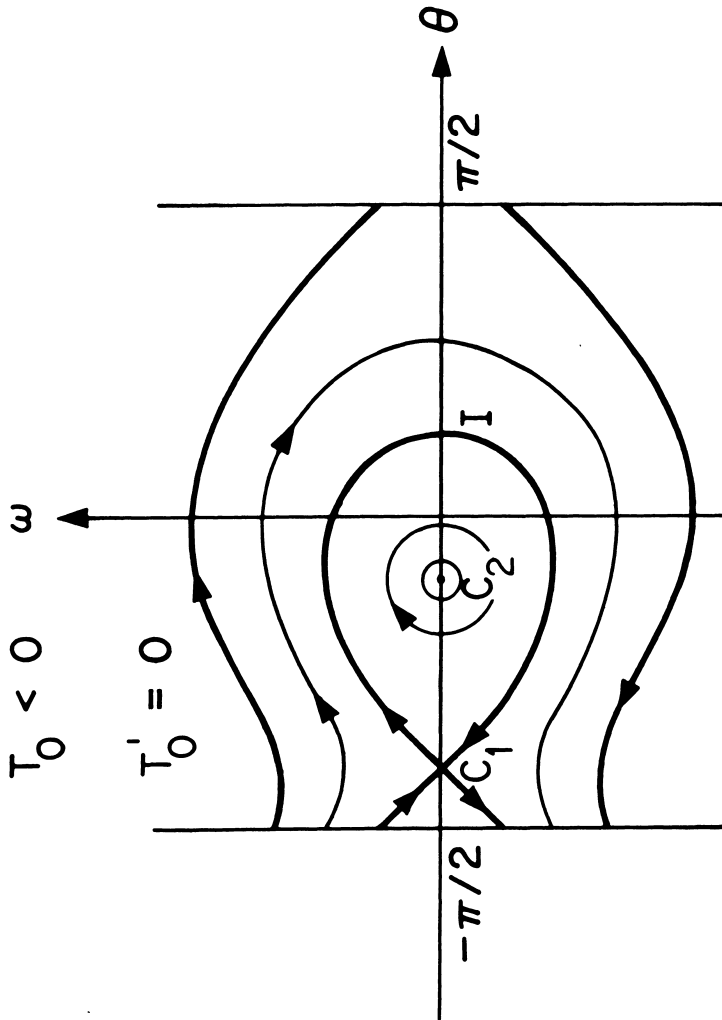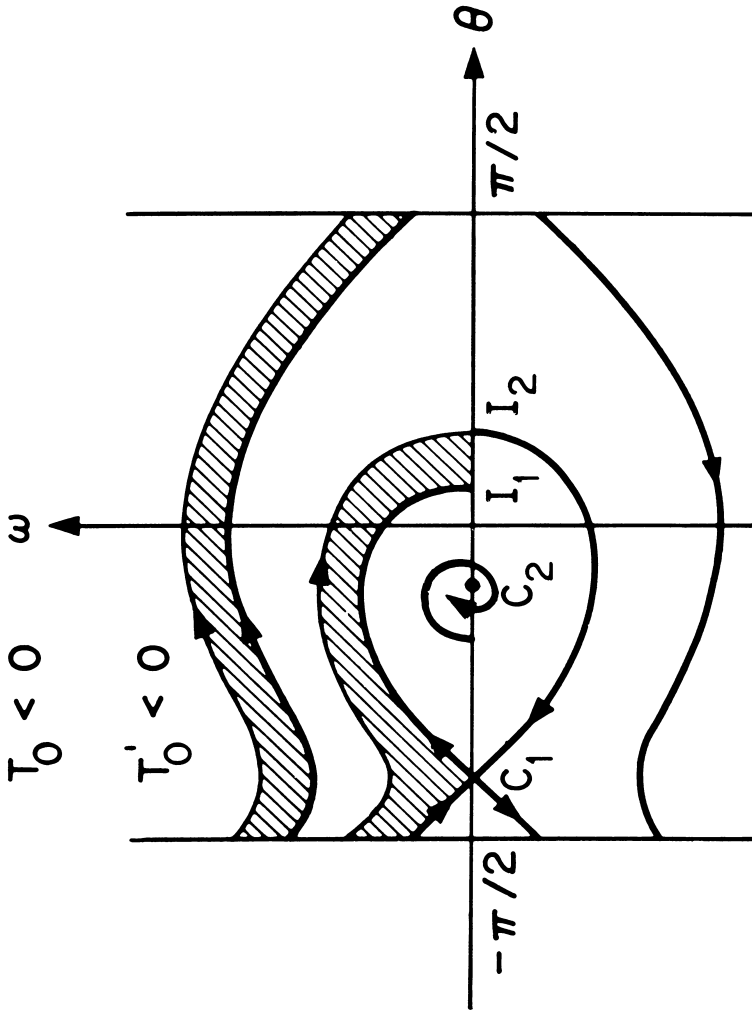MOMENT OF INERTIA

ORIENTATION OF MERCURY'S AXIS OF MINIMUM
MOMENT OF INERTIA

CENTRO INTERNAZIONALE MATEMATICO ESTIVO

(C. I. M. E.)

V. SZEBEHELY

"APPLICATIONS OF THE RESTRICTED PROBLEM OF THREE BODIES
IN SPACE RESEARCH"

# APPLICATIONS OF THE RESTRICTED PROBLEM OF THREE BODIES
## IN SPACE RESEARCH

by

V. Szebehely    (University-Yale )

## 1 . Statement of the problem

The problem under discussion is the restricted circular planar pro-
blem of three bodies (problème restreint) . Two  bodies (assumed to be point
masses and called primaries) revolve around their center of mass in circu-
lar orbits under  the influence of their  mutual gravitational attraction. A
third body (attracted by the previous two but not influencing their motion)
moves in the plane defined by the two revolving bodies . The problem is to
determine the motion of this third body .

The literature  often refers to the restricted problem when the prima-
ries move on conic sections and in order to further specify their motion
we speak of the circular restricted problem. Conventions,  however,  are not
well established  since some authors exclude all noncircular motion of the
primaries when speaking about the restricted problem.

If the initial position vector and velocity vector of the third body is
in the plane of the motion of the primaries,  there will  be no force direc-
ted out of this plane and the motion of the third body will take place in
this plane. It is generally accepted that the term "restricted problem" refers
to the three degrees of freedom (i.e. three dimensional motion)case, so the
further specification,  "planar" is necessary .

## 2. Equations  of  motion

Let the masses of the two primary bodies be    $m_1$ and  $m_2$ , their
angular velocity  n, and their  distance $\ell$  (see figure 1) . Then

(1)
$$k^2 M = n^2 \ell^3$$

V. Szebehely

where $M = m_1 + m_2$ and $k$ is the Gaussian constant of gravitation .

The center of mass of the system is located on the line connecting $m_1$ and $m_2$ , and its distance from $m_2$ and $m_1$ , is respectively

(2)
$$a = \frac{m_1 \ell}{M} \; , \quad b = \frac{m_2 \ell}{M} \; , \quad \ell = a + b \; .$$

Taking the origin of a fixed inertial coordinate system $(X, Y)$ at the center of mass , and using $t^*$ for time, the equations of motion referred to this system will be

(3)
$$\frac{d^2 X}{dt^2} = \frac{\partial F}{\partial X} \; , \quad \frac{d^2 Y}{dt^2} = \frac{\partial F}{\partial Y} \quad ,$$

where $F$ is Poincaré's "force function" or the negative potential energy and it is given by

(4)
$$F = k^2 \left( \frac{m_1}{\rho_1} + \frac{m_2}{\rho_2} \right)$$

with $\rho_1$, and $\rho_2$ being the distances between the primaries and the third body.
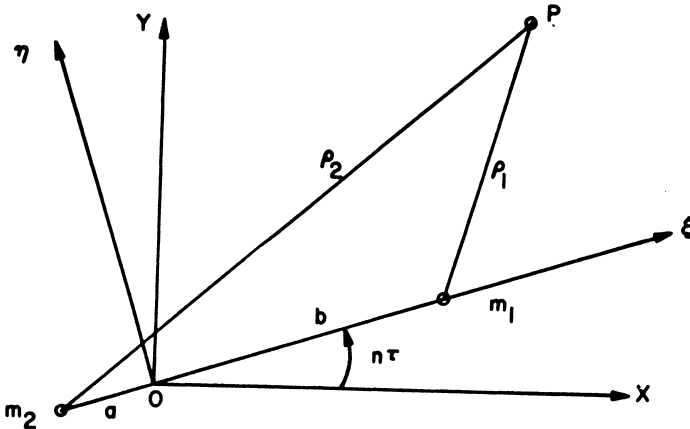


Figure 1 . Planar Circular Restricted Problem in Fixed $(X, Y)$ and Rotating $(\xi, \eta)$ Coordinate Systems

V. Szebehely

Introducing a uniformly rotating coordinate system ( $\xi$, $\eta$ ) with origin at the mass center so that $m_1$ and $m_2$ are located on the $\xi$ axis with coordinates (b, o) and (-a, 0), the equations of motion become

$$\frac{d^2 \xi}{dt^{*2}} - 2n \frac{d\eta}{dt^*} = \frac{\partial F^*}{\partial \xi},$$

(5)

$$\frac{d^2 \eta}{dt^{*2}} + 2n \frac{d\xi}{dt^*} = \frac{\partial F^*}{\partial \eta},$$

(6)   where   $F^* = F + \frac{1}{2} n^2 (\xi^2 + \eta^2)$

(7)   and   $\rho_1^2 = (\xi - b)^2 + \eta^2, \rho_2^2 = (\xi + a)^2 + \eta^2.$

The introduction of nondimensional quantities simplifies the equations .
Let

$$x = \xi / \ell , \quad y = \eta / \ell , \quad r_1 = \rho_1 / \ell , \quad r_2 = \rho_2 / \ell ,$$

(8)

$$t = n t^* , \text{ and } \mu = \frac{m_2}{M} .$$

The equations of motion in nondimensional form are :

$$\frac{d^2 x}{dt^2} - 2 \frac{dy}{dt} = \frac{\partial \overline{\Omega}}{\partial x} ,$$

(9)

$$\frac{d^2 y}{dt^2} + 2 \frac{dx}{dt} = \frac{\partial \overline{\Omega}}{\partial y} ,$$

where   $\overline{\Omega} = \frac{1}{2} (x^2 + y^2) + \frac{\mu}{r_2} + \frac{1 - \mu}{r_1} = \frac{F^*}{\ell^2 n^2}$   (10)

(11)   and $r_1^2 = (x - \mu)^2 + y^2 , \quad r_2^2 = (x + 1 - \mu)^2 + y^2 .$

Equations  9 ,  10 and 11 represent the problem in conventional non-

V. Szebehely

dimensional quantities. The corresponding physical picture is as follows (see Figure 2) : the two primary bodies are located on the x axis which rotates with a uniform angular velocity. The coordinates of the primaries are $P_1(\mu, o)$ and $P_2(\mu-1, o)$. The masses are $1-\mu$ and $\mu$, with $0 \leqq \mu \leqq 1$. The distance between the primaries, their total mass, their angular velocity and the gravitational constant are unity.

The Jacobi integral is obtained by multiplying the first of equations (9) by. 2 dx/dt, the second by 2 dy/dt and adding :

$$(12) \qquad \left( \frac{dx}{dt} \right)^2 + \left( \frac{dy}{dt} \right)^2 = 2 \, \overline{\Omega} - \overline{C} \quad ,$$

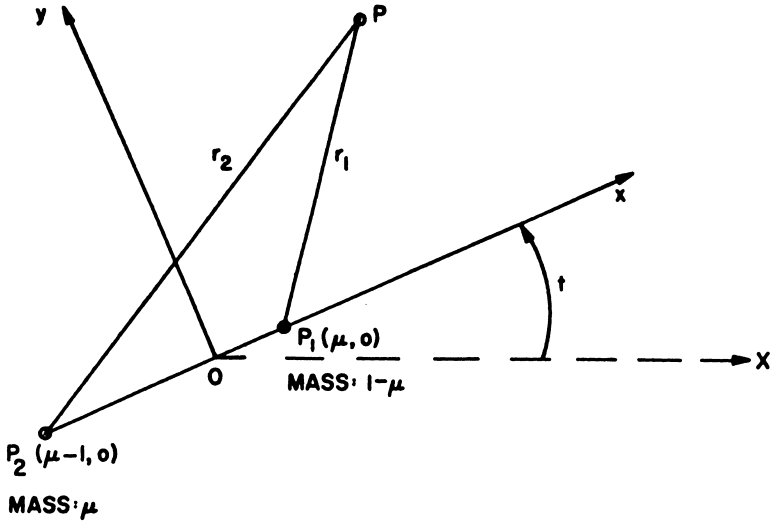where $\overline{C}$ is the constant of integration.



Figure 2. Planar Circular Restricted Problem in Nondimensional Rotating System

V. Szebehely

We note that according to equations (9) the partial derivatives of $\overline{\Omega}$ completely determine the problem, therefore by adding a constant to the expression for $\overline{\Omega}$ - as given equation (10) - will not change equations (9), but it will influence the value of $\overline{C}$ in equation (12). A symmetrical form of $\overline{\Omega}$ is obtained by adding the constant quantity:

$$\frac{1}{2} \mu(1 - \mu)$$

to $\overline{\Omega}$ and let $\Omega = \overline{\Omega} + \frac{1}{2} \mu(1- \mu)$, or

(13) $\qquad \Omega = \frac{1}{2} \left[ (1-\mu) \; r_1^2 + \mu r_2^2 \right] + \frac{1-\mu}{r_1} + \frac{\mu}{r_2}$ .

The new Jacobi constant $C$ is related to the previous one by $C = \overline{C} + \mu(1-\mu)$ .

The final set of equations, subject to the investigation, with the notations

$$\Omega_x = \frac{\partial \Omega}{\partial x}, \qquad \Omega_y = \frac{\partial \Omega}{\partial y}, \qquad \dot{x} = \frac{dx}{dt}, \quad \text{etc.}$$

becomes :

(14) $\qquad \ddot{x} - 2\dot{y} = \Omega_x, \qquad \ddot{y} + 2\dot{x} = \Omega_y$

(15) $\qquad (\dot{x})^2 + (\dot{y})^2 = 2\Omega - C$ ,

where $\Omega$ is given by equation (13) .

It is remarked that Levi-Civita gave the above equations in canonical form by introducing the canonical variables :

$$q_1 = x - \mu, \qquad p_1 = \dot{x} - y,$$

$$q_2 = y, \qquad p_2 = \dot{y} + x - \mu.$$

The equations of motion are

V. Szebehely

$$\dot{q}_i = \frac{\partial F}{\partial p_i} \quad ; \quad \dot{p}_i = -\frac{\partial F}{\partial q_i} \quad , \quad (i = 1, 2)$$

where $\quad F = \frac{1}{2}(p_1 + q_2)^2 + \frac{1}{2}(p_2 - q_1)^2 - \Omega(q_1, q_2)$ .

The $F = $ constant equation, of course, will correspond to the Jacobi integral .


### 3 . R e d u c t i o n s

#### 3.1 Reduction to the third order

The fourth order system of equations (14) can be reduced to a third order system by making use of the Jacobi integral, (equation 15). One way to accomplish this is by introducing the angle (z) between the tangent of the orbit and the positive x axis as a new dependent variable :

$$(16) \qquad\qquad z = \text{arc}\ \ \text{tg}\ \frac{dy}{dx} \ .$$

We will show that the equations of motion become :

$$(17) \qquad \begin{aligned} \dot{x} &= \bigwedge(x, y)\ \cos\ z \\ \dot{y} &= \bigwedge(x, y)\ \sin\ z \\ \dot{z} &= -2 - \bigwedge_x\ \sin z + \bigwedge_y\ \cos z, \end{aligned}$$

$(18)$ where $\quad \bigwedge = \sqrt{2\,\Omega - C}$ .

The proof of the first two equations of (17) is as follows. If $\underline{s}$ denotes the arc length,

$$\frac{dx}{ds} = \cos\ z, \ \frac{dy}{ds} = \sin\ z, \quad \text{and}\quad \text{so}$$

$$\dot{x} = \frac{dx}{ds}\ \dot{s} = \bigwedge \cos\ z, \dot{y} = \bigwedge \sin z,$$

since the absolute value of the velocity vector $(\dot{s})$ is obtained from the

Jacobi integral as

$$\left(\frac{ds}{dt}\right)^2 = (\dot{x})^2 + (\dot{y})^2 = 2\Omega - C = \bigwedge^2 \ .$$

The proof of the third equation of (17) requires the evaluation of $\dot{z}$ from equation (16) :

$$\dot{z} = \frac{\ddot{y}\ \dot{x} - \dot{y}\ \ddot{x}}{(\dot{s})^2} \ \ .$$

Using equations (14) to eliminate $\ddot{x}$ and $\ddot{y}$ , and the first two equations of (17) to eliminate $\dot{x}$ and $\dot{y}$ , we obtain the desired result .

We note that an alternate approach would be to use the first two of equations (17) as definitions of a transformation without reference to equation (16) . In this way we have $\dot{x} = \bigwedge \cos z, \dot{y} = \bigwedge \sin z$ and consequently $\ddot{x} = \dot{\bigwedge}\cos z - \dot{z}\bigwedge\sin z$ and $\ddot{y} = \dot{\bigwedge}\sin z + \dot{z}\bigwedge\cos z$, from which $\bigwedge \dot{z}$ can be obtained ; $\bigwedge \dot{z} = \ddot{y} \cos z - \ddot{x} \sin z$. By the same elimination process as before we can obtain the third equation of (17) .

Equations (17) represent the third order version of the problem in form of three first order differential equations with x, y, z as dependent and t as the independent variable. The significant fact is noted that these equations can be written as

$$\dot{x} = \phi(x, y, z),$$

(19)     $$\dot{y} = \psi(x, y, z),$$

$$\dot{z} = \chi(x, y, z),$$

that is, the right hand members do not contain the time. This fact allows a physical interpretation of the equations by an analogy and also assures further reduction of the order by elimination of the time.

V. Szebehely

## 3.2 Flow analogy

Consider a flow field with velocity $\bar{v} = \bar{v}(\bar{r}, t)$, where $\bar{r}$ is the position vector and t is the time. This velocity vector gives a description of the flow field since at every point $\bar{r}$ in the field, at any time t, the velocity can be evaluated - excepting singular points. A flow is called steady if

$$\frac{\partial \bar{v}}{\partial t} = 0 \quad ,$$

i.e. if none of the velocity components depend explicitly on the time. The velocity components defined by equations (19) can be interpreted therefore as the description of a steady flow field.

The continuity equation of hydrodynamics is

$$\frac{\partial \rho}{\partial t} + \text{div} \ (\rho \bar{v}) = 0$$

where $\rho$ is the fluid density. For an incompressible fluid $\rho$ = constant and the continuity equation becomes

$$\text{div} \ \bar{v} = 0 \ ,$$

or using the notation of equations (19)

$$\phi_x + \psi_y + \chi_z = 0 \ .$$

Since the $\phi$, $\psi$, $\chi$ velocity components as given by equations (17) satisfy this equation, the dynamical problem is analogous to the three dimensional steady flow of an incompressible fluid.

It is remarked that the flow is <u>not</u> a potential flow since

$$\text{curl} \ \bar{v} \neq \bar{o}$$

as computation of six partial derivatives ( $\phi_y$, $\phi_z$, $\psi_x$, $\psi_z$, $\chi_x$, $\chi_y$)
will convince the reader .

The function $\Lambda$ contains C, therefore equations (17) will determi-
ne a flow field for a givev C value. By changing the constant of integra-
tion (C) , the streamline picture will change. For a given C equations
(17) will determine the totality of motions of the dynamical system and
also the corresponding streamlines , provided that the inequality .

$$(\frac{ds}{dt})^2 = 2\,\Omega\,(x,y) - C \gtreqless O$$

is satisfied.

The streamline representation is singular when $\Lambda$ = o or when
$\Lambda \longrightarrow \infty$ . The first case corresponds to zero velocity , the second
to collisions at $P_1$ or $P_2$ .

Poincare's original flow analogy being similar to but not identical
in details with the previous one should be mentioned also. Let $x = x_1$ ,
$y = x_2$, $\dot{x} = x_3$ and $\dot{y} = x_4$ . Then the equations of motion can be written
as

$$\dot{x}_1 = x_3$$

$$\dot{x}_2 = x_4$$

$$\dot{x}_3 = 2\,x_4 + \frac{\partial\,\Omega(x_1,x_2)}{\partial\,x_1}$$

$$\dot{x}_4 = -2x_3 + \frac{\partial\,\Omega(x_1,x_2)}{\partial\,x_2}$$

or $\dot{x}_i = F_i(x_1,\ x_2,\ x_3,x_4)$ , $i = 1, \ldots 4$ .

It can be seen that div $\overline{v} = o$ , i.e.

$$\sum_{i=1}^{4} \frac{\partial\,F_i}{\partial\,x_i} = o\ ,$$

V. Szebehely

so we are dealing with the four dimensional steady stream line flow of
an incompressible fluid. The actual motion of the particle corresponds to
those streamlines which lie on the

$$x_3^2 + x_4^2 - 2 \Omega (x_1, x_2) + C = o$$

three dimensional hypersurface.

    We note that both Poincaré and Birkhoff contributed to the develop-
ment of the above discussed streamline flow analogy which, however,
has not resulted in any significant new information as yet. Some suggestions
regarding further developments and possible uses are given in the last cha-
pter.

### 3.3 Reduction to the second order.

    Equations (19) can be written as

(20) $$\frac{dx}{\phi} = \frac{dy}{\psi} = \frac{dz}{\chi} ,$$

from which

(21) $$\frac{dz}{dx} = \frac{\chi}{\phi} .$$

    On the other hand equation (16) gives

(22) $$\frac{dz}{dx} = \frac{d}{dx} \ arc \ tg \ y' = \frac{y''}{1 + (y')^2} ,$$

where $y' = \frac{dy}{dx}$ .

    Substituting for $\chi$ and $\phi$ in equation (21) their expressions as
given by equations (17), eliminating z by equation (16) and equating
equations (21) and (22) results in a second order differential equation
describing the dynamical system (excepting at points where the transforma-
tions are singular) :

(23) $$y'' = \frac{1 + (y')^2}{\Lambda} ( \Lambda_y - \Lambda_x y' - 2 \sqrt{1 + (y')^2} ) .$$

V. Szebehely

It is noted that elimination of time and use of the Jacobi integral can be combined and the above second order differential equation can be obtained in a single step directly from the original fourth order system. The general solution of (23) will contain the following three constants of integration : C which is included in $\wedge$ and the two constants which enter when (23) is integrated. The complete solution of the original fourth order system requires the determination of the time dependence of the variables. This process will result in the fourth integration constant. To establish the time dependence, we write the Jacobi integral as

$$(\dot{x})^2 \left[ 1 + (y')^2 \right] = \wedge^2 \ ,$$

where the $\dot{y} = \dot{x} \, y'$ relation was used . The time is evaluated from the last equation as

$$t = \int \frac{\sqrt{1 + (y')^2}}{\wedge} \ dx + C_4 \ ,$$

Here $y'$ and $\wedge(x,y)$ are functions of $x$ only since $y(x)$ has been obtained from equation (23) .

## 4 Regions of Motion

The second important application of the Jacobian integral is to establish regions of motion or to find the so-called "forbidden areas" of the plane $x, y$. The Jacobian integral may be written as

$$v^2 = 2 \, \Omega \, (x, y) - C \ .$$

At a given point $(x_o, y_o)$ along an orbit, let the speed relative to the synodic coordinate system be $v_o$ . Then the value of the Jacobian constant along this orbit may be computed from

V. Szebehely

$$C_o = 2\Omega(x_o, y_o) - v_o^2 \ .$$

Since this value must be constant for any point $(x, y)$ on the orbit, we have

$$C_o = 2\Omega(x, y) - v^2 \ ,$$

where $v$ is the speed at the point $x, y$. Solving the above equation for $v^2$, we have

$$v^2 = 2\Omega(x, y) - C_o \ .$$

If now have the point $x, y$ is selected such that $2\Omega(x, y)$ is larger than the precomputed value of $C_o$ then $v^2 > 0$ and motion is possible at such a point. If , on the other hand

$$2\Omega(x, y) < C_o$$

then $v^2 < 0$ and motion cannot take place at the point $x, y$.

The regions of possible motion are, therefore, separated from the forbidden regions by curves which are given by the equation

$$2\Omega(x, y) = C_o \ .$$

Along these curves the speed is zero and for this reason these cuves are called curves of zero velocity of Hill's curves.

To establish the forbidden regions the curves $\Omega = $ constant must be constructed. In what follows these "equipotential" or niveau curves and the $\Omega(x, y)$ function will be analysed in some detail.

According to the value of $C$ there are six distinctly different cases representing basically different problems as well physical situations. To review these cases we start with $C \rightarrow \infty$ and decrease its value to $C = 3$. For a detailed discussion the reader is referred to, for instance, C. L. Charlier's book (see bibliography) .

V. Szebehely

### Case I

As seen from

$$(24) \qquad \Omega = \frac{1}{2}\left[(1-\mu)\, r_1^2 + \mu r_2^2\right] + \frac{1-\mu}{r_1} + \frac{\mu}{r_2}$$

and from the definition of the zero velocity curves $(2\Omega - C = o)$, large values of C imply large $\Omega$ values . This occurs if one of the following conditions is approached :

$$r_1 \rightarrow o,\ r_2 \rightarrow o,\ r_1 \rightarrow \infty,\ \text{or}\ r_2 \rightarrow \infty.$$

Let $r_1 = \varepsilon$ , $r_2 \cong 1$, in order to investigate the first possibility . For this case

$$\Omega \simeq \frac{3}{2}\mu + \frac{1-\mu}{\varepsilon}$$

and the velocity square from the Jacobi integral becomes :

$$(25) \qquad v^2 \simeq \frac{2(1-\mu)}{\varepsilon} + 3\mu - C\ .$$

So the zero velocity line is approximately a circle around $P_1$ with radius

$$(26) \qquad \varepsilon = \frac{2(1-\mu)}{C - 3\mu} = r_{10}\ .$$

This oval shrinks as C increases, so if a zero velocity oval corresponding to a given $C_1$ value is constructed around $P_1$ , another zero velocity oval with $C = C_2$ will be inside of the first one if $C_2 > C_1$ . If a particle has given velocity and position vector , i.e. its state of motion is defined giving for its Jacobi constant a value equal to $C_1$ then it can move only inside of the $C_1$ oval. This follows if we consider that on the $C_1$ oval

$$V^2 = 2 \ \Omega(x, y) - C_1 = 0$$

and inside of this oval there are $2 \ \Omega(x, y) - C_2 = 0$ curves for various $C_2$ values with $C_2 > C_1$. Now if the particle has $C_1$ as its Jacobi constant and it is located on a $2 \ \Omega(x, y) = C_2$ curve, then the square of its velocity must be larger than zero on this curve, since $C_2 > C_1$. The particle with $C_1$ can not go outside the $C_1 = 2 \ \Omega$ oval since it would encounter points for which the zero velocity oval is associated with $C_3 < C_1$ which would require a decrease in the particle's $V^2$ from zero, i.e. it would require imaginary velocity .

In fact it is generally true that motion always takes place on that side of the zero velocity surface where the constant $\Omega$ lines have positive gradients, i.e. where C is increasing.

A formal computation gives the same result for the case under discussion. From equation (26) the zero velocity surface has the equation

(27) 
$$0 = \frac{2(1 - \mu)}{r_{10}} + 3\mu - C .$$

The square of the velocity at a point Q which is at a distance $r_1$ from $P_1$ is :

(28) 
$$V^2 = \frac{2(1 - \mu)}{r_1} + 3\mu - C$$

according to equation (25) .

Subtracting (27) from (28) we obtain

$$V^2 = 2(1 - \mu) \left( \frac{1}{r_1} - \frac{1}{r_{10}} \right) , \quad \text{from which}$$

it can be observed that if

$$r_1 < r_{10} \quad \text{then} \quad V^2 > 0 \quad \text{and} \quad \text{if}$$

$r_1 > r_{10}$ then $V^2 < 0$, i.e. in the first case the point $Q$ is inside the zero velocity surface (for $C$) and the positive $V^2$ value assures the possibility of motion, while in the second case $Q$ is outside and motion is not possible.

It is noted that for large $C$ a zero velocity oval can also be established around $P_2$ which is approximated by a circle of radius

$$(29) \qquad r_{20} = \frac{2\mu}{C - 3(1 - \mu)} \ .$$

Comparing (26) and (29) we see that

$$r_{10} > r_{20} \quad \text{if} \quad 0 \leq \mu < 1/2 \ ,$$

$$r_{20} > r_{10} \quad \text{if} \quad \frac{1}{2} < \mu \leq 1, \quad \text{and}$$

$$r_{10} = r_{20} \quad \text{if} \quad \mu = 1/2 \ .$$

In other words the larger one of the two principal masses is surrounded by a larger zero velocity oval. The two ovals around $P_1$ and $P_2$ are identical for $\mu = 1/2$.
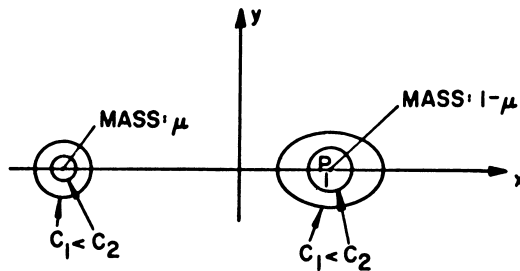
Figure 3 shows case I.



Figure 3. Zero Velocity Surfaces in Case I, for Large $C, \mu < \frac{1}{2}$ .

V. Szebehely

Since a particle moving within the zero velocity oval associated with either $m_1$ or $m_2$ will never leave this oval, case I is often referred to as the satellite case and hence the statement "once a satellite always a satellite."

### Case II.

Returning to equation (24) we observe that the $r_1 \to \infty$ (and consequently $r_2 \to \infty$) condition also results in large $\Omega$ and so in large C values .

Let $r_1 \simeq r_2 \simeq r$ for this case. Then

$$\Omega \simeq \frac{r^2}{2} \quad \text{and the velocity becomes}$$

(30) $$V^2 = r^2 - C.$$

The zero velocity ovals can be approximated by circles with radii

$$r_o = C.$$

These circles expand as C increases, therefore motion is possible outside the zero velocity line. If a particle moves outside such a zero velocity oval , it will never cross it, i.e. the particle will never change its planetoid or comet characteristics and it will never become a satellite. Hence, "once a comet, always a comet".

### Case III .

As C is decreased, the two ovals described under Case I will increase in size and will touch at $R_2$. At the same time the large outside oval of Case II shrinks. This is shown on Figure 4, where because of symmetry with respect to the x axis, only the upper half is drawn.
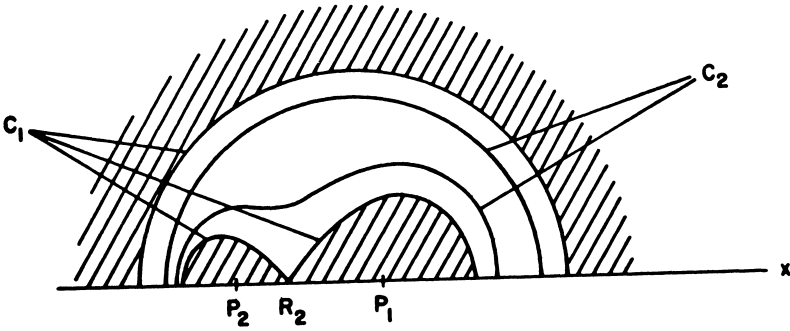
V. Szebehely



Figure 4 . Case III for $\mu < \frac{1}{2}$ , $C_1 > C_2$ .

As C is decreased further the two inside ovals unite (see the $C_2$ curve on Figure 4) permitting motion inside of the curve. In this case therefore satellite exchange might take place since particles are allowed to wander from the neighborhood of $P_1$ to the neighborhood of $P_2$. The practical significance of this case in connection with earth-lunar trajectories is apparent. It is noted, nevertheless, that an interchange of particles between the inside region of the $C_2$ curve enclosing $P_1$ and $P_2$ and the outside region of the large oval (corresponding to the same $C_2$) is still not permissible. The shaded areas show the regions of possible motions for $C_1$.

## Case IV.

Further decrease of C will cause the large outside oval and the inside figure (see the $C_2$ curve on Figure 4 ) to become in contact at $R_3$ or at $R_1$ depending on the value of $\mu$ . The contact point is established on the side of the smaller mass, i.e. for $\mu > \frac{1}{2}$ at $R_3$ and for $\mu < \frac{1}{2}$ at $R_1$ .

Figure 5 shows the permitted regions of motion and the zero velocity curves. Since the zero velocity curves are always symmetric to the

x axis, only the upper half is shown.



Figure 5 . Case IV for $\mu < 1/2$ .

## Case V.

As C is decreased further the figure opens up at $R_1$ (or at $R_3$) allowing a communication between the external and internal areas for the first time (see Figure 6 , Curve $C_1$). The next step is when the central portion (A B) of the horseshoe like forbidden area narrows and its width becomes zero (Curve $C_2$ on Figure 6) intersecting the x axis at $R_3$ ($R_1$) . Figure 6 shows the opening ($C_1$) and the narrowing processes ($C_2$) . In the shaded area motion is possible.

V. Szebehely



Figure 6 . Case V for $\mu < 1/2$ , $C_1 > C_2$ .

## Case VI.

As  C  is  further  decreased, the $C_2$ curve of  Figure  6  separates
at  $R_1$ and the two areas (one above  and the other  one below the x axis)
start shrinking toward $R_4$ and  $R_5$ . When C reaches its  minimum  value
(C = 3) the fobidden  areas  shrink  to zero. Figure 7 shows   this process.
When  C = 3 , motion is possible everywhere.



Figure 7 . Case   VI for $\mu < 1/2$,  $C_1 > C_2 > C_3 = 3$ .

A summary of the regions of motion is shown on Figure 8 for various C values. The figure which is correct only in a topological sense was constructed for $\mu < 1/2$ .



Figure 8 . Summary of the Hill Curves for $\mu < 1/2$, $C_2 C_2 > \ldots > C_8 = 3$

V. Szebehely

## 5 Regularization

### 5.1 General results.

The purpose of regularization is to eliminate the singularities which occur at the collision between the third body and either primary. One of the essential differences between the problems of celestial mechanics and problems of space dynamics is the collision and the close approach problem .

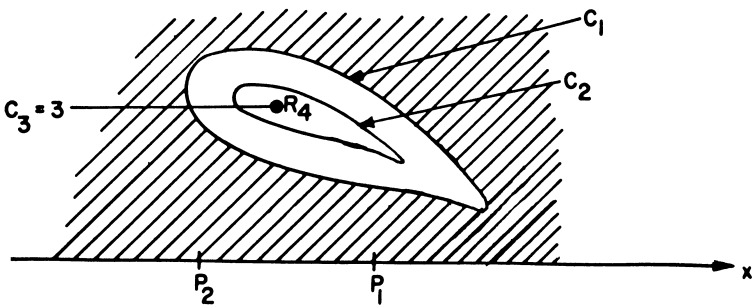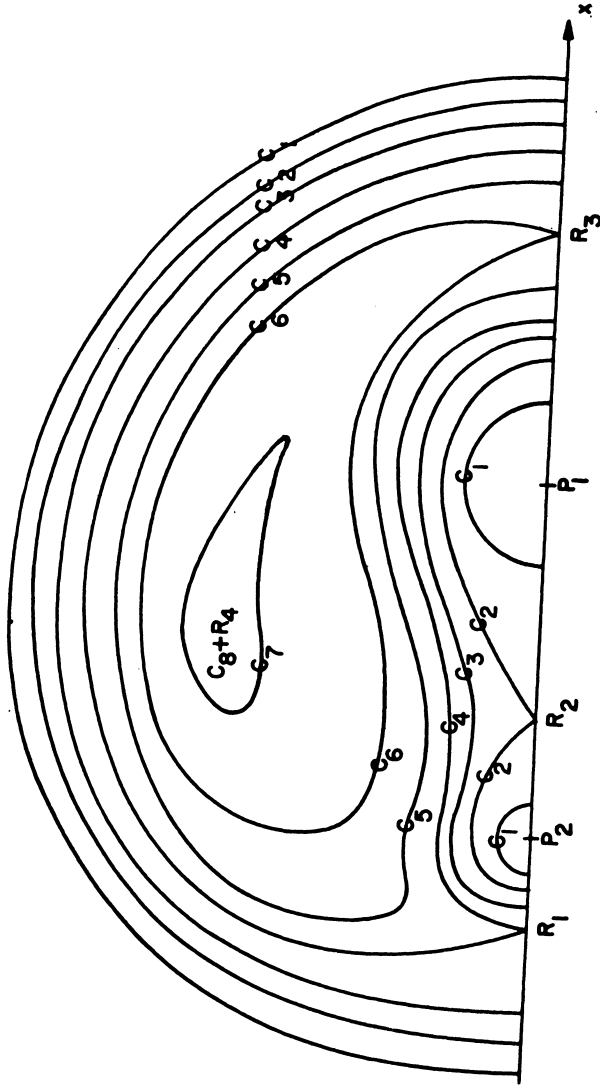Bodies participating in planetary and lunar theories of classical celestial mechanics do not experience collisions, whereas trajectories connecting the vicinities of celestial bodies are of vital interest in space dynamics. It is to be emphasized that close approaches are just as important and at the same time just as complicated to handle as are actual collisions between bodies. From a "practical" point of view, regularization is of importance since Earth-Moon trajectories connect two singularities, sometimes from the analytical, and sometimes from the numerical points of view. Finally, as we will see later, regularization will allow us to establish orbits near singularities ; an important advantage when studying the existence of periodic motions.

In order to eliminate the singularities mentioned above we consider complex transformations given by

$$(31) \qquad\qquad z = f(w)$$

where $z = x + i\,y$ and $w = u + i\,v$ and we find the equations of motion in the $w$ plane . Transformation of the dependent $(x, y)$ variables will be combined with the transformation of the independent variable $(t)$ and relations will be established between the original variables describing the motion $(x, y, t)$ in the $z$ plane and the new variable $(u, v, \tilde{t})$ which describe the motion in the $w$ plane.

The equations of motion (14) and the Jacobi integral (15) can be written as

(32) $$\ddot{z} + 2 i \dot{z} = \text{grad}_z U$$

(33)     and     $$\left| \dot{z} \right|^2 = 2 U ,$$

where $U \equiv \Omega - C / 2$ and the $\text{grad}_z$ operator is defined by

(37) $$\text{grad}_z U = U_x + i U_y .$$

The dot denotes time derivative with respect to the previously introduced nondimensional time (t), so for instance

$$\dot{z} = \frac{dx}{dt} + i \frac{dy}{dt} .$$

When the differential equation (32) representing the motion in the z plane is transferred to the w plane we consider the $z = z(t)$ solution in the z plane transferred to the $w = w(\bar{t})$ solution in the w plane.

According to this $z(t) = z \left\{ w \left[ \bar{t}(t) \right] \right\}$ and so

$$\dot{z} = \frac{dz}{dw} \frac{dw}{d\bar{t}} \frac{d\bar{t}}{dt}$$

or simply

(35) $$\dot{z} = f' w' \dot{\bar{t}}$$

which also defines the notation .

The second derivative becomes :

(36) $$\ddot{z} = f' w' \ddot{\bar{t}} + f' w'' (\dot{\bar{t}})^2 + f'' (w')^2 (\dot{\bar{t}})^2 .$$

The right hand member of equation (32) is transformed to :

(37) $$\text{grad}_z U = \frac{1}{\bar{f'}} \text{grad}_w U ,$$

V. Szebehely

where bar denotes conjugate .

Equation (37) is proved as follows :

(38) $\qquad \text{grad}_w U = U_x x_u + U_y y_u + i(U_x x_v + U_y y_v)$ .

The Cauchy-Riemann relations transform the imaginary part into

$$- U_x y_u + U_y x_u \, ,$$

which allows a rearrangement of terms in equation (38) :

$$\text{grad}_w U = (x_u - i y_u) (U_x + i U_y) , \qquad \text{q.e.d.}$$

Substituting $\dot{z}$ from (35) , $\ddot{z}$ from (36) and $\text{grad}_z U$ from (37) into (32) we obtain the equation of motion in the w plane :

(39) $\qquad w'' + w' \dfrac{\ddot{t} + 2 i \dot{t}}{(\dot{t})^2} = - \dfrac{(w')^2 f''}{f'} + \dfrac{\text{grad}_w U}{|f'|^2 (\dot{t})^2}$ ,

The Jacobi integral (33) takes the form :

(40) $\qquad |w'|^2 = \dfrac{2 U}{|f'|^2 (\dot{t})^2}$ .

Equation (39) can be written in a slightly different form by using equation (40) . We note that

(41) $\qquad \text{grad}_w (U |f'|^2) = U \, \text{grad}_w |f'|^2 + |f'|^2 \, \text{grad}_w U$ .

Here the first term on the right side can be written as

$$|w'|^2 |f'|^2 (\dot{t})^2 f' \, \bar{f}''$$

since

V. Szebehely

$$U = \frac{1}{2} \ |w'|^2 \ |f'|^2 \ (\dot{t})^2$$

from the Jacobi integral, and since

$$grad \ |f'|^2 = 2 \ f' \ \bar{f}''$$

This latter equation is proved as follows:

Leg $g(w)$ an analytic function. Then

$$g_u = \frac{dg}{dw} = \frac{1}{i} \ g_v \ \text{and} \ \bar{g}_u = \overline{(\frac{dg}{dw})} = - \frac{1}{i} \ \bar{g}_v$$

Furthermore,

$$grad_w \ |g|^2 = grad_w \ (g \ \bar{g}) \ \text{or}$$

$$= g \ grad_w \ \bar{g} + \bar{g} \ grad_w \ g.$$

Here

$$grad_w g = g_u + i \ g_v = \frac{dg}{dw} - \frac{dg}{dw} = o$$

and

$$grad_w \ \bar{g} = \bar{g}_u + i \ \bar{g}_v = 2(\overline{\frac{dg}{dw}}) = 2 \ \bar{g}'.$$

So

$$grad_w \ |g|^2 = 2 \ g \ \bar{g}'$$

and writing $f'$ for $g$ we have

$$grad_w |f'|^2 = 2f' \ \bar{f}'' \qquad q.e.d.$$

Solving (34) for $grad_w \ U$ and substituting this result into (39) we obtain

$$(42) \ w'' + w' \frac{\ddot{t} + 2 \ i\dot{t}}{(\dot{t})^2} = - \frac{(w')^2 f''}{f'} - \frac{|w'|^2 \ f' \ \bar{f}''}{|f'|^2} + \frac{grad_w \ (U|f'|^2)}{|f'|^4 \ (\dot{t})^2} \ .$$

If now we select the $\bar{t}(t)$ function so that

V. Szebehely

(43)
$$\dot{\bar{t}} = \frac{d\bar{t}}{dt} = \frac{1}{|f'|^2} \ ,$$

then

$$\ddot{\bar{t}} = \frac{d}{dt} \frac{1}{|f'|^2} = -\frac{1}{|f'|^6} \ (\bar{f}' \ f'' \ w' + f' \ \bar{f}'' \ \bar{w}') \ ,$$

and the equation of motion (42) becomes :

(44)
$$w'' + 2 i w' |f'|^2 = \mathrm{grad}_w \ (U|f'|^2) \ .$$

The Jacobi integral is

(45)
$$|w'|^2 = 2 |f'|^2 \ U \ .$$

It is remarked that in the z plane according to the Jacobi integral the square of the velocity $(v^2)$ goes to infinity as U, which in turn goes to infinity as $r_1$ or $r_2 \rightarrow o$. On the other hand for instance equation (42) suggests that $|f'|^2$ should be selected so that $U |f'|^2 \rightarrow$ finite as $r_1$ or $r_2 \rightarrow o$. This means that $v |f'| \rightarrow$ finite as $r_1$ or $r_2 \rightarrow o$.

The Jacobi integral in the w plane (40) can be written as

$$|w'|^2 = \frac{2 \ U |f'|^2}{\left( |f'|^2 \ \dot{\bar{t}} \right)^2} \ ,$$

and so if $|f'|^2$ is selected so that $U |f'|^2 \rightarrow$ finite at the singularities, it is also required to have

$$|f'|^2 \ \dot{\bar{t}} \rightarrow \text{finite and nonzero.}$$

This requirement is clearly satisfied if (43) holds, in fact

$$|f'|^2 \ \dot{\bar{t}} \rightarrow 1.$$

It is of further interest that equation (42) is nonlinear in w', which nonlinearity disappears when equation (43) is satisfied.

Summarizing we state :

V. Szebehely

The equations of motion in complex form can be written as

$$\ddot{z} + 2 i \dot{z} = \text{grad}_z \; U$$

and the Jacobi integral as

$$|\dot{z}|^2 = 2 U \quad.$$

Applying the transformations

$$z = f(w) \quad \text{and} \quad dt = |f'|^2 \; d\bar{t} \quad,$$

the equations of motion and the Jacobi integral become :

$$(44) \qquad \frac{d^2 w}{d\bar{t}^2} + 2 i |f'|^2 \frac{dw}{d\bar{t}} = \text{grad}_w \; (U |f'|^2),$$

$$(45) \qquad \left| \frac{dw}{d\bar{t}} \right|^2 = 2 U |f'|^2.$$

The real and imaginary parts are

$$(46) \qquad \frac{d^2 u}{d\bar{t}^2} - 2 |f'|^2 \frac{dy}{dt} = \Omega^\star_u$$

$$\frac{d^2 v}{d\bar{t}^2} + 2 |f'|^2 \frac{du}{d\bar{t}} = \Omega^\star_v$$

and

$$\left( \frac{du}{d\bar{t}} \right)^2 + \left( \frac{dv}{d\bar{t}} \right)^2 = 2 \Omega^\star$$

where

$$(48) \qquad \Omega^\star = U |f'|^2 = ( \Omega - C/2) |f'|^2.$$

In the following three short chapters we introduce three specific transformation functions as important examples. The reader will not find it too difficult to construct additional regularizing transformations of interest.

V. Szebehely

The significance of the Levi-Civita transformation is in its simplicity. In spite of its disadvantage - regularization of only one of the two singularities - it will be extremely useful in studying the motion near one of the singularities. The Hill model of the Earth moon system furnishes an excellent example for the usefulness of this transformation.

The second and third transformation functions regularize both singularities. The Birkhoff transformation accomplishes this by a rational fraction function, while the Thiele transformation uses a trigonometric function. The complications caused by the many valuedness of the Thiele transformation are compensated by the fact that the functions used are well tabulated . For analytical investigations the double valued Birkhoff transformation is recommended, while for detailed hand calculations the Thiele transformation seems to be better suited.

## 5.2 The Levi-Civita transformation

Let

(49)
$$z = f(w) = \mu + w^2$$

or

$$x - \mu = u^2 - v^2 \quad \text{and} \quad y = 2 u v.$$

The new time variable is related to the old by

$$dt = 4 \, d\bar{t} \, (u^2 + v^2) \,,$$

since

$$|f'|^2 = 4 |w|^2 \,.$$

The equations of motion become

$$\frac{d^2 u}{d\bar{t}^2} - 8(u^2 + v^2)\frac{dv}{d\bar{t}} = 4\left[(\Omega - \frac{C}{2})(u^2 + v^2)\right]_u$$

(50)

$$\frac{d^2 v}{d\bar{t}^2} + 8(u^2 + v^2)\frac{du}{d\bar{t}} = 4\left[(\Omega - \frac{C}{2})(u^2 + v^2)\right]_v$$

and the Jacobi integral is

(51) $\qquad \left(\frac{du}{d\bar{t}}\right)^2 + \left(\frac{dv}{d\bar{t}}\right)^2 = 8(\Omega - \frac{C}{2})(u^2 + v^2).$

The term on the right side becomes :

(52) $\quad (\Omega - \frac{C}{2})(u^2 + v^2) = \frac{1}{2}(u^2 + v^2)^3 + \frac{\mu^2}{2}(u^2 + v^2) + \mu(u^4 - v^4) +$

$$+ 1 - \mu - \frac{C}{2}(u^2 + v^2) + \frac{(u^2 + v^2)}{\left[(u^2 + v^2)^2 + 1 + 2(u^2 - v^2)\right]^{1/2}}.$$

The geometry of the transformation is summarized as follows (see Figure 9 ) :

1. The point $P_1$ located at $(\mu, o)$ in the z plane is transformed to the origin of the w plane (o, o).

2. The point $P_2$ located at $(\mu - 1, o)$ in the z plane goes into $w_{1,2} = \pm i.$

3. The origin of the z plane (o, o) goes over to $w_{1,2} = \pm i\sqrt{\mu}.$

4. The upper z plane (y>0) goes into the first quarter of the w plane.

5. The z plane goes into the upper half of the w plane (v>o) .

6. The lower half of the w plane (v<o) corresponds to the second leaf of the z plane.

7. The transformation takes the one leaf of the w plane into the two leaves of the z plane.

V. Szebehely



Figure 9. The Levi-Civita Transformation.

The significance of the transformation is that it eliminates the singularity at $P_1$. Since $w = o$ corresponds to $P_1$ we find from (52) that at this point

$$(\Omega - \frac{C}{2}) (u^2 + v^2) = 1 - \mu$$

and so the velocity at $w = o$ (according 51) will be finite, $2 \sqrt{2(1-\mu)}$.

The derivatives of $(\Omega - C/2) (u^2 + v^2)$ with respect to $u$ and $v$ are also finite at this point, in fact according to (50) at $P_1$

$$\frac{d^2 u}{d\bar{t}^2} = \frac{d^2 v}{d\bar{t}^2} = o .$$

It is noted that the Levi-Civita transformation eliminates only the singularity at $P_1$; at $P_2$ we still have a singularity, since

$$(\Omega - \frac{C}{2}) (u^2 + v^2) \longrightarrow \infty$$

as $w \longrightarrow \pm i.$ .

## 5.3 The Birkhoff transformation

Let

(53)
$$z = f(w) = \frac{w^2 + \mu(1-\mu)}{2w + 1 - 2\mu}$$

and so

(54)
$$\left| f'(w) \right| = \frac{\rho_1 \rho_2}{2 \rho_3^2} \quad , \quad \text{where}$$

$$\rho_1 = |w - \mu| = \sqrt{(u - \mu)^2 + v^2}$$

(55)
$$\rho_2 = |w + 1 - \mu| = \sqrt{(u + 1 - \mu)^2 + v^2}$$

$$\rho_3 = \frac{1}{2}|2w + 1 - 2\mu| = \sqrt{(u + \frac{1}{2} - \mu)^2 + v^2} \quad .$$

The equations of motion are

$$\frac{d^2 u}{dt^2} - \frac{1}{2}\left(\frac{\rho_1 \rho_2}{\rho_3^2}\right)^2 \frac{dv}{dt} = \Omega^*_u$$

$$\frac{d^2 v}{dt^2} + \frac{1}{2}\left(\frac{\rho_1 \rho_2}{\rho_3^2}\right)^2 \frac{du}{dt} = \Omega^*_v$$

where

$$\Omega^* = (\Omega - C/2)\left(\frac{\rho_1 \rho_2}{2 \rho_3^2}\right)^2$$

or

(57)
$$\Omega^* = \frac{\rho_1^2 \rho_2^2}{32 \rho_3^6}\left[(1-\mu)\rho_1^4 + \mu \rho_2^4\right] + \frac{1}{2\rho_3^3}\left[(1-\mu)\rho_2^2 + \mu \rho_1^2\right] - \frac{C}{2}\left(\frac{\rho_1 \rho_2}{2 \rho_3^2}\right)^2 \quad .$$

The geometry of the transformation is summarized as follows :

1. $P_1$ and $P_2$ are invariant points of the transformation.

V. Szebehely

2. The z plane's origin becomes $\pm$ $i\sqrt{\mu(1-\mu)}$ in the w plane.

3. The $w = \mu - \frac{1}{2}$ point corresponds to $z \longrightarrow \infty$.

4. The $w \rightarrow \infty$ point corresponds to $z \rightarrow \infty$.

At $P_1$ in the w plane $\rho_1 = 0$, $\rho_2 = 1$, $\rho_3 = \frac{1}{2}$ and $\Omega^* = 4(1-\mu)$, and so the velocity is $2\sqrt{2(1-\mu)}$. At $P_2$, $\rho_1 = 1$, $\rho_2 = 0$, $\rho_3 = \frac{1}{2}$ and $\Omega^* = 4\mu$, and so the velocity is $2\sqrt{2\mu}$. At $P_3$, $\rho_1 = \rho_2 = \frac{1}{2}$, $\rho_3 = 0$ and $\Omega^* \rightarrow \infty$ and so the velocity is not finite. The same statements apply to $\Omega_u$ and $\Omega_v$ and so we find that the Birkhoff transformation eliminates both singularities in the w plane. The only finite singularity (introduced by the transformation) in the w plane is the $w = \mu - \frac{1}{2}$ point, which however corresponds to $z \longrightarrow \infty$. The only other singularity in the w plane is $w \longrightarrow \infty$, which also corresponds to $z \longrightarrow \infty$. So we conclude that the equations of motion in the w plane are regular as long as the point is not rejected to infinity in the z plane.



Figure 10. The Birkhoff Transformation

## 5.4 The Thiele transformation

Let

V. Szebehely

$$z = f(w) = \mathcal{\mu} - \cos^2 \frac{w}{2}$$

(58)

$$\text{or} \quad z = \mathcal{\mu} - \frac{1}{2}(1 + \cos w) .$$

We observe that for $\mathcal{\mu} = \frac{1}{2}$, $z = \frac{1}{2}\cos w$, the simplicity of which might explain partly why the Copenhagen school in their $\mathcal{\mu} = \frac{1}{2}$ undertakings favored this transformation.

The real and imaginary parts of (58) are

$$x = \mathcal{\mu} - \frac{1}{2} - \frac{1}{2}(\cos u)(\mathrm{ch} v)$$

$$y = \frac{1}{2}(\sin u)(\mathrm{sh} v)$$

which equations show two advantages of this transformation. Firstly the (x y) net is transformed into confocal ellipses and hyperbolas and secondly the functions occuring are well tabulated and well suited for hand calculatons.

The principal values of the transformation give $P_1(\pi, o)$ and $P_2(o, o)$ in the w plane and observing that

$$\left| f'(w) \right|^2 = r_1 r_2$$

we see that both singularities are eliminated, since

$$\Omega^* = r_1 r_2 (\Omega - C/2) .$$

## 5.5 Numerical results for lunar trajectories .

We now use Birkhoff's transformation and show results of numerical integration of equations (56). These solution cannot be obtained by integrating the original equations of motion (see equation 14) because all solutions go through the singularities located at the primaries.

V. Szebehely



Fig.11 . A typical fast earth-to-moon trajectory with
         consecutive collisions.  Solid curves are
         referred to the fixed and rotating physical
         (x,y) systems, dotted curves to the regularized
         (u,v) system. $\mu$ = 1/82.45.



Fig. 12  A typical low energy earth-to-moon trajectory
         with consecutive collisions.  Solid curves
         are referred to the fixed and rotating physical
         (x,y) systems, dotted curves to the regularized
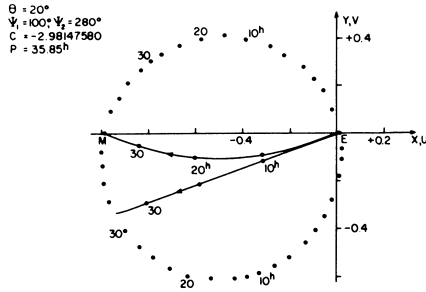         (u,v) system. $\mu$ = 1/82.45.

V. Szebehely



Fig. 13 A typical lofted earth-to-moon trajectory with
consecutive collisions. Solid curves are referred to the fixed and
rotating physical (x, y) systems, dotted curves to the regularized
(u, v) system. $\mathcal{M}$ = 1/82.45 .

The selection of a particular value of C (less than a critical value
given by the Jacobian integral) allows a trajectory to be constructed, which
connects $P_1$ and $P_2$ in the w plane. The totality of such trajectories, for
all possible values of C, form the underline{family} of orbits with consecutive colli-
sions. We define underline{a group} as a subset of the family in the following way.
Consider a trajectory connecting $P_1$ and $P_2$ which is obtained using a va-
lue of the Jacobian constant $C_1$. Changing $C_1$ to $C_2 = C_1 + \Delta C$ results
in another trajectory also connecting $P_1$ and $P_2$. If, as $\Delta C \rightarrow 0$, the
second trajectory approaches the first one, i.e., if the members of the
set can be generated by analytic continuation, then we speak about a group
of orbits with consecutive collisions.

A convenient dependent parameter is the pseudofiring angle (denoted
$\psi$, on Figs. 11-13) that determines the direction of the $w_1'$ complex ve-
locity vector at $P_1$. For a given value of C, this angle can be found
by solving a two-point boundary value problem in the w plane., using an

V. Szebehely

iteration procedure that begins with a two-body approximation. [1]
Consider, for instance, a very large negative value of C which corresponds
to a trajectory with a high relative energy. This follows from the Jaco-
bian integral, which at a fixed point in the physical plane gives

$$(\text{velocity})^2 = 2\,\Omega - C = \text{cons} - C$$

and therefore the velocity squared and the relative energy are very large
when C is large and negative . If the fixed point is not a singular point,
then the velocity becomes infinite as $C \longrightarrow -\infty$ .

Consider now a rectilinear traje-
ctory with consecutive collisions in the synodic, physical plane. This orbit
is possible only when $C = -\infty$ and the velocity is infinite along the entire
trajectory from one primary to the other. This orbit is just that part of the
x axis which lies between $P_1\,(\mu, 0)$ and $P_2(\mu-1, 0)$ . Increasing C from
$-\infty$ generates the group of consecutive collision trajectories that are descri-
bed in the next section.

The points of the curve of Fig. 14 represent the trajectories belonging
to the group mentioned previously. The insert shows the synodic coordinate
system $(x, y)$ and the firing angle $\theta$, in the synodic system , which is used
to present the results. The mass ratio is that for the earth-moon system,
$\mu = 1/82.45$. The $C = C(\theta)$ curve is asymptotic to the C axis, i.e., $\theta \not> 0$,
$C \longrightarrow -\infty$ as mentioned before. The orbits corresponding to the maximum
point of the curve $(C_o \simeq 2, \theta_o \simeq 60^o)$ seem to represent a change in the
basic characteristics of the members of the family, since for $\theta < \theta_o$,

1) Measuring all angles counterclockwise from the line $P_1 P_2$ , a trajectory in
the physical plane with a firing angle $\theta$ at $P_1$ is mapped into two trajecto-
ries in the regularized plane with pseudofiring angles of $\theta/2 \pm 90^o$ .

V. Szebehely

the maximum distance of the space probe from the earth is the same as the earth-moon distance. When $\theta > \theta_o$, the maximum distance of the probe from the earth becomes larger than the earth-moon distance.

At $\theta = 360^o$, the physical situation is identical to the $\theta = 0$ case; nevertheless, Fig. 14 indicates another value of C, i.e. $C(0) \neq C(2\pi)$. The explanation of this lies in



Fig. 14 . The Jacobian constants and firing angles for a group of simple , consecutive collision earth-to-moon trajectories in the framework of the restricted problem of three bodies. $\mu = 1/82.45$.

the definition of a "group" of trajectories. By analytic continuation, the curve was produced that is shown on Fig. 14 . At $\theta = 2\pi + \Delta\theta$, ( $\Delta\theta > 0$ ), another member of our group, can be obtained close to C = 1 , but of course it is also true that trajectories with very large negative values of C are also available.

It is of considerable practical and historical interest to mention that one of the points on the $C(\theta)$ curve of Fig. 14 can be (and as a matter of fact is) obtained by modifying a proposed circumlunar Apollo mission trajectory . The coordinates of this point are C = 1.805 and $\theta = 57^o.97$. After establishing this orbit, whose transit time is 79 hr as compared to the original orbit's transit time of 76 hs, the curve can be completed by va-

rying one of the parameters, C or    $\theta$, and differentially correcting the
other. Cowell  integrations of  the regularized differential equations (56) of
motion  were performed, using the Gauss-Jackson method of numerical in-
tegration with seventh differences and variable integration step sizes. Com-
putation of most of the orbits started at the center of the earth, but it
makes no difference at which primary the trajectories are started if the
moon moves in a retrograde  direction around the earth when the orbits a-
re started at the moon.

Although the trajectories are computed in a rotating coordinate system,
as is the usual procedure for the restricted  problem so that the Jacobian
integral can be  used for a check on the computations, it is easier to see
the motion of the space probe in a fixed frame . This is because 1) as it
will be shown in  the next section, our trajectories are slightly perturbed
conic sections in a  fixed frame, 2) it is customary to study lunar trajecto
ries in  a  fixed, geocentric frame, and  3) for several of the orbits, espe-
cially those with long transit times,  the path of the probe in the rotating
system is due essentially to the motion of the moon (i.e. , of the system)
and it is very difficult to separate the effects of the trivial geometry from
that of the dynamics.

Figure 11 shows a typical fast trajectory with C = -2.98. The probe
reaches the moon in  about  36 hrs, during which  time the moon has tra-
veled  $20^\circ$ . The four lines of Fig. 12 represent the same orbit in three dif-
ferent coordinate systems.

. The straight line with a slight curvature at its left end refers the
motion to  a fixed, geocentric system where the moon initially is at
(-1.0,0.0). The solid curve, connecting the points  E and M shows the tra-
jectory in the rotating system.  This is the way the orbit looks to an observer
moving with the earth-moon system. The dots represent the two branches of

V. Szebehely

the trajectory in the rotating, regularized (w) system. Midcourse times are denoted in hours along all four curves to give an idea of how the points are transformed in the coordinate systems.

As the energy is further decreased by increasing  C, the space probe takes longer to reach  the lunar distance,  and so the trajectory must point in a direction  farther  ahead of the moon. This way the curve is continued up to  where an  orbit is obtained corresponding  to a point in  the vicinity of  C = 2, $\theta = 60^O$ . Near the top of the curve of Fig. 14 are  such orbits, one of which is illustrated in  Fig. 12 . This low-energy orbit has  just enough energy to travel out far enough  so  that  the moon's attraction can pull the probe into a collision. For C  much  larger  than 2, the probe cannot suffi- ciently overcome the attraction of the earth, so it falls  back toward  it without reaching  the center of the moon.

It must be pointed out that  this does not  mean that a  probe with C > 2  cannot  go  from  the center of the earth to the center  of the moon; rather, it means  that such  an  orbit is  not contained in our group of trajectories. In  fact, an analysis of the zero velocity  curves for $\mu = 1/82.45$ shows  that  if  C < 3.2 orbits connecting the earth and the moon do  exist. However, as Egorov  points out, the probe may have to make several hun- dred revolutions around the earth,  in a highly eccentric ellipse with its apogee distance very slowly increased by the moon's perturbations,  before reaching the moon.

The  rest of the curve of Fig. 14 ($60^O < \theta \leq 360^O$) represents trajectories on which the space probe, still traveling essentially on geocentric, rectili- near ellipses, reaches the lunar radius before the moon has traveled  $\theta^O$ . Since  C < 2 , the probe reaches an apogee  greater than  1.0 , and it collides with the moon as it falls back toward the earth. Figure 13 is a typical trajec-

V. Szebehely

tory of this kind. With  C = 1.43, the probe leaves the earth  in  a  direc-
tion just opposite to that of the moon ( $\theta$ = 180$^{\text{o}}$), passes the moon's  di-
stance in about  three days, and reaches apogee (r = 1.431) in  about eight
days. After this, it falls back toward the earth  and reaches the center of
the moon half a sidereal month after leaving the earth. Also, this figure
shows, for the sophisticated reader who compares the two dotted trajecto-
ries in the regularized system, that since the Birkhoff  transformation maps
the region of infinity to the region midway between the two primaries
($\mu$ - $\frac{1}{2}$, 0), one of the dotted curves goes near this point. For the other
regularized trajectory, the region of infinity is unchanged by the transforma-
tion.

To the right  of the point corresponding to the Fig. 13, trajectories
on the C ($\theta$) curve are higher energy orbits with lower values of C. These
have greater apogee distances and must therefore  start in a direction still
farther in front of the moon. At the end of this group is the trajectory that
starts the probe heading right for the moon, but by the time it gets to
the lunar distance, the moon has moved on about  25$^{\text{o}}$ so the probe reaches
apogee (r = 2.1)  and then falls back to hit the moon about one month after
leaving the earth.

The trajectories  discussed are by no means the only ones that pass
through the centers of both  the primaries,and  they are only  a small per-
centage of the entire family. When the curve of Fig. 14 is continued to  the
right  ($\theta$> 2 $\pi$), it branches, and the trajectories  for the approximate ran-
ge  360$^{\text{o}}$ < $\theta$  < 420$^{\text{o}}$ encounter strong perturbations from  the moon as they
pass  it  on  the way out to apogee. If these  trajectories are allowed to
collide with the moon  on  the way out, the curve shown on  Fig. 14 is,  of
course,  just repeated. For  $\theta$ > > 2 $\pi$ there exist orbits that are associa-
ted with  several revulutions of the moon around  the earth before collision
and also orbits that are characterized by several close encounters of the

V. Szebehely

of the probe with the moon or the earth  before collision.

   When new numerical methods are introduced an important question
of considerable practical importance must  often be  investigated. This is
the question related to the speed and the associated economy of the new
procedure. In the examples mentioned in this section , this question comes
up quite naturally since before the numerical integration is performed, one
must  transform to  the new variables and after completion of the integra-
tion one must revert to the original (physical) variables. We also shuld
note that  the transformed differential equations of motion using the new va-
riables are of more complicated forms than the original equations. The eli-
mination of the singularities simplifies the equations from a mathematical
point of view but the algebraic representations become more complicated. Is
it then practical to become involved with a more complicated system and
with transformations  or is it more economical to integrale the equations in
their original form ? If we deal with collision trajectories, the answer
must be on  the side of regularization, since without this collision orbits
can not be computed. In these experiments one computes orbits with and
without regularization and one compares the results regarding speed, econo-
my and accuracy .

   Such  comparisons are shown in the following two tables. As the title
indicates we first study a 72 hour trajectory connecting the surface of the
earth with the surface of the Moon . The numerical integration required
1260  steps using the original variables and only 340 steps in the regulari-
zed system. The step-size used is variable in this method of numerical
integration and  the program automatically determines the largest allova-
ble step-size for a given local accuracy. Only two such step-size changes
are required  in the regularized system and 15 when the original varia-
bles are used  as shown in the second  line of the first table. The total
computational  times  also  favor the regularized method  as  shown
on the third  line .    The last four lines in  the tables  show the errors

V. Szebehely

obtained when the reversibility condition is checked . We integrate from
the Earth to the Moon and  then we reverse the integration and record
the  differences  between the  initial positions  and velocities on one hand and.
the final return positions and velocities on the other hand. As seen in the
tables, the regularized system exhibits consistently smaller errors than
the original system.

The second table gives the results for a slightly different  trajectory
utilizing a parking  orbit around the Earth.

| | Physical | Regularized |
|---|---|---|
| No. of integration steps | 1260 | 340 |
| No. of step size changes | 15 | 2 |
| Computation time (secs.) | 40 | 16 |
| $\delta x_s$ | $-138 \times 10^{-13}$ | $-396 \times 10^{-15}$ |
| $\delta y_s$ | $-105 \times 10^{-13}$ | $-77 \times 10^{-17}$ |
| $\delta \dot{x}_s$ | $414 \times 10^{-11}$ | $244 \times 10^{-13}$ |
| $\delta \dot{y}_s$ | $405 \times 10^{-11}$ | $122 \times 10^{-12}$ |

Errors in a 72-hour trajectory computed from the surface of the earth to the surface of the moon

| | Physical | Regularized |
|---|---|---|
| No. of integration steps | 1423 | 489 |
| No. of step size changes | 12 | 3 |
| Computation time (secs.) | 43 | 22 |
| $\delta x_s$ | $-645 \times 10^{-14}$ | $-128 \times 10^{-14}$ |
| $\delta y_s$ | $574 \times 10^{-14}$ | $215 \times 10^{-14}$ |
| $\delta \dot{x}_s$ | $-164 \times 10^{-11}$ | $-292 \times 10^{-12}$ |
| $\delta \dot{y}_s$ | $-201 \times 10^{-11}$ | $-717 \times 10^{-12}$ |

Errors in a 72-hour trajectory computed from a parking orbit around the earth to the surface of the moon

V. Szebehely

6. Additional Numerical results

The behavior of a dynamical system may be characterized in a variety of ways. If the dynamical system is non-integrable an immediate approach is to attempt to find the totality of certain special orbits. We must restrict the search to periodic, almost - periodic or asymptotic orbits since orbits of more general characteristics can not be established on a computer. The greatest popularity is enjoyed by periodic orbits and certain higly specialized orbits of engineering interest. We may look at the numerical as an experimental approach and it must complement the theoretical considerations.

This experimental approach consists of the computation of a large number of trajectories, preferably with systematically varied initial conditions. One of the pioneers of this approach was G. DARWIN (1897) who, without the aid of digital computers and using a 1 : 10 mass ratio for the primaries, computed and analysed several classes of orbits. Similar remarks apply to. E. STRÖMGREN (1935) usign 1 : 1 mass ratio and giving an even more complete set of orbits. The unfortunate aspect of these large undertakings is that since at the mass ratio 1:27 the character of certain orbits change, extension of DARWIN's and of STRÖMGREN'S results to cases of importance in dynamical astronomy or in space dynamics cannot be made readily since the mass ratios of interest are considerably smaller than 1:27 . The study of periodic orbits made by MOULTON (1920) resulted in highly specialized orbits. None of these studies covered the entire four dimensional manifold and even in their respective ranges of investigations many questions remained unresolved.

Recent work along the line of experimental establishment of the totality of orbits in the restricted problem fall in two categories. The first group of activities is oriented along the lines of celestial mechanics and

V. Szebehely

concerns itself with the study of sets of special orbits applicable to problems of celestial mechanics. MESSAGE (1959) and RABE (1961, 1962) are outstanding examples, the first finding periodic orbits for the planetary (exterior) case and using the Sun-Jupiter system as the model, the second establishing periodic orbits numerically around the equilateral libration points .

This latter work is of special interest since THURING (1951) offered a proof of analytical nature for the non-existence of periodic orbits around the libration points. existence question of periodic orbits around these points is still open since neither THÜRING's analytical nor RABE's numerical "proofs" are entirely acceptable. LAGRANGE (1772)[*] has shown the existence of stationary solutions (in the rotating frame) for the restricted problem - in fact for the general problem of three bodies these solutions also exist in a more general form - namely the so-called collinear and equilateral configurations. The standard, linearized treatment of these solutions (cf. e.g. MOULTON, 1914) reveals instability for the latter provided the mass ratio of the primaries is less than 1 : 27 . Periodic solutions can be found for the linearized equilateral case without difficulty, therefore, the above mentioned RABE-THÜRING controversy is a question of non-linear effects. The astronomical and cosmogonical implications of this problem are numerous. The well established Trojan group of asteroids at the equilateral point of the Sun-Jupiter system and the conditionally established Kordylewski clouds for the Earth-Moon system are examples.

Just as those whose prime interest is in the astronomical applications concern themselves with the satellite and planetary orbits, workers in space dynamics study orbits which connect the vicinity of one primary with the neighborhood of the other. These studies are of very limited extent in spi-

---

[*] The historically oriented reader will not fail to observe that the above date coincides with EULER's publication of his second lunar theory and therefore also with the beginning of the history of the restricted problem.

V. Szebehely

te of the ease with which high speed digital computers furnish trajectories. This phenomenon is explained by the "practical" orientation of the investigators, who are, by necessity, limiting themself to the goal of obtaining an

acceptable trajectory for the specific mission on hand. The result is a formidable mosaic of points in the four dimensional manifold of the initial conditions, with high density in certain isolated regions where trajectories were needed for specific purposes and with large, completely unexplored regions. To attempt to undertake a systematic investigation regarding the totality of possible orbits of the restricted problem has not only theoretical significance but it has great practical ramifications . Without comprehensive information on all possible trajectories, the orbit selection process cannot be elevated from its present state of trials and errors method. Since free trajectories can be, and in practical cases are combined by orbit modification techniques, it cannot be expected that out of the very large number of possible combinations a purely experimental process with all its presently existing limitations can select the most desirable orbit.

It is inconceivable to refer to all the papers at this point which report on experimentally established trajectories, for two reasons. Firstly, these publications seldom are included in the recognized open literature and secondly they offer such small ranges of the initial conditions that no general and systematic conclusions can be drawn from them. HUANG, R. NEWTON, Arenstorf, Broucke, etc. present a special set of periodic orbits which enclose and approach the vicinity of both primaries . The range of initial conditions is limited; nevertheless, this work is of definite interest for space mechanics.

The above-mentioned Lagrangian or libration points have also space mechanics applications. Besides the use of space probes at the Earth-Moon equilateral points as solar flare observational stations or for other possible scientific or space purposes, an interesting mathematical relation might be

V. Szebehely

mentioned between the linearized equations of motion of such probes and one type of the rendez-vous problem. The category of rendez-vous problems in modern space mechanics encompasses a very large area since orbits with with predetermined end conditions can always be looked upon as rendez-vous. When the meeting of two, gravitationally non-interacting space probes is analyzed and their relative motion in a central force field is studied, the pertinent differential equations, from a mathematical point of view, become identical with the linearized equations of motion in the vicinity of the equi-lateral libtration points. This special rendez-vous problem, therefore, is so-lidly embedded in classical celestial mechanics. The other large group of rendez-vous problems might be referred to as two-point boundary value pro-blems in which the initial conditions are to be established so that certain mission requirements be satisfied by the orbit.

The above-review of the experimental approach to orbit classification is concluded by observing that a thorough restudy of the classical results obtained by DARWIN and STRÖMGREN, combined with studies of smaller than 1: 27 mass ratios and including astronomically significant cases as well as well as types of orbits of interest in space explorations is long overdue. This undertaking, of course, will have to be combined with theo-retical guiding principles, possible along topological lines.

## REFERENCES

1 . Birkhoff, G.D. , "Sur le problème restreint des trois corps, " Two memoirs , both published in Annali della R. Scuola Normale Superiore di Pisa . The first in series 2, Vol. 4, pp. 267-306 (1935) ; the second in series 2, Vol. 5, pp. 1-42 (1936) .

2.      Birkhoff, G.D., "Dynamical systems with two degrees of freedom, " Trans. Am. Math. Soc., Vol. 18, pp. 199-300 (1917) .

3.      Birkhoff, G.D. , "The restricted problem of three bodies, " Rendiconti de Circolo Matematico di Palermo, Vol. 39, pp. 1-70 (1915) .

4. Brouwer, D. and G. Clemence, "Methods of Celestial Mechanics", Academic Press, 1961.

5.      Brouwer, L.E.J., "Über eineindeutige, stetige Transformationen von Flächen in sich, " Math. Annalen, Vol. 69, pp. 176-180 , (1910) .

6.      Charlier, C.L. , "Die mechanik des Himmels, " Leipzig, von Veit Co. , First Vol. 1902, Second Vol. 1907 .

7.      Darwin, G.H. Acta Math. , Vol. 21, pp. 99-242 (1897) .

8.      Hill, G.W. , "Researches in the lunar theory, " Am. J. of Math. Vol. 1, pp. 5-26, 129-147, 245-260 (1878) .

9.      Klose, A., "Topologische Dynamik der interplanetaren Massen, " Vierteljahrsschrift der Astronomischen Gesellschaft, vol. 67, pp. 61-102. (1932) .

10.     Levi-Civita, T. , Acta Math. , Vol. 42 , pp. 99-144 (1919) .

11.     Levi-Civita, T., " Sur la résolution qualitative du problème restreint de trois corps, " Acta Mathematica, Vol. 30, pp. 305-327 (1906) .

12.     Levi-Civita, T., " Traiettorie singolari ed urti nel problema ristretto dei corpi, " Annali di Matematica, Ser. 3, Vol. 9, pp. 1-32, (1904) .

13.     Levi-Civita, T. , Ann. di Mat. , Ser. 3, Vol. 5. pp. 221-309, (1901) .

14.     Moulton, F.R., Proc. Math Congr., Cambridge, England, Vol. 2, pp. 182-187 (1913) ; also, Periodic Orbits, Carnegie Inst. Wash. (1920) .

15.     Poincaré, H., "Les méthodes nouvelles de la mécanique céleste, " Paris , 1892, 1893, 1899.

16.     Strömgren, E. , "Connaisance actuelle des orbites dans le problème des trois corps, "Bull. Astr. (2) , Vol. 9, pp. 87-130 (1935), and Publ. Kbh. Obs. No 100, pp. 1-44 (1935) .

17.     Sundman, C.F. , Acta Math. , Vol. 36, pp. 105-192 (1912).

18.    Szebehely, V., "Theory of Orbits, The restricted problem of three  ׁodies" , Academic Press, N.Y. and London, 1967.

19.    Thiele, T.N. Astron, Nachrichten Vol. 138, pp. 1-10, (1895) .

20    Wintner, A., "The analytical foundations of celestial mechanics", Princeton Univ. Press, 1947 .

CENTRO INTERNAZIONALE MATEMATICO ESTIVO

(C.I.M.E.)

G. A. WILKINS

THE ANALYSIS OF THE OBSERVATIONS OF

THE SATELLITES OF MARS

# THE ANALYSIS OF THE OBSERVATIONS OF
# THE SATELLITES OF MARS

by

G. A. Wilkins

(Royal Greenwich Observatory)

Summary

The practical aspects of a general procedure for the analysis of
positional observations of satellites (or other astronomical bodies) are
discussed and illustrated by reference to experience gained during an analysis
for the satellites of Mars.

## 1. Introduction

The principal purpose of this paper is to discuss the practical aspects
of a general procedure for the analysis of observations of the positions of
satellites (or other astronomical bodies) and to illustrate the use of such a
procedure by reference to a current analysis for the satellites of Mars. The
paper differs in emphasis from the original two lectures in that the discussion
of the procedure has been extended while the discussion of the results has been
shortened since they have already been given elsewhere (Wilkins, 1967).

There are two main reasons for analysing observations of the apparent
positions of the satellites of a planet. Firstly, from the characteristics
of the orbits we can deduce information about the gravitational field of the
planet, and hence about the total mass of the planet and its departures from
spherical symmetry, i.e. we can deduce the 'dynamical shape' and principal axis
of the planet for comparison with the apparent shape and axis of rotation.
Secondly, we can compare the observations with the predictions given by a
theoretical model of the satellite system and so determine whether the model
provides an adequate representation; if it does not then we can look for other
perturbing effects that may not have been properly taken into account, e.g.,
resonance effects caused by mutual perturbations of the satellites. Once an

adequate model has been developed we can investigate the stability of the system and speculate on its origin and future, e.g., whether any of the satellites have been recently captured or will eventually escape.
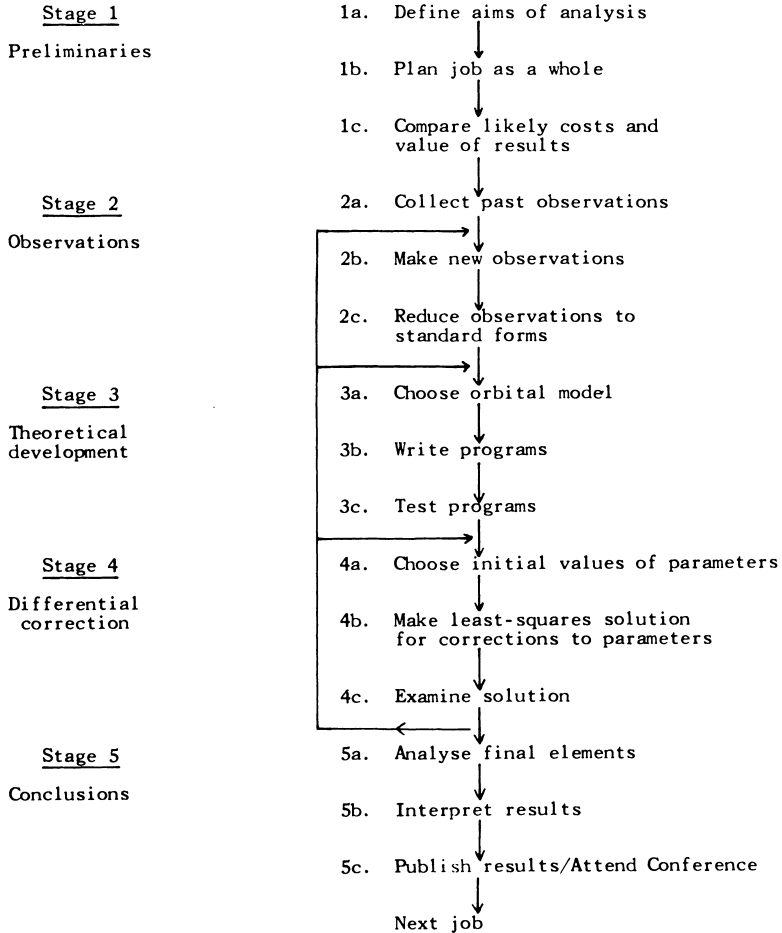
The details of the method of analysis will naturally vary according to the complexity of the model that is required to represent the observations. The procedure described here should be of fairly general application in a wider field than that of the study of satellite orbits, although much more sophisticated methods that demand very powerful computing systems would be appropriate for some applications.

The satellites of Mars form a comparatively simple system and the main difficulty arises from the lack of suitable observations. The initial reason for making a new analysis was to investigate the existence, or otherwise, of secular accelerations, such as could be caused by drag or tidal reactions, in the motions of the satellites. The results have, however, proved to be of interest in connection with the properties of Mars itself.

## 2. A general procedure for the analysis of observations

We may divide the task of analysing observations into a number of stages as indicated in the flow-chart in Figure 1. For simplicity the chart has been drawn to suggest that the various stages are carried out in sequence, but in practice we may carry out some of the stages (for example, 2 and 3) concurrently.

Figure 1.        Flow-chart of general procedure for analysis of observations

Stage 1

Preliminaries

1a.  Define aims of analysis

1b.  Plan job as a whole

1c.  Compare likely costs and
value of results

Stage 2

Observations

2a.  Collect past observations

2b.  Make new observations

2c.  Reduce observations to
standard forms

Stage 3

Theoretical
development

3a.  Choose orbital model

3b.  Write programs

3c.  Test programs

Stage 4

Differential
correction

4a.  Choose initial values of parameters

4b.  Make least-squares solution
for corrections to parameters

4c.  Examine solution

Stage 5

Conclusions

5a.  Analyse final elements

5b.  Interpret results

5c.  Publish results/Attend Conference

Next job

Stage 1. Preliminaries. The successful execution of a job may be largely determined by the amount and quality of the effort put into the preliminary planning of the work. First of all we must define as precisely as possible the aims of the job in order to make certain that the following stages can be designed in such a way that these aims can be realised. The aims should not be restricted too closely, but should allow for the possible investigation at some later time of aspects that are subsidiary to the principal aim - after all it is unexpected effects that often prove to be of the greatest interest and value.

Once the aims have been defined our next step is to plan the execution in sufficient detail from beginning to end to ensure that no fundamental difficulties are overlooked and that the output of each stage is likely to be adequate in both extent and quality for use as input to the succeeding stages. It is, in fact, often useful to consider the stages in the reverse order to that in which they will be executed. It is almost invariably wasteful to embark on the execution of the early stages until the job has been studied in fair detail right up to the final stage.

Once this planning stage is completed we should be able to make reasonable estimates of the resources (of men and equipment) that will be required to complete the job. This is the time at which the decision to go ahead, or to redefine the aims, or to abandon the project, should be taken. In most cases the decision to go ahead will depend on the approval of a higher authority for the required expenditure, but even if such approval is not necessary, or has been given at an earlier stage, a conscious decision should always be made at this time.

Stage 2. Observations. In celestial mechanics it is rarely possible for one person or group to make all the necessary observations and to analyse them as well, and so the next step is to collect together all the suitable observations that have been made. The preliminary studies should have indicated the period over which observations are needed and the type and quality of observations to

be used; they may even have shown that new observations must be made if the aims of the analysis are to be achieved.

The collection of observations will normally involve a literature search in order to find out where the observations have been published. The volumes of the Astronomischer Jahresbericht provide an invaluable guide for the years 1899 onwards, and a systematic search covering theoretical as well as observational papers will usually serve to draw attention to most available observations. Observatories often publish routine observations only in special series of volumes and these must normally be examined separately.

Even when the observations have been traced a considerable amount of detective work may still be required in order that the precise basis of the printed values may be determined. For example, even the time scale or coordinate system may be ambiguous, and quite frequently observers fail to state what corrections have been applied to the measures. If such information cannot be determined from other evidence the observations must be rejected, or at best used with low weight. When an observational value is transcribed from the published form to that to be used in the analysis it should be accompanied by auxiliary indications that define precisely the significance of the value. It is, of course, necessary to check carefully that any errors made in the transcription process are corrected.

The next step is to 'reduce' each of the observations to one of a number of possible standard forms that are acceptable to the main stage of the solution. These standard forms should be chosen so as to keep this reduction stage as short as is consistent with avoiding the necessity for a long or complex input stage in the main program. Again great care is necessary to ensure that no errors are made in the reduction to standard form, especially if the corrections are calculated or applied by hand. Herrick (1960) has given a useful review of the separation between the functions of the observer and the orbit analyst in the reduction of observations; he also discusses the choice of reference systems for the analysis.

<u>Stage 3. Theoretical development</u>. The first step in the theoretical development involves the choice of a suitable model for the system to be analysed. The form of the model largely depends on what perturbations are considered to be significant in relation to the accuracy of the observations. This choice involves also the selection of a suitable set of parameters, or elements, whose values are to be determined more precisely by a differential-correction procedure. These parameters specify not only the initial conditions for the orbit being analysed but also other relevant characteristics of the system (e.g., the masses of the perturbing bodies). The model is defined by a set of equations, or procedures (e.g. numerical integration), by which we can compute for any instant of time ($t$) a predicted value of each observational quantity ($\alpha$) given adopted values of the parameters ($p_1$, $p_2$, ... $p_n$). This means that for each type of observation we have a functional relationship of the form

$$\alpha = \alpha\,(p_1,\ p_2,\ \ldots\ p_n;\ t), \tag{1}$$

and for each observation $\alpha_i$ ($i = 1$, 2, ... $N$) we can compute the residual

$$\Delta\alpha_i = (\alpha_i)_{obs} - \alpha_i\,(p_1,\ p_2,\ \ldots\ p_n;\ t_i) \tag{2}$$

These residuals can be considered to be due partly to observational error and partly to the differences between the adopted and true values of the parameters. If the chosen model does not adequately represent the actual physical system there will also be a contribution to each residual from this cause. However, if we ignore this possibility and assume that the observational errors are randomly distributed, we can adopt the principle of least squares and find values of the parameters that make the sum of the squares of the residuals $\Delta\alpha_i$ a minimum.

Let $$p_j{}' = p_j + \Delta p_j,$$

and let

$$\Delta'\alpha_i = (\alpha_i)_{obs} - \alpha_i\,(p_1{}',\ p_2{}',\ \ldots\ p_n{}';\ t_i).$$

Then, if each $\Delta p_j$ is sufficiently small, we have $N$ equations of condition of the form

$$\Delta'\alpha_i = \Delta\alpha_i - \sum_{j=1}^{n} \frac{\partial\alpha_i}{\partial p_j}\,\Delta p_j,\ (i = 1,\ 2,\ \ldots\ N). \tag{3}$$

Thus the problem is reduced to one of solving in the least-squares sense a system
of $N$ linear equations in $n$ unknowns. This can be done by standard techniques by
forming from the conditional equations a set of $n$ normal equations which can then
be solved directly. It will be noticed that in addition to computing the predicted
values of $x_i$ it is necessary to compute the $n$ partial derivatives $\partial x_i / \partial p_j$ for each
observation. This can be done either by first analytically differentiating the
fundamental equations (1) or by direct numerical differentiation, i.e., by re-
evaluating the equations (1) after introducing small increments to each of the para-
meters in turn. The choice will depend on the complexity of the system of equations
that represent the model and on the characteristics of the computer that is to be
used for the job. Except for comparatively simple systems direct numerical
differentiation, which involves much less programming effort, is now normally used.
It should be noticed, however, that the increments must be sufficiently small that
second-order terms are negligible, but sufficiently large that the changes in $x_i$
due to differences in the propagation of rounding-errors in the computation are
much smaller than those due directly to the changes in the parameters.

When the necessary analytical formulae have been developed the corresponding
computer program must be written and tested to make certain that it does correspond
precisely to the specification of the model. If analytical formulae have been
developed for the partial derivatives, they should be tested by direct numerical
differentiation. Particular care should be taken over units since it is unlikely
that the observations, the theoretical formulae and the computer routines will all
use the same system.

Stage 4. Differential corrections. The initial values for the parameters will
normally correspond as closely as possible to the results of previous investigations.
The least-squares solution will then give differential corrections and their
standard errors as indicated by the internal evidence of the observations.
The corrected set of parameters must be checked by using them as a new set of
initial parameters and then verifying that the resultant further differential
corrections are not significant in comparison with the standard errors. Several

such iterations may be required if the assumption that the second-order terms
are negligible is not valid for the initial choice of elements. The convergence
(or otherwise) of the successive solutions gives a good guide to the validity
of the final solution. The final residuals should always be examined since
departures from a random distribution will indicate the presence of deficiencies
in the model or systematic errors in the observations - but the converse is not
true!

    If the solution fails to converge satisfactorily, or if some of the standard
errors of the differential corrections are very high, it may be desirable to
modify the model or the selection of parameters. For example, if the inclination
of an orbit is small the longitude of the node and the argument of pericentre
(measured from the node) will be poorly determined, but their sum, the longitude
of the pericentre, may be well determined. The presence of other less obvious
correlations is indicated by the presence of large off-diagonal elements in the
inverse of the matrix of the normal equations. It may even be found that one
or more of the aims of the investigation must be abandoned until new observations
have been made. The main point to be emphasised is that, except in routine
cases, reliance should not be placed on any standard computer routine that
automatically iterates to a least-squares solution, especially if the number of
unknowns is large.

Stage 5. Conclusion. If the preliminary planning has been done well, the analysis
of the final elements should be straightforward, i.e., given the set of parameters
that best represent the observational data, and their standard errors, it should
then be possible to compute the corresponding values of other related parameters.
The interpretation of the results may, of course, give rise to a completely new and
interesting set of problems and may itself show where further investigation is
required.

3.   The analysis for the satellites of Mars

    The general procedure for the analysis of observations that has been

described in section 2 may be illustrated by the methods used and the results so far obtained in the current analysis of the observations of Phobos and Deimos. We first discuss the background to the analysis and then consider each of the stages of the analysis in turn.

Background.   The two satellites of Mars were discovered in August 1877 by Hall (1878), who used the new 26-inch refractor at the U.S. Naval Observatory, and who named them Phobos and Deimos (after Fear and Terror, the companions of the God of War).  They move in nearly circular orbits which lie close to the equatorial plane of Mars.  They are difficult to observe except at favourable oppositions, since they are faint, fast-moving, and always lie comparatively close to the planet.  The radii of the orbits of Phobos and Deimos are only 2.7 and 6.9 times the radius of Mars.  Their periods of revolution are $7^h$ $39^m$ and $30^h$ $18^m$, while the period of rotation of Mars is $24^h$ $37^m$ - as viewed from the planet, Phobos would therefore rise in the west and set in the east.  Oppositions of Mars occur at intervals of about 2 years and 2 months, and favourable oppositions, when the geocentric distance of Mars may be as small as 0.38 a.u., occur at intervals of 15 or 17 years, the next being in August 1971.  The apparent magnitudes of the satellites at opposition are about 11 and 12, compared with a magnitude of, say, -2 for Mars itself; at such times the angular radius of Mars is only about 10".  Attempts to find a third satellite have so far proved unsuccessful.

    The bulk of the observations of the positions of the satellites are of the apparent coordinates of each satellite with respect to the centre of Mars, and have been made visually with filar micrometers.  The measurements may have been made directly in terms of position angle and distance or indirectly by measuring in rectangular coordinates the distances from the tangents to the limbs of the planet.  There are a much smaller number of observations of the positions of one satellite relative to the other; such observations are easier to make, particularly by photography, and should be less liable to systematic errors.  Since the equatorial plane of Mars is inclined at about 24° to the ecliptic the apparent orbits of the satellites do not always intersect the apparent disc of the planet.

At some oppositions, however, the phenomena of occultation, transit, and eclipse occur, but as yet, owing to the difficulties of observation, there are no published series of such observations.

The satellites were observed fairly regularly during favourable oppositions up to 1909 and definitive sets of orbital elements were determined by Struve (1911). Further observations were made at the U.S. Naval Observatory during the 1920's, and were used by Burton (1929) to obtain corrections to Struve's elements. The predictions in The Astronomical Ephemeris are, however, still based on Struve's elements since they have proved to be adequate for finding purposes up to 1956 when the last-known series of observations was made. Both Struve and Burton determined independent sets of elements for each satellite, but Woolard (1944) drew attention to the inconsistencies between the elements for the two satellites. Shortly afterward Sharpless (1945) claimed to have shown that Phobos has a definite secular acceleration and that Deimos might have a small secular deceleration. Since no satisfactory explanation of such accelerations could be found Dr. G. M. Clemence suggested in 1957 that I should make a new analysis of the observations to see if Sharpless' suggestion could be confirmed and if the inconsistencies between the sets of elements could be reduced.

It is of interest to note that Hall, Burton, Woolard, Sharpless and Clemence were all on the staff of the U.S. Naval Observatory, and that I started the analysis while on a short tour of duty at the Observatory in Washington. I made a preliminary, hurried solution using the IBM 650 computer at the Yale University Observatory early in 1958, but the results were inconclusive. Other work and the lack of a suitable computer at Herstmonceux resulted in the job being left unfinished for a number of years. Provisional results were obtained in 1964 and 1966 by using an IBM 7090 computer in London. The programs were written in Fortran IV and have recently been adapted to run on the ICT 1909 computer at Herstmonceux.

Stage 1. The analysis was started with the principal aim of confirming, or otherwise, Sharpless' suggestion that the satellites had secular accelerations.

In all previous studies the observational data were treated in a piecemeal fashion. The observations made at each opposition by each observer, or group of observers, were analysed separately to determine the set of elements for an elliptic orbit that best represented those observations. These sets of elements were then combined in order to determine the secular motions of the elements. Sharpless found that when this was done it was not possible to fit new observations made in 1939 and 1941 by using a simple linear expression for the mean longitude, but that a quadratic term was required.

This procedure is obviously open to the objection that the final elements were only indirectly fitted to the original observations and so the inconsistencies obtained might be due to the differing methods used to reduce and analyse the observations. It was therefore decided to collect together all the available observations for each satellite and to fit them by simple, independent models in which some of the elliptic elements, including the mean motion in longitude, were assumed to have secular changes. This would give new estimates of the secular accelerations for comparison with Sharpless' values and the other parameters could be tested for consistency.

It is now clear that such an approach is unsatisfactory for a number of reasons, Firstly, it introduces more parameters to be determined than there are independent parameters in the system. This means that the apparent fit to the observations may be improved at the expense of introducing possible inconsistencies between the parameters. Secondly, the parameters which are now of prime interest (i.e., the constants defining the gravitational field of Mars) are not explicitly determined, but have to be obtained indirectly from the inconsistent values of the orbital parameters. Thirdly, it is not possible to include directly in the solutions the observations of the positions of one satellite relative to the other, even though they are probably more accurate than the observations relative to Mars. Fourthly, the form of the solution assumes that any changes in the mean motion are regular (and can be adequately represented by a quadratic term), whereas it is possible that an irregular variation is present.

It is, of course, easy to make such criticisms in retrospect but the point
to be emphasised is that some, if not all of them, could have been avoided if
more consideration had been given to the procedure to be used in stage 5 before
stages 2 and 3 were planned and executed.

Stage 2. The only real difficulty in the second stage arose from the fact that
some observers failed to indicate the precise significance of the results that
were published (e.g., whether corrections for the phase of the planet had been
applied). The details of some 2800 measures were punched on to cards. In order
to simplify the input stage of the main program the times of observation were
reduced to the Ephemeris Time scale (by adopting a standard table of values of
E.T. - U.T.), allowance was made for light-time, and the apparent coordinates
of Mars were associated with each observation.

The accuracy with which secular motions can be determined is, of course,
improved as the time-span of the observations increases. It was at first
thought that the observations made from 1877 - 1956 would be adequate for the
job in hand, but it is now clear that additional observations would be useful
both in order to put closer limits on any irregularities in the mean motions and
so that the gravitational field parameters for Mars can be determined more
accurately for use in space-probe missions and other studies.

Stage 3. The orbital model that was used in the analysis was based on the
assumption that each satellite can be considered to be moving in an ellipse whose
elements are subject only to secular (i.e. cumulative) changes. The periodic
perturbations due to the Sun and to the non-sphericity of Mars were neglected,
as were also the perturbations of the other planets and the mutual perturbations
of the two satellites.

For each orbit the principal secular perturbation is due to the oblateness
of Mars and causes a retrograde motion of the normal to the orbital plane about
the principal axis of the planet; the principal perturbation due to the Sun is a
similar motion about the normal to the plane of the orbit of Mars around the
Sun. The net result is that the normal to the orbital plane precesses about an

axis lying between the principal axis of the planet and the normal to the orbit
of Mars. This axis is close to the principal axis of the planet since the
perturbations due to the oblateness are very much greater than those due to the
Sun. The precession is such that the orbital plane maintains a constant
inclination to the so-called Laplacian plane, and the line of intersection of
the two planes (i.e., the line of nodes) regresses (i.e., moves in the opposite
sense to the motion of the satellites) at a constant rate. In addition the line
of apsides in the orbital plane advances at nearly the same rate. The precessions
of the Laplacian planes for the two satellites have been ignored.

The orbital model for each satellite was specified by 12 elements since the
position of the Laplacian plane and the secular rates were treated as independent
parameters. The fundamental orbital elements were referred to the reference
system defined by the equator and equinox of 1950.0, but for each observation
it was necessary to reduce the elements to the equator of date before calculating
the apparent position of the satellite. In view of the simplicity of the model
it was possible to evaluate the partial derivatives of the observed quantities
with respect to the parameters directly from analytical expressions. The normal
equations were solved by the use of a standard matrix inversion subroutine.

The experience gained in testing and using the program indicated the need
for a number of modifications to the initial simple least-squares solution
that was planned. For example, since we are looking for secular effects we are
endeavouring to find the coefficients in expressions of the form

$$x(t_i) = a_0 + b_0 (t_i - t_0) + c_0 (t_i - t_0)^2, \qquad (4)$$

where $t_0$ is some arbitrary date (e.g., near 1950.0). When, as in this case, the
majority of observations are at one end of the range the normal equations are
found to be very badly conditioned, i.e., they are difficult to solve correctly
and the standard errors of the solutions are high. The ill-conditioning can,
however, be overcome by re-expressing equation (4) in the form

$$x(t_i) = a_m + b_m (t_i - t_m) + c_m (t_i - t_m)^2, \qquad (5)$$

where $t_m$ is the (weighted) mean date of the observations. This mean date is not
known in advance but an adequate value can be obtained after one or two iterations.
It was, therefore, necessary to introduce extra program steps to compute $a_m$, $b_m$, $c_m$
from $a_0$, $b_0$, $c_0$ and vice-versa.

The initial program was also modified so that a solution could be made for
any weighted selection of the full set of 12 parameters. The weighting factors
are useful in improving the conditioning of the normal equations since it is
desirable that the elements on the principal diagonal are of the same order of
magnitude. Since the inclination ($i$) of the orbital plane to the Laplacian plane
and the eccentricity ($e$) of each orbit are small an automatic weighting was
introduced by solving for $\sin i \, dN$ and $e \, dP$ where $dN$ and $dP$ are the required
corrections to the longitudes of the node and pericentre.

Stage 4. So far some six solutions using all the observations up to 1928 have
been made; some of them were simple iterations, but in others the effects of
eliminating some of the unknowns, or changing the limits for the rejection of
observations, were tried. The adopted solution is given in Table 1. The errors
that are quoted are based on a consideration of both the formal standard errors
given by the least-squares solutions and the self-consistency of the various
solutions.

It was found that no significant improvement in fit to the observations was
obtained when a secular acceleration term was included, even though the formal
solution for Phobos indicated a non-zero value. The reason for this ambiguous
result is almost certainly that the bulk of the observations were made before
1909 and the remainder were made in the 1920's. In an attempt to resolve the
difficulty the program was adapted so that the later observations made in 1941 and
1956 could be used and so that some indication of the residual errors in orbital
longitude over the whole period could be obtained.

The few observations available for 1941 and 1951 were entirely in the form
of positions of Phobos relative to Deimos, and the assumption was made that both
sets of orbital elements were correct except that the mean longitude of Phobos

might be subject to fluctuation. The program was used to find for each
observation the value of the zero of longitude that would give the minimum
residual; there was a considerable scatter but the general trend was plain for
each group of observations.

The earlier observations fall naturally into groups, corresponding to
measures of one coordinate (e.g. position angle), by one observer, for one
opposition. By assuming that all the other elements were correct it was
possible to estimate directly the value of the zero of mean longitude that
would give the best fit in the least-squares sense to each group of observations.
The differences between these estimates and the original value could show the
presence of a steady rate of change of the mean motion or a fluctuation in rate
such as would occur if air drag were significant but varied greatly with solar
activity. No such effect was found; the scatter (see Wilkins, 1967) was greater
than that given by Sharpless, but this may be merely due to the larger number of
groups of observations that were used. The 1956 observations gave a much smaller
residual (2°) than the value (6°) that would be given by an acceleration of the
magnitude suggested by Sharpless.

As far as is known there have been no direct observations of the satellites
since 1956; attempts to see and photograph them at Herstmonceux have failed.
Rakos (1965), however, attempted to observe a series of partial eclipses of Phobos
early in 1965; a preliminary analysis suggests that the mean longitude must be
altered by about 30° if the present model is to be reconciled with the observations.
The reality of such a large effect must be in doubt until it has been confirmed
by further observations or the observational data and the subsequent analysis have
been subjected to close and independent scrutiny.

Stage 5. The final set of parameters are substantially the same as those obtained
by Struve and Burton. They can be analysed to determine the mass ($m$) of Mars
and the quantity $J_2 R_e^2$, where $J_2$ is the coefficient of the second-harmonic in the
expression for the gravitational potential of Mars and $R_e$ is the equatorial radius
of Mars. Adoption of a value of $R_e$ (which is not well determined by visual

observations) leads to a value of $J_2$ and hence to the 'dynamical flattening'
($f$) and mean density of the planet. (See Wilkins, 1967, for the derivation and
discussion of the results.) The value for the mass of Mars is consistent with
the more accurate value that has recently been obtained (Gaugler, 1967) from an
analysis of the orbit of the space-probe Mariner IV. The value of the flattening
is only half that obtained by direct visual observations of the planet (Dollfus,
1967); the reason for this discrepancy is not known. (For example, see Öpik, 1962,
and Runcorn, 1967).

Conclusions. In spite of the elapsed time since its inception, this analysis of
the observations of the satellites of Mars cannot be regarded as definitive and
so an improved analysis is being made. The principal change will be the use of
a much more sophisticated model of the system in which the positions of the
Laplacian planes and the secular motions of the motions of the node and pericentre
will be computed theoretically in terms of the parameters $J_2$ and higher-order
coefficients, and in which the principal periodic perturbations will also be
included. It is thereby hoped to reduce the root-mean-square residual for the
visual observations from 0.''55 to a value of, say, 0.''3 such as is obtained if a
group of observations for a single opposition are fitted by a simple elliptic
orbit. Correspondingly the standard errors of the derived parameters and the
scatter in the. residuals in orbital longitude should be reduced, and it should
be possible to remove the present uncertainty about the secular motions of these
elusive satellites.

Table 1

## Orbital parameters for satellites of Mars (from analysis made in 1964)

| Elements for equator and equinox of 1950.0 | Phobos | Deimos |
|---|---|---|
| Epoch | J.D. 2414600.5 | J.D. 2414800.5 |
| Longitude of node of fixed plane | 46°9 ± 0°1 | 46°40 ± 0.05 |
| Inclination of fixed plane to equator | 37°57 ± 0°07 | 36°64 ± 0.03 |
| Argument of node of orbital plane at epoch | 80° ± 5° | 358°3 ± 0°9 |
| Daily motion of node of orbital plane | -0°438 ± 0°001 | -0°0180 ± 0°0003 |
| Inclination of orbital plane to fixed plane | 0°9 ± 0°1 | 1°80 ± 0°02 |
| Mean longitude at epoch | 227°1 ± 0°1 | 333°87 ± 0°03 |
| Daily mean motion in longitude | 1128°8443 ± 0°0001 | 285°16192 ± 0°00001 |
| Secular acceleration in longitude | – | – |
| Longitude of pericentre at epoch | 211° ± 3° | (300° ± 20°)* |
| Daily motion of pericentre | +0°436 ± 0°001 | (+0°016 ± 0°003)* |
| Apparent semi-major axis at unit distance | 12."91 ±0."01 | 32."36 ± 0."01 |
| Eccentricity of orbit | 0.018 ± 0.001 | 0.0 ± 0.0003 |
| Mean epoch of observations | 1899.0 | 1898.7 |
| Number of observations (1877 - 1928) | 1440 | 1300 |
| Standard error of observations of unit weight | 0."52 | 0."55 |

* Formal values only; the final value of the eccentricity of the orbit of Deimos
  was much less than the standard error and so the orbit may be assumed to be
  circular.

## References

Burton, H. E. 1929.  Elements of the orbits of the satellites of Mars. *A.J.*, 39,
155 - 164.

Dollfus, A. 1967.   In *The Mantles of the Earth and Terrestrial Planets*, ed.
S. K. Runcorn, London: John Wiley. pp. 85 - 92.

Gaugler, E. A. 1967. In *The Mantles of the Earth and Terrestrial Planets*, ed.
S. K. Runcorn.  London: John Wiley. p. 397.

Hall, A. 1878. Observations and orbits of the satellites of Mars with data for
ephemerides in 1879.   Washington: Government Printing Office.

Herrick, S. 1960.   Observation requirements for precision orbit determination.
*Univ. of California Los Angeles Astronomical Papers*, 3,
23 - 40.

Öpik, E. J. 1962.  Surface properties of Mars and Venus.  *Progress in Astronautical
Sciences*, Vol. 1, 298 - 307.

Rakos, K. D. 1965. Results communicated privately by R. L. Duncombe.

Runcorn, S. K. 1967. In *The Mantles of the Earth and Terrestrial Planets*, ed.
S. K. Runcorn.  London: John Wiley. pp. 425 - 430.

Sharpless, B. P. 1945. Secular acceleration in the longitudes of the satellites
of Mars.  *A.J.*, 51, 185 - 6.

Struve, H. 1911.  Uber die Lage der Marsachse und die Konstanten im Marssystem.
*Sitz. Ber. Preuss. Akad. Wiss., Sitz. Phys. Math. Classe vom
30 Nov. 1911*.

Wilkins, G. A. 1967. In *The Mantles of the Earth and Terrestrial Planets*, ed.
S. K. Runcorn, London: John Wiley, pp. 77 - 84.

Woolard, E. W. 1944. The secular perturbations of the satellites of Mars.
*A.J.*, 51, 33 - 6, 1944.