

ORIGIN OF THE MOON— THE COLLISION HYPOTHESIS

D. J. Stevenson

Division of Geological and Planetary Sciences, California Institute of Technology, Pasadena, California 91125

INTRODUCTION

In 1871, during his presidential address to the British Association in Edinburgh, Sir William Thompson (later Lord Kelvin) discussed the impact of two Earth-like bodies, asserting that “when two great masses come into collision in space, it is certain that a large part of each is melted” [see Arrhenius (1908, p. 218) for the complete quotation]. Although he did not go on to speculate about lunar origin, it must have been remarkable to see one of the creators of the bastions of nineteenth century conservative science discuss such an apocalyptic event and the debris issuing from it. It is equally remarkable that until recently, lunar origin myths have usually centered around three possibilities (fission, capture, and binary accretion) that exclude any important role for giant impacts. The Origin of the Moon Conference held in Kona, Hawaii, on October 13–16, 1984, saw a megaimpact hypothesis of lunar origin emerge as a strong contender, not because of any dramatic new development or infusion of data, but because the hypothesis was given serious and sustained attention for the first time. The resulting bandwagon has picked up speed (and some have hastened to jump aboard). Most significantly, efforts have been made to simulate giant impacts using three-dimensional hydrodynamic codes on supercomputers. Although all this effort is promising, a sober reflection on the problem after two years suggests that a lot more work is needed. It is not yet clear whether the collision hypothesis satisfies the observational facts.

A definition is in order. By the impact or collision hypothesis, I mean any theory that seeks to derive the Moon-forming material from the outcome of one or more collisions between the Earth and other Sun-

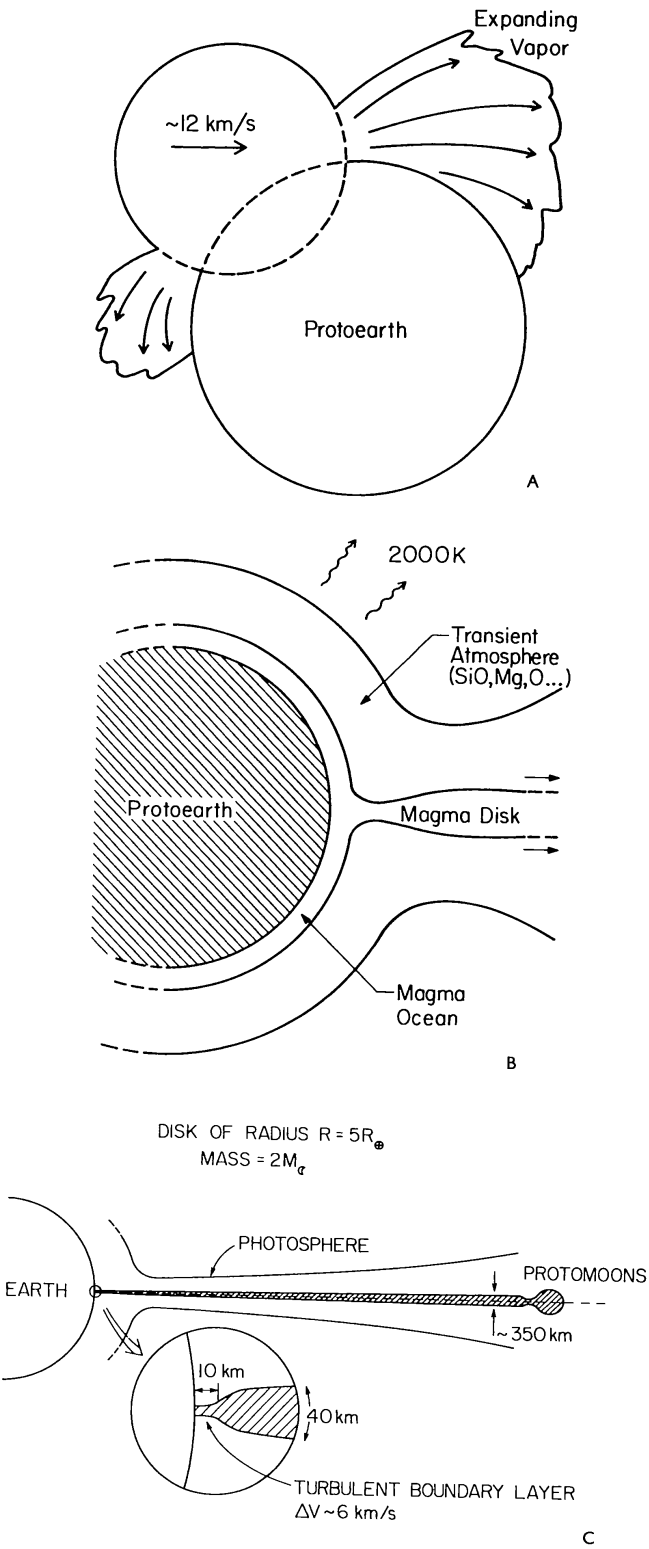
orbiting bodies. For reasons that will become apparent, the impacting body or bodies must be large—larger than the Moon and perhaps even larger than Mars. Notice that this definition does *not* assume that the formation of the Moon was necessarily a singular event. Among proponents of the collision hypothesis, there are those who think that a single event overwhelmingly dominated and those who think that a few (or even many) impact events were needed. There are even versions of the collision hypothesis that are not very different from extreme versions of one of the alternative origin scenarios of capture, fission, and binary accretion!

The mainstream view (if one can be said to exist) of one, or at most a few, oblique impacts ejecting material into Earth orbit owes its origin largely to the ideas of Hartmann & Davis (1975) and Cameron & Ward (1976). Hartmann & Davis were among the first to emphasize the possibility of very large “planetesimals” as part of the population of impacting bodies during Earth accretion, a possibility that is consistent with computer simulations by Wetherill (1980, 1985). Cameron & Ward were the first to assess the physical outcome of very large impacts and the important issues posed by the angular momentum budget; this has led to recent numerical simulations (Cameron 1985b, Benz et al 1986a,b). As often happens in this kind of interdisciplinary science, many other characters had (and have) roles to play and are introduced in due course. The biggest danger in this review, however, lies not in the possibility of inadequate attribution but in attempting to review an area where many of the important calculations are in progress or have not yet been done (perhaps cannot be done). This review proceeds by advancing 10 propositions that I believe embody the most important issues confronting the theory. These propositions may or may not be true, but they form a framework for asking the right questions and for organizing the presentation. Figure 1 summarizes the main features of the impact hypothesis embodied in these propositions.

TEN PROPOSITIONS

1. *The other theories of lunar origin are inadequate.* Fission is dynamically implausible; capture and binary accretion have both dynamical and cosmochemical problems.

Figure 1 Possible sequence of events leading to lunar formation. (A) A giant impact causes an expanding flow of liquid and vapor away from the impact site, carrying part of the angular momentum of the projectile. (B) A disk has formed, consisting of a liquid sublayer and a gas “atmosphere.” (C) The disk has spread, and protomoons form at its extremity.



2. *Large impacts occurred during planetary accumulation.* There is no good reason to suppose that the masses of the impacting bodies were always much less than the masses of the resulting planets.
3. *Large impacts lead to qualitatively different outcomes than small impacts at the same velocity. In particular, orbital injection of material may occur.* The failure of scaling arises mainly because of the essential three dimensionality of large (oblique) impacts. Gradients in the gravitational field become important and hydrodynamic effects (especially pressure gradients) can operate over distances comparable to the radius of the Earth. These factors may be essential to the issue of orbital injection efficiency.
4. *The Earth's escape velocity is neither much less than nor much greater than the impact velocity needed for substantial vaporization of rock.* This is important because impact velocities are comparable to escape velocity, and vaporization is needed if pressure-gradient acceleration plays an essential role in orbital injection.
5. *The postimpact Earth may be like a brown dwarf star for about 100 yr.* An immense amount of energy may be dumped into the Earth, causing a transient global magma ocean and a transient atmosphere of silicate vapor. The Earth may radiate from an extended photosphere at $T \sim 2000$ K ; such radiation would have been detectable by infrared astronomers orbiting nearby stars! The consequences for Earth evolution may provide a test for the theory.
6. *The material injected into orbit is very hot and probably consists of two phases (liquid with gas bubbles or 'foam'). It may form a disk and spread rapidly.* Evolution times of a two-phase disk are very short ; material may spread out to and beyond the Roche limit in 10^2 – 10^3 yr, before it has a chance to cool and solidify.
7. *The Moon-forming material may be derived primarily from the mantles of the projectile and the Earth and therefore iron poor.* Both projectile and target should be hot and well-differentiated bodies. Although not yet convincingly demonstrated by computer simulations, the outer layers of each may contribute most to the protomoon(s). The relative contributions of projectile and target are very uncertain.
8. *The newborn Moon or protomoon(s) are hot because of impact energy rather than because of their gravitational self-energy.* This is because everything happens so quickly that radiative and convective cooling are inadequate to allow solidification prior to Moon formation. The newborn Moon is then at least partially molten ; this is relevant to lunar thermal history.
9. *Despite incongruent vaporization and localized differentiation, the system may be almost "closed."* The net hydrodynamic outflow to

infinity may be small, although important for devolatilization. However, major-element fractionation may be unimportant, except to the extent that *physical* separation (e.g. core formation or separation of liquid iron from liquid silicate) dictates the material available for the Moon.

10. *One Moon arises because the largest impact was the last important impact for supplying lunar material.* Accretion is hierarchical, and the largest impact may occur late. The resulting protomoon undergoes more rapid tidal evolution than smaller, earlier protomoons and sweeps up these smaller bodies.

Although not included in this list, an additional proposition is that the giant-impact hypothesis may have implications for comparative planetary science, including Earth-Venus dissimilarities, the absence of a substantial moon around Mars, and the obliquities of Saturn, Uranus, and Neptune. These are assessed briefly.

The next section deals with observational data: What properties of the Moon do we seek to explain? This is followed by sections on planetary accretion and on the various lunar origin scenarios (defense of Proposition 1). The section following these (Physics of Large Impacts) motivates Propositions 3–6, and subsequent sections discuss recent and ongoing numerical simulations, efforts to understand the postimpact evolution, and the chemical aspects of the hypothesis.

OBSERVATIONAL DATA

A nice review of these data has been recently provided by Wood (1986). This section, although slanted differently, is accordingly kept to the bare essentials.

Mass and Angular Momentum

The lunar mass, one eightieth of the Earth's mass, seems an anomalously large fraction of the total Earth-Moon mass compared with other planet-satellite systems. However, this is of questionable significance. There are few terrestrial planets with which to compare, and two of these (Mercury, Venus) have been greatly affected by solar tides. It is probably inappropriate to compare the Earth-Moon system with outer solar system planets because the latter may have different satellite origins (e.g. Stevenson et al 1986) and certainly have a large gas component. Actually, the ratio of the total Jovian satellite system mass to the mass of the *core* of Jupiter (or the rock and ice component of Jupiter) is probably not much

less than the Moon to Earth mass ratio. Consequently, there is no strong reason to suppose that the Earth-Moon system is “special.”

For similar reasons, the angular momentum budget of the Earth-Moon system may not be particularly special or anomalous. We have simply too few other systems with which to compare. Nevertheless, this angular momentum (equivalent to that of an Earth rotating with a ~ 4 -hr period) is a very important constraint on origin models that is not always readily satisfied.

Bulk Chemistry

The Moon's mean density is $3.344 \pm 0.002 \text{ g cm}^{-3}$ (Bills & Ferrari 1977). If one constructs a body of lunar mass assuming cosmic Mg/Fe and Mg/Si ratios and appropriate combined oxygen, then the resulting mean density depends on the form of the iron (metal or oxide or substituting for Mg in silicates), but it is always at least 10% greater than observed. The only reasonable way to explain this discrepancy is by reducing the iron content by at least a factor of three relative to the cosmic abundance. This argument is independent of, but supported by, evidence that the Moon either has no iron-rich core or has a core that is at most ~ 400 km in radius (corresponding to $\lesssim 2\%$ of the mass). Constraints on the lunar core arise from a variety of arguments, many of which are geophysical (see Newsom 1984).

The depletion of iron is not in dispute, but its interpretation is still unclear because there is no consensus on the similarity of lunar bulk chemistry and Earth mantle composition. A refractory-rich Moon has had many advocates (Anderson 1972, Cameron 1972, Taylor & Jakes 1974, Ganapathy & Anders 1974), but the idea that the Moon is iron poor because of noncondensation of metallic iron in the region of lunar formation seems implausible, based on both cosmochemical and petrological considerations (Ringwood 1979). Trace-element abundances may be more diagnostic (see below), but they are controversial. Comparisons of the Moon with the *whole* mantle of the Earth are even more uncertain because the lower mantle is not yet well characterized (but see, for example, Jeanloz & Thompson 1983).

Volatile Depletion

Lunar soils and rocks are strongly depleted in volatiles, even more so than the Earth's mantle. It is widely believed that this depletion results from some highly energetic process accompanying lunar formation. Certainly, it cannot be attributed solely to the small size of the Moon, since at least one solar system body of comparable mass and size (Io) has a substantial volatile component. It is likely, however, that the volatile depletion of the Moon is not due to a single event and at least partially predates lunar formation (Taylor 1986). It is also possible that the degree of volatile

depletion has been overestimated: There may be a significant component of volatile material deep within the Moon. We should be wary of a lunar origin scenario that extracts volatiles *too* efficiently. We should also be careful about terminology; molecules or atoms that are volatile in one thermodynamic or chemical environment may be involatile in another. Accordingly, each physical scenario has to be modeled directly and not loosely categorized merely by the degree of volatile depletion.

Trace Elements

Most attention has been given to siderophile elements (those that preferentially partition into a metallic iron phase and are therefore believed to be concentrated in planetary cores). Most, but by no means all, siderophiles are also volatile. Ringwood (1979) and Wänke and coworkers (Wänke & Dreibus 1986) have advanced the view that the similarity of siderophile *patterns* in the Earth and Moon argues for coGenesis (e.g. the derivation of the Moon from the Earth's mantle). However, there are differences in the patterns (Drake 1983, Kreutzberger et al 1986), so the inferences are unclear. In fact, it would be surprising if the patterns were extremely similar, since there may have been some further differentiation (e.g. lunar core formation) and further accretion of Sun-orbiting debris after the main lunar formation event(s).

The trace-element questions are both complex and important. It is probably a fair assessment at present to say that the data argue neither conclusively for nor conclusively against deriving lunar matter from the Earth's mantle or from the mantle of a body that has undergone geochemical differentiation similar to that of the Earth. Much more discussion of these issues can be found in several chapters of Hartmann et al (1986).

Primordial High Temperatures

The anorthositic highlands of the Moon have been frequently attributed to fractional crystallization from a primordial magma ocean (reviewed by Warren 1985). In fact, the magma ocean concept arose more out of geochemical convenience than from compelling physical or chemical arguments. Nevertheless, essentially global melting or extensive partial melting to a depth $\gtrsim 100$ km seems to be needed, although this melting need not have been uniform in space and time. It is questionable whether this could have been achieved from the gravitational energy of lunar formation (Wetherill 1975, Kaula 1979). The "magma ocean" is therefore a significant constraint on lunar formation models.

Orbital Evolution

The gradual increase of the Earth-Moon distance has long been known and has been directly measured by astronomical methods (Lambeck 1980,

Ch. 10). However, the backward extrapolation in time of tidal theory is highly uncertain, even aside from the well-known “problem” that the current specific tidal dissipation (or reciprocal of the quality factor Q) is higher than the average over geologic time. It is relatively easy to construct a model that brings the Moon back to the Roche limit at $\sim 4.5 \times 10^9$ yr before present (e.g. Walker et al 1983; see also Lambeck 1986, Walker & Zahnle 1986), but it is not possible to predict the configuration of this primordial orbit. In particular, a near-equatorial orbit cannot be excluded, even though specific calculations (e.g. Goldreich 1966) suggest an inclined orbit. This lack of prediction arises because of incomplete knowledge of tidal dissipation in both the Earth and the Moon, and because of the possibility that one or both bodies were significantly affected by later impacts.

PLANETARY ACCUMULATION

Modern ideas of solar system formation are guided by astrophysical observations and our improved understanding of planetary properties, and they are aided by the rapid recent developments in computing facilities. There are two very distinct views of planetary formation that currently receive the most attention. The less popular view, strongly advocated by Cameron (1985a), involves “giant, gaseous protoplanets” that arise through gravitational instability of the solar nebula. In the terrestrial zone, these bodies are believed to lose their gaseous component, leaving the rock and iron nuclei that are the building blocks of the terrestrial planets. Cameron asserts that there would have been many more nuclei than the current number of terrestrial planets, so that the subsequent evolution must have involved giant impacts between these (Mars-sized?) building blocks. Cameron’s theory has received insufficient quantitative development to be assessed fully, but it appears to have some potential problems for explaining both the terrestrial and the giant planets (see, for example, Stevenson et al 1986). For our present purposes, it is sufficient to note that this theory probably provides the kind of impact that Cameron wishes to invoke for lunar origin (Cameron 1985b, 1986).

The more popular view of planetary formation assumes that the growth of solid bodies is by a sequential hierarchical process: condensation of dust grains, aggregation into larger clumps, formation of \sim kilometer-sized planetesimals (possibly by gravitational instability), and progressive growth of larger planetesimals leading eventually to the planets. There are two main versions of this scenario: gas free and gas rich (see review by Wetherill 1980). In the gas-free scenario, which is the focus of the work done by Safronov (1966, 1969) and carried on by Wetherill, most of at

least the later stages of terrestrial planetary accumulation occur in the absence of any primordial, hydrogen-rich nebula. The gas-rich scenario is mainly the work of Hayashi and collaborators (Hayashi et al 1985) and assumes that there is still sufficient gas, even at the later stages, to affect dynamical and thermal conditions of accumulation. Although the Hayashi group favor a nonimpact lunar origin, their theory is not necessarily inimical to an impact origin.

In any event, the central issue is the spectrum of planetesimal masses. Can we model the formation of the Earth by runaway growth of a single large embryo that sweeps up much smaller bodies, or are the impacting bodies not much smaller than the protoplanet? There is no unequivocal answer to this question, but current understanding tentatively favors the latter possibility. There are several steps and issues involved in reaching this assessment. First, one must understand the early evolution of a "gas" of small planetesimals, which may be initially monodisperse (i.e. about equal mass). This is best treated by kinetic theory, or "particles in a box," simulations. Although there are uncertainties in the collision physics, most controversy has centered around the correct treatment of gravitational stirring and scattering. Greenberg and coworkers have proposed "run-aways," in which certain embryos grow much faster than neighboring bodies because the gravitational cross section can be much larger than the physical cross section when the encounter velocities are small compared with the escape velocity from the embryo (Greenberg 1982). A series of calculations by Wetherill and coworkers (Wetherill & Cox 1984, 1985, Stewart & Wetherill 1986, Wetherill & Stewart 1986) indicates that although runaway is conceivable, a substantial embryo is needed to initiate the process. The presence or absence of gas is not important, and the uncertainties in modeling the collisions do not affect the conclusion. It is important to realize that even if runaway occurs, one is left with a very large number of bodies, probably of lunar size, and not with just four planets. Impacts between large bodies would still take place.

In the absence of runaway, it is usually possible to approximate the mass spectrum as a power law, $N(m) \propto m^{-\alpha}$, where $N(m) dm$ is defined as the number of planetesimals with masses between m and $m + dm$. In many models, as discussed by Wetherill (1980), α lies between 5/3 (for which each mass decade contributes an equal amount of surface area) and 2 (for which each mass decade contributes an equal amount of mass). The value of α may be affected by the presence of gas but still be in this range. Suppose we were standing on the protoearth, maintaining an inventory of incoming bodies as they impacted. Table 1 shows what the accumulated inventory might look like. Each line in the table refers to roughly a decade in mass (so that "10² Iapetus" means "roughly 100 bodies each within a

Table 1 Illustrative inventories of planetesimals required to form the Earth

$\alpha = 2^a$	$\alpha = 5/3^b$
1 Mars	5 Mars
10 Moons	3 Moons
10^2 Iapetus	12 Iapetus
.	.
.	.
.	.
10^9 1-km planetesimals	10^5 1-km planetesimals
Total = 1 Earth mass	Total = 1 Earth mass

^a Equal amount of mass in each logarithmic mass interval down to 10^{-10} Earth masses.
^b Equal amount of cross-sectional area contributed by each logarithmic mass interval (lower mass cutoff unimportant for total mass).

factor of three of the mass of Iapetus”). Strictly speaking, this table contradicts Safronov’s model (Safronov 1966, 1969), which has formed the basis for much of the work on planetary accumulation for more than a decade. Safronov argued that the second biggest body in a given accumulation zone is only $\sim 10^{-2}$ or 10^{-3} of the largest mass because of the enhanced (gravitational) cross section of the largest body. However, this is an overinterpretation of the Safronov model, which artificially isolates a preferred embryo planet in a specified zone of accretion and does not, therefore, make any pretense of explaining why there are 4 (rather than, say, 100) terrestrial planets.

There has never been a computer simulation that goes all the way from kilometer-sized planetesimals to fully grown planets. Ideally, one should use the output from kinetic (“particles in a box”) calculations, in which the Keplerian character of the orbits is not important, as the input for an orbital simulation of later stages. Currently, this is computationally prohibitive. However, numerical simulations support many of the most important features of Safronov’s (1969) model, especially the notion of a “steady-state” velocity. Safronov showed that the relative velocity of two planetesimals is of the order of the escape velocity from the largest planetesimal in the swarm. Guided by this theory and tests of its validity, Wetherill has performed simulations of planetary accumulation that come closest to satisfying the stated goals. Although other simulations exist (Lecar & Aarseth 1986) and analytical work continues (Horedt 1985, Hayashi et al 1965), the Wetherill simulations (Wetherill 1980, 1985) are closest to a realistic description of terrestrial planetary

formation because they are fully three dimensional and involve only minor approximations in the gravitational physics, except for neglecting the effect of possible resonances.

One interesting feature of Wetherill's results is that they indicate that many aspects of the final outcomes are insensitive to most of the initial conditions, so the precise dovetailing of early to later stages may not be needed. In particular, his results are *not* strongly affected by whether he begins with a monodisperse system or a mass spectrum with $\alpha = 1.83$. Wetherill's most recent simulations, largely motivated by the giant-impact hypothesis of lunar origin, yield terrestrial planetary "systems" with characteristics rather similar to those observed (Wetherill 1985). He typically begins with an initial swarm of 500 bodies distributed between 0.7 and 1.1 AU, with initial eccentricity randomly distributed between 0 and 0.05 and inclinations between 0 and 0.025 radians, as expected from Safronov's theory (Safronov 1969). As the Monte Carlo simulation proceeds, bodies grow and eccentricities and inclinations increase. Figure 2 shows the most important aspect of these simulations for the question of lunar origin—namely, that many large impacts occur on the body that eventually becomes Earth. These bodies may be a few times the mass of Mars in some cases. Typical impact velocities are $\sim 9 \text{ km s}^{-1}$ (slightly less than Earth escape velocity because of the finite size of the projectile and the fact that the Earth has not yet reached its final mass). These impacts

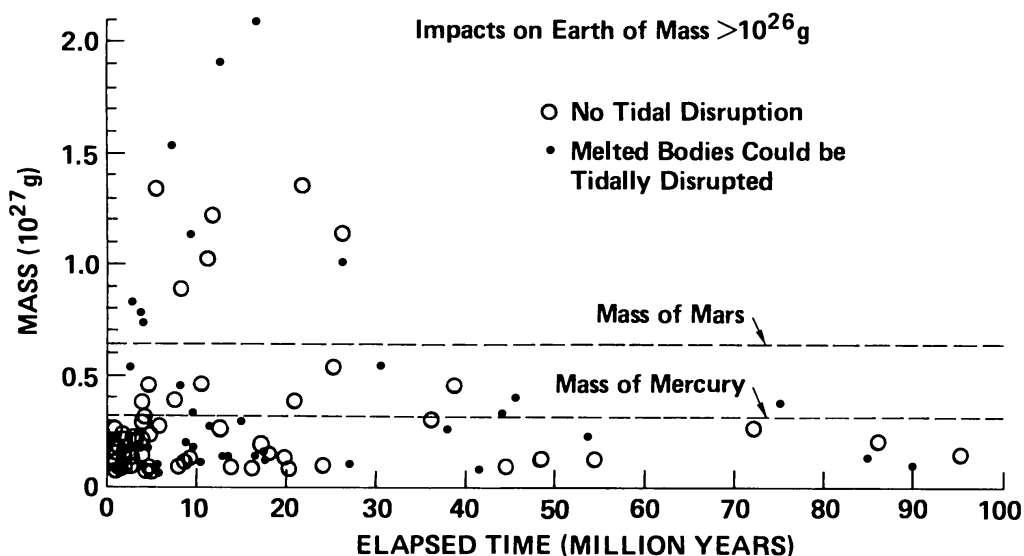


Figure 2 Combined results of 10 simulations of terrestrial planetary accumulation (Wetherill 1985), showing the time and size of giant impacts on Earth. In 5 of the simulations (labeled by open circles), tidal disruption of bodies previously impact melted is allowed for; in the other 5 simulations (closed circles), there is no tidal disruption.

are sufficiently common that they cannot be a special occurrence for Earth.

The final orientation of planetary spin axes (obliquities) and the orbital characteristics provide a measure of the role of large impactors. If all of the impacting mass is in the form of small bodies, then planetary orbital eccentricities, inclinations, and obliquities should all be small (leaving aside the subsequent, nonsecular perturbations among the planets). Wetherill's simulations provide reasonable eccentricities and inclinations because the large excursions excited by the largest impactors tend to be damped by the incoherent stream of smaller impactors. Analytical calculations (Harris & Ward 1982) suggest that the mass spectrum must be fairly "soft" (meaning much mass in smaller bodies, or $\alpha = 2$ in Table 1). Even in this case, Mars-sized impactors occur. The situation for obliquities is not so clear: Giant impacts would tend to randomize their values, yet the terrestrial planets tend to have a preferred (prograde) sense of spin. However, we are dealing with the statistics of small numbers, and two planets (Venus, Uranus) do have reversed spin. The pattern of obliquities generally supports the existence of large impacts, but the quantification of this is uncertain.

ORIGIN SCENARIOS

There are many reviews of lunar origin (most recently Boss 1986, Wood 1986). The main emphasis here is on the inadequacies of alternative ("conventional") origins. First, a philosophical point: It often seems that proponents of one or another theory adopt the attitude that their favorite theory is perfectly acceptable because it cannot be conclusively disproved. In this area of science, this is an invalid criterion. The origin of the Moon is a problem that involves many aspects and complications; it cannot be addressed by "sterilizing" it into an abstract dynamics problem in fission or capture or whatever. It is only by assessing the broad ramifications that a meaningful probability can be assigned. It is in this context that fission, intact capture, and possibly coformation (binary accretion) fail as satisfying explanations of lunar origin. I briefly discuss each of these alternatives in turn, and I also consider less clearly defined intermediate cases or modifications.

Fission

In the original Darwinian version (Darwin 1880), the protoearth rotates so rapidly that it is dynamically unstable to fission. There are two main problems with this: How do you create this state, and how do you explain the fact that the current angular momentum of the Earth-Moon system is lower than that needed for fission by a factor of at least three? There is

nothing in astrophysical or solar system experience to suggest that the requisite high angular momentum can be supplied gradually during the planetary accumulation process. On the contrary, the net angular momentum influx during accretion is a subtle and small effect (as mentioned at the end of the last section; see also Harris 1977). On the second issue concerning the total angular momentum budget, this problem could in principle be avoided by initiating fission at an earlier stage (well before the protoearth approaches its final mass) or by guaranteeing that a large amount of mass escapes to infinity, carrying away excess angular momentum (much as satellites outside the Saturnian ring system can “soak up” the outward angular momentum flux present in the rings by gravitational torques). These possibilities cannot be disproved, but neither do they arise naturally in any self-consistent, fully developed theory of planetary accumulation. More recent work on fission has focused on the dynamics and the possibility of disk formation (Durisen & Scott 1984), but the fundamental objections remain unanswered.

“Fission” has also been invoked in a highly modified form by Ringwood (1966, 1979), and the word has been used, perhaps inappropriately, by Stevenson (1984a) to describe the spin-out of a superrotating atmosphere immediately after a giant impact. Although Ringwood’s scenario has some dynamical difficulties, many aspects of his idea are remarkably similar to a possible aftermath of a giant impact. Since these “fission” proposals embody the physics of impact, they are more properly discussed later, when large impacts are described.

Capture

If collisions occur between the protoearth and bodies at least as large as the Moon, then close encounters are even more common. However, both Safronov’s theory and the numerical simulations indicate that the encounter velocity (i.e. the velocity at infinity) is significant (typically up to a few kilometers per second), so that a substantial amount of kinetic energy must be dissipated. This proves to be not possible by tidal dissipation unless the encounter distance is so close that the body is disrupted. (This case is discussed further below.) One could envisage a scenario involving several, successive nondisruptive encounters and damping of the excess energy, but such a model would be contrived, since the problem is not just one of three bodies (Sun, Earth, Moon) but involves other scattering bodies that will guarantee incoherence and hence a predominance of scattering (i.e. an *increase* of encounter velocity rather than decrease). Gas drag could also be invoked (Nakazawa et al 1983), but this model is also contrived because it requires a very small encounter velocity and because the gas responsible for capturing the Moon must subsequently be removed

rather quickly if the Moon is not to spiral inward and accrete onto the Earth. It is possible, however, that disruption does not occur if the viscosity of the planetesimal is too high (Mizuno & Boss 1985).

It seems almost superfluous to point out that no satisfactory explanation of lunar composition has arisen in models where the Moon is made elsewhere in the solar system. However, *disruptive* captures involving encounters within ~ 2 Earth radii and orbital injection of some debris might have happened, and this material may be preferentially from the mantle. Even though disruptive capture would only be important for small encounter velocities (Öpik 1972, Wood & Mitler 1974, Kaula & Beachey 1986, Hayashi et al 1985), it might contribute nonnegligibly to the Moon, provided disruption can occur at all (see Mizuno & Boss 1985).

Coformation (Binary Accretion)

Ruskol (1960), motivated by the ideas of Schmidt (see Ruskol 1982), pointed out that when two planetesimals collide within the Hill sphere of a planet, at least some of the debris may have both low enough energy and high enough angular momentum to end up in orbit. A circumplanetary disk can eventually form, fed by the debris of these collisions and, eventually, from collisions between the disk and later planetesimals. This disk can evolve, with one or more satellites forming at or beyond the Roche limit. This model has been developed further by Harris & Kaula (1975) and considered anew for lunar origin by Weidenschilling et al (1986).

The main virtue of this model is that it invokes a process that almost certainly happens, provided only that one accepts any of the hierarchical accumulation pictures of planetary formation. However, the model has three or four problems. First, it has difficulty explaining the iron depletion of the Moon. The possibility that the disk is a “compositional filter,” which selectively excludes iron because of its greater strength or density, has been suggested, but this scenario appears to require very restrictive conditions to work, if it ever works. Second, the disk is cold and particulate; the Moon grows at the same rate as the Earth and is never very hot. Consequently, one cannot easily explain the hot, primordial Moon (or the “magma ocean”). Third, the disk is fed by roughly equal amounts of positive and negative angular momentum; it seems difficult to ensure that the material has the requisite angular momentum to make the Moon. The fourth point, related to the third, concerns the evolution of the disk: The natural time scale for redistributing angular momentum within the disk is short compared with the formation time of Earth. Is it possible to even maintain the disk?

In the final analysis, the significance of coformation may rest mainly on the mass spectrum of the planetesimals. As Stevenson et al (1986) discuss,

it is only possible to put a significant amount of material into orbit by the Ruskol mechanism if most of the incoming mass is in small planetesimals. Given the other problems outlined above, it still seems likely that co-accretion is not *the* explanation of lunar origin, but at best a contributor of mass.

PHYSICS OF LARGE IMPACTS

The literature on impact physics is extensive but disappointing. It is extensive because impacts are an important planetary process and because there are obvious connections with the physics of explosions. It is disappointing because much of the work is empirical or “modular” (meaning that the person or persons responsible often do not understand what is going on in all parts of the calculation or interpretation because they are connecting together independently developed algorithms or procedures). Compilations of work in this area include Kinslow (1970), Roddy et al (1977), and the relevant subsection of Silver & Schultz (1982). The problem of a collisional lunar origin is even more challenging because it would have involved an impact far beyond any human experience. This means that any attempt to “scale” known impacts is probably doomed. A better understanding of the fundamentals is needed.

There are three sets of issues confronting this better understanding :

1. *Thermodynamics (equation of state)*. How does the material respond to high shock pressures? In particular, what is the irreversible entropy production and the extent of melting and vaporization following shock release?
2. *Constitutive law (rheology)*. What is the relationship between stress and strain during postimpact expansion? What viscosity (turbulent or otherwise) or small-scale instabilities characterize the macroscopic flow?
3. *Dynamics (equation of motion)*. How does the shocked material flow, subject to nonuniform gravity, pressure gradients, and “viscous” stresses?

I discuss each of these issues in turn, but it is first valuable to motivate the important questions. If we wish to explain lunar origin by one or more giant impacts, then some material must be emplaced in an Earth orbit. In impact events, one normally thinks of the positive acceleration of debris as confined to the immediate vicinity of the impact site. Upward-moving debris is subsequently subject only to the r^{-1} gravitational potential of the (approximately) spherical Earth and suffers one of two fates. If it has a total (gravitational plus kinetic) energy that is positive, then it escapes on a hyperbolic trajectory. If it has a total energy less than zero, then it

traverses a closed elliptical orbit that must eventually reimpact the Earth. This is illustrated in Figure 3*a*. Actually, reimpact of all negative energy debris also occurs if one allows for the higher order field of an oblate (rapidly rotating) Earth. The important point is that one needs a “second burn”—that is, some way of raising the periapse of negative energy material above the surface of the Earth. Each of the three issues listed above is connected to possible ways of achieving this “second burn.”

Figure 3*b* shows one way to do this. Following a large impact in which significant vaporization occurs, the outflowing material is subject to pressure gradients as well as to gravity. This hydrodynamic effect increases the kinetic plus gravitational energy of the material but, more importantly, can also increase its angular momentum and lift the periapse above the surface. In order for this to happen, a significant pressure gradient must act over a substantial fraction of a planetary radius. (This is quantified more fully below.)

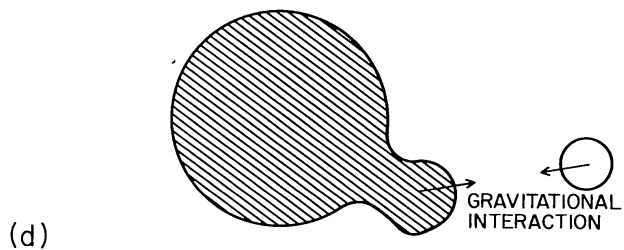
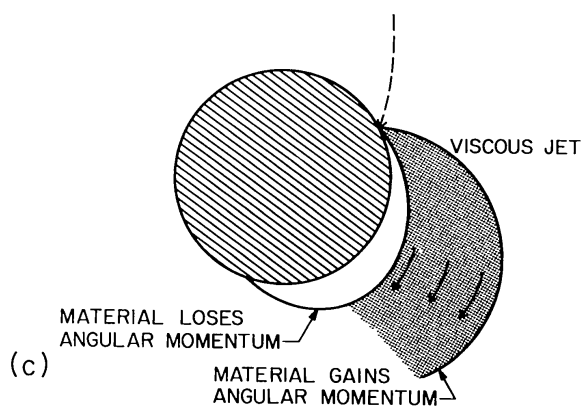
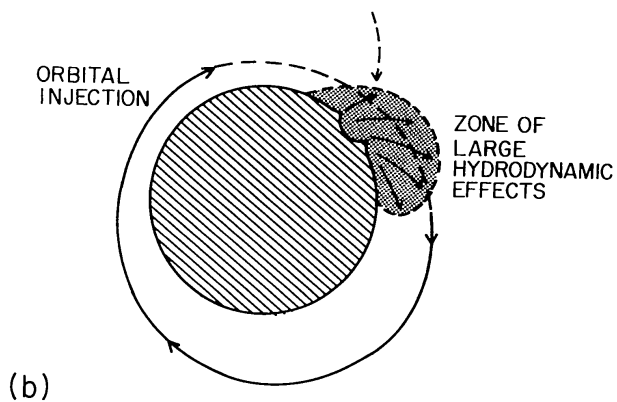
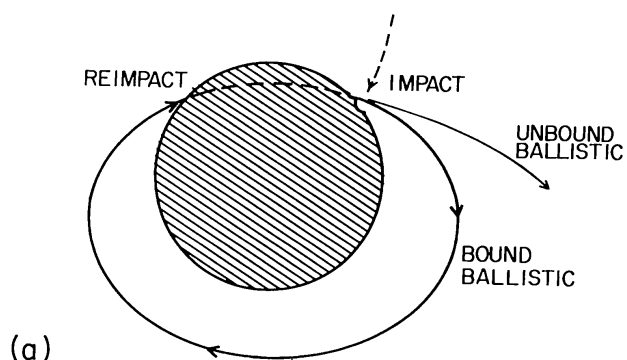
Figure 3*c* shows another way of achieving orbital injection. “Viscous” stresses within the jet of outgoing material redistribute angular momentum, allowing some material to achieve orbit at the expense of other material that loses angular momentum and falls back to Earth. This is a less plausible scenario than pressure-gradient acceleration because it requires a very large effective viscosity (sufficient to redistribute a large amount of angular momentum in just one orbital period).

Figure 3*d* shows a third way of achieving orbital injection, which involves gravitational torques (or, more generally, a severe time-dependent perturbation of the gravitational potential). Either a bulge on the planet or a separate body could transfer angular momentum to a more distant body, leading to orbital injection. This occurs in some recent numerical simulations (Benz et al 1986*b*) and is discussed further below and in the next section.

Thermodynamics of Impact

Contrary to what one might suppose from a cursory inspection of the literature, the thermodynamics of impacts are not well understood,

Figure 3 Schematic outcomes for outflow from a giant impact. (*a*) The material follows a purely Keplerian trajectory and must either escape or reimpact. (*b*) Pressure gradients near the impact site increase the angular momentum of (negative energy) material, which allows it to achieve orbital injection. (*c*) Redistribution of angular momentum within a viscous jet allows some material to gain angular momentum and achieve orbital injection at the expense of other material that reimpacts. (*d*) The gravitational torque exerted by the bulge or discrete body on the more distant protomoon allows the latter to gain angular momentum and avoid reimpact.



especially when substantial vaporization takes place. The general principles are known, but the detailed quantification is incomplete. During impact, material is very rapidly shocked to high pressure and temperature; irreversible entropy production occurs at this stage. This material then expands approximately isentropically from the peak pressure state. It is, of course, not a thermodynamic question as to whether the expansion phase is well approximated as isentropic; this depends on opacity, turbulence, etc. For the moment, isentropy is assumed. Often the peak pressure state is supercritical (meaning that it cannot be characterized as being either gas or liquid); during expansion it eventually hits a phase boundary, and the extent of vaporization is then well defined. It is popular in the literature to state the degree of vaporization at a nominal 1-bar pressure. This is arbitrary and somewhat misleading, since (as we shall see) there is nothing special about a pressure of 1 bar in terms of understanding giant impacts on the protoearth. More importantly, the degree of vaporization can be substantially different if one chooses (say) a kilobar or a millibar as the nominal state. In the brief historical survey below, the nominal state is 1 bar.

Stanyukovich (1950) was first to estimate vapor production and obtained $\sim (V_{\text{imp}}/14 \text{ km s}^{-1})^2$ in units of projectile mass, assuming cold (terrestrial rock) starting material. Notice that the kinetic-energy content of the projectile has to be over an order of magnitude larger than the heat of vaporization of rock ($\sim 10^{11} \text{ erg g}^{-1}$). There are two reasons for this low yield: (a) Much of the energy remains kinetic (at least initially), and (b) much of the energy is stored in the internal energy of compression. Subsequent calculations, based on more substantial data, predict even less vaporization (Ahrens & O'Keefe 1972, O'Keefe & Ahrens 1977, 1982), except when the target is a magma ocean (Rigden & Ahrens 1981). All of these calculations are potentially misleading because they assess vaporization in terms of the shock process alone; no contribution due to gradual degradation of kinetic energy is computed. In a giant impact, all of the impact energy must eventually be accounted for in assessing vaporization. Whereas in small impacts it is valid to neglect energy release from the rain-out of debris far from the impact site, there is no such place as "far from the impact site" on Earth if the projectile is the size of Mars!

It is instructive to perform a computation of prompt, irreversible entropy production in a shock event to show where the uncertainties arise. For this purpose I consider SiO_2 , where the data base is far more complete than for more appropriate starting materials (e.g. Mg_2SiO_4 or MgSiO_3). Approximate calculations for more realistic silicate assemblages indicate that this is not a serious deficiency [mainly because the vapor pressure of

MgO is not much different than that for SiO₂ (Krieger 1967)]. The entropy production upon shock compression is given by

$$\Delta S = \int_{T_s}^{T_h} C_v dT/T, \quad (1)$$

where T_h is the Hugoniot temperature (the temperature behind the shock wave), T_s is the temperature that the material would have at the same density if the compression had been isentropic, and C_v is the specific heat at constant volume (usually not constant). To get a rough idea as to how ΔS scales, we can use the empirical fact that at very high (megabar) pressures, we have

$$T_h \simeq T_1 P, \quad (2)$$

where P is the peak pressure and T_1 is some constant. This behavior is very frequently observed [e.g. Lyzenga et al (1983) for SiO₂]. We also have

$$T_s \simeq T_0(\rho/\rho_0)^\gamma, \quad (3)$$

where T_0 is the initial temperature, ρ is the shock pressure, ρ_0 is the initial density, and γ is the Gruneisen parameter (here assumed constant). Moreover, we have

$$P = \rho_0 U_s U_p \propto \rho_0 V_{\text{imp}}^2, \quad (4)$$

where U_s is the shock velocity and U_p is the particle velocity [again assuming very high pressure ($P \gg$ bulk modulus)]. If we let $\beta \equiv d \ln P / d \ln \rho$, then we have

$$\Delta S \simeq (1 - \gamma/\beta) C_v \ln (V_{\text{imp}}^2) + \text{constant}. \quad (5)$$

To estimate vapor production we then compare this expression with ΔS_v , the entropy difference between liquid and vapor. It is interesting, but perhaps counterintuitive, that ΔS is only a weak (logarithmic) function of impact velocity, since we would have expected the total vapor production to scale as the kinetic energy of the projectile (e.g. O'Keefe & Ahrens 1982). The resolution of this apparent conflict lies in the realization that at very high impact velocities, a volume much larger than the projectile volume is at least partially vaporized, and it is the integral over all this volume (which scales roughly linearly with peak shock pressure) that is relevant. For our *present* purposes, it is likely to be the local vapor content that is important, since we wish to consider the role of pressure-gradient acceleration.

A detailed calculation was carried out using the following data and theory: The equation of state and shock temperatures are from Lyzenga & Ahrens (1980) and Lyzenga et al (1983). These data were also used to

estimate γ and C_v . An extended Debye model was used for C_v (Kieffer 1979). The Dulong-Petit limit of $C_v = 6 \text{ cal (mole-atom-K)}^{-1}$ appears to be exceeded at $T_h \gtrsim 6000 \text{ K}$ ($P \gtrsim 1 \text{ Mbar}$), probably because of electronic excitation. The Gruneisen parameter is also mildly temperature dependent. It was also assumed that the initial condition was a 50% molten Earth at $T \sim 1800 \text{ K}$. The motivation for this choice is work done on planetary accretion (Kaula 1979, 1980, Stevenson 1981, 1983a, Davies 1985) suggesting that the heat retention due to prior impacts is sufficiently large to guarantee a hot target. The actual surface of the protoearth may be cold (because of radiative cooling), but the thermal boundary-layer thickness would be \sim few kilometers, negligible compared with the excavation depth of the projectile. The protoearth may even have had a magma ocean, as assumed by Rigden & Ahrens (1981) and Hofmeister (1983), but the entropy difference between solid and liquid is small compared with that needed for vaporization, so the prior presence of a magma ocean is not important. It is important that the target be hot; the same calculation with an initial temperature of 300 K yields much less vaporization. The hot initial state means that the shock data on highly porous silica (aerogel) of Holmes et al (1984) were especially useful.

The results from numerical integrations of Equation (1) are shown in Figure 4, superimposed on the SiO_2 phase diagram. This entropy-pressure representation, popularized in volcanology by Kieffer (1982), seems particularly suitable for understanding postimpact expansion. Each vertical line represents an isentropic pressure release path from the peak pressure achieved for a given impact velocity under the assumptions of normal impact and equal material properties for target and projectile. The phase boundary was computed using JANAF thermochemistry tables (*JANAF Thermochemical Tables* 1971) and allowing for the following species in the gas phase: SiO_2 , SiO , Si , O_2 , and O . In fact, SiO and O_2 dominate. Pressure corrections for the vapor-phase chemical potentials were made assuming ideal gas partition functions with fully excited internal degrees of freedom. The resulting vapor pressure agrees with the calculation of Krieger (1967) but *not* with some data (Ruff & Schmidt 1921). The source of the discrepancy is not known. The extrapolation to the critical point is highly uncertain, but it can be crudely estimated using corresponding-states arguments (see also Ahrens & O'Keefe 1972).

These results are revealing in several ways. First, the fact that entropies are additive means that a lever rule can be applied to determine the extent of vaporization at a given pressure on the release path. For example, material subject to a 10 km s^{-1} impact is about 20% vapor, 80% liquid *by mass* at $P = 1 \text{ kbar}$ ($T \sim 4000 \text{ K}$). As expansion proceeds, the liquid boils, producing more vapor. (This is directly analogous to the boiling of

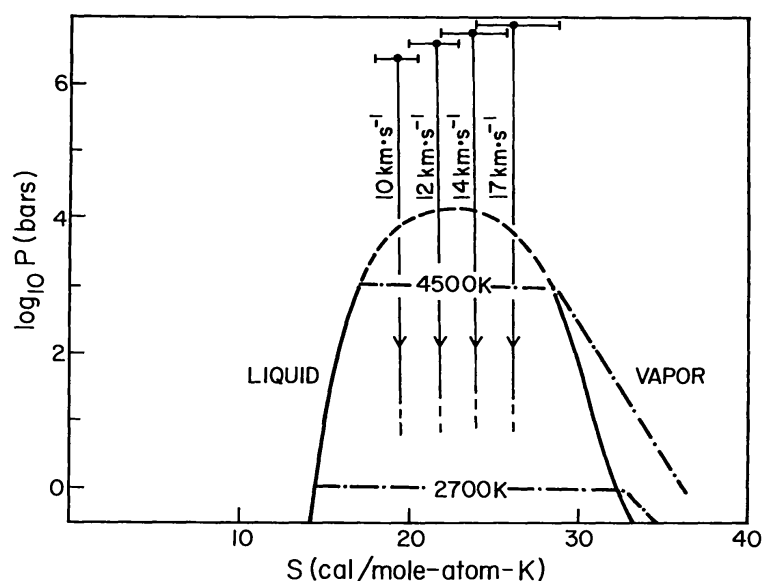


Figure 4 Phase diagram for SiO₂, with superimposed isentropic pressure-release lines for postimpact flows corresponding to the impact velocities indicated. The critical region (where liquid and vapor become indistinguishable) is not well known and is indicated by a dashed line. The error bars at the top of each trajectory indicate the uncertainties in entropy computation. By applying the lever rule, this diagram indicates the degree of vaporization for a given irreversible entropy production.

water during isentropic expansion to lower pressure and temperature.) If the postimpact entropy is larger than the critical value of ~ 22 cal (mole-atom-K)⁻¹, then the expansion involves condensation of silica droplets (i.e. the mass fraction of liquid increases as expansion proceeds). This applies for impact velocities $\gtrsim 14$ km s⁻¹. In any event, a substantial fraction of the mass (and almost all the volume) is eventually in vapor form, even for an impact at 10 km s⁻¹, a value that coincidentally is comparable to the escape velocity from Earth or the impact velocity on Earth (cf. Proposition 4). A second interesting feature of Figure 4 is indicated by the error bars associated with each vertical curve. These attempt to incorporate all the shock data and theoretical extrapolation uncertainties entering computation of ΔS . The uncertainties become very large at high impact velocity because large extrapolations of shock data are required and electronic corrections to γ and C_v become important, yet uncertain.

Of course, this calculation does not deal with the actual conditions of an oblique impact because of its highly idealized nature of a plane, normal shock, but it does provide a useful guide to the possible extent of prompt vaporization. The differences between oblique and normal impacts are discussed in the next section.

There is another way of analyzing vaporization, on purely energetic

grounds, that relies on the expectation that the impact is so large that the postimpact Earth can be treated as spherically symmetric, with the surface layer heated more than the interior. [Clearly, this would be a ridiculous assumption for small impacts, even for the body that supposedly caused biological trauma at the end of the Cretaceous (see Silver & Schultz 1982).] Calculations were carried out assuming injection of a specified total amount of energy ΔE (some fraction of the kinetic energy of the projectile), which is then partitioned among internal energy (mostly thermal) and increased gravitational energy (because the planet puffs up). As a rough guide for what to expect, a Mars-sized projectile impacting at 10 km s^{-1} has an energy of $3 \times 10^{38} \text{ erg}$, which, if completely converted to heat, would increase the average temperature of an Earth mass by roughly 5000 K (if we assume no latent heat buffering). Figure 5 shows estimated temperature profiles for an arbitrary heat injection of $1.5 \times 10^{38} \text{ erg}$ and a range of energy emplacements varying from uniform to that obtained by assuming a pressure drop-off away from the impact site as r^{-2} [a likely upper bound to the rapidity of the drop-off (e.g. Orphal et al 1980)]. These calculations show that the postimpact Earth certainly has a deep magma ocean with a supercritical near-surface layer that merges continuously with a vapor atmosphere. The cooling time of this system is defined by

$$\tau_{\text{cool}} = \frac{\Delta E}{\sigma T_e^4 4\pi R^2}, \quad (6)$$

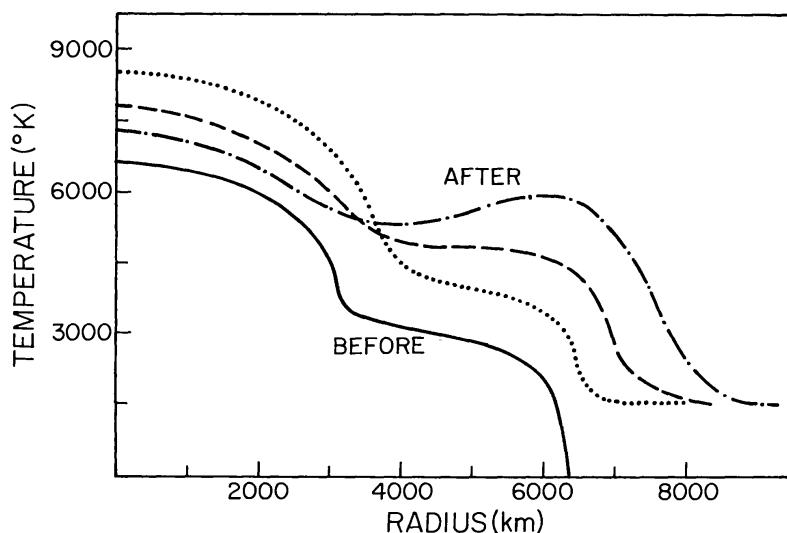


Figure 5 Several possible outcomes arising from injecting $1.5 \times 10^{38} \text{ erg}$ into Earth because of a giant impact. The preimpact state (solid line) is assumed to be 50% molten because of previous smaller impacts. Three possible postimpact states vary from uniform energy deposition (dotted line) to a concentration in the mantle (dash-dot line). In each case, the planet “puffs up” and an extended transient atmosphere ($T \sim 2000 \text{ K}$) develops.

where T_e is the effective temperature, σ is the Stefan-Boltzmann constant, and R is the radius of the photosphere. The radiating level is probably determined by the opacity of dust grains that condense (Thompson & Stevenson 1986) and corresponds to $T_e \sim 2000$ K. For this choice, we obtain $\tau_{\text{cool}} \sim 10^3\text{--}10^4$ yr. Notice that this is an enormous time compared with dynamical times. This is a fundamental distinction between giant impacts and small impacts (such as the event at the end of the Cretaceous). Perhaps the most interesting point of Figure 5 is that it demonstrates how traumatic this giant impact would have been for Earth. Perhaps the best evidence for the event will be found in the deep mantle in the form of residual evidence for magma ocean differentiation and layering (e.g. Ohtani 1985).

Clearly, a description of the postimpact behavior requires consideration of two-phase media (a liquid with bubbles or a gas with droplets). One property of a two-phase medium is that it is not completely specified by two thermodynamic variables in the usual way. For example, a medium with specified P and T (these two variables being linked by the Clausius-Clapeyron equation for vapor equilibrium) can have any liquid mass fraction between zero and unity. This means it can have an enormous range of average densities and average specific entropies for the same P and T . It is also possible to have an enormous range of pressures for a given internal energy. Unfortunately, this thermodynamic requirement is not satisfied by all equations currently in use. In particular, it is not satisfied by the Tillotson equation of state (Tillotson 1962, Allen 1967), despite its plethora of adjustable constants. This is the equation used initially by Benz et al (1986a,b), and thus it casts doubt on the accuracy of their vaporization estimates and pressure-gradient effects. More realistic equations of state are now being used. Regrettably, it is not possible at the time of this review to assess whether this greatly modifies the results.

Constitutive Properties

The focus here is on “viscosity,” due both to microscopic processes and to fluid dynamical processes. This is less well understood than the thermodynamics but possibly also less important. In fact, numerical simulations of impacts do not usually include viscosity except as a numerical artifice to promote stability of the code. We return to this “artificial” viscosity in the next section; here, we pose the question, What is a dynamically interesting viscosity and might it exist? A dynamically interesting viscosity would be one for which an element of material ejected from the impact site is subjected to a viscous couple in one orbital period large enough to significantly increase its angular momentum and to make orbital injection possible. Of course, this would be done at the expense of other

fluid elements that lose angular momentum. Clearly, a complete quantitative answer requires a detailed model of jets, but we can get a rough idea by assuming Keplerian differential rotation of a disk of material and by borrowing from the physics of accretion disks (Lynden-Bell & Pringle 1974). In steady state, we have

$$F \frac{dh}{dR} = \frac{dg}{dR}, \quad (7)$$

where F is the total mass flux radially outward at radius R of a disk of material with local specific angular momentum $h = R^2\Omega$ (Ω is the angular velocity). The viscous couple is g , and we assume that the radial velocity is $\sim 0.1R\Omega$ (to raise the periaipse of outer material in the disk in one orbital period). Then, since

$$F = 0.1\pi R^3\sigma\Omega, \\ g = 2\pi R^3\nu\sigma\Omega, \quad (8)$$

where σ is the surface density of the disk (mass per unit area), Equation (7) implies that we need a viscosity

$$\nu \gtrsim 0.01R^2\Omega, \quad (9)$$

or $\nu \gtrsim 10^{14}\text{--}10^{15} \text{ cm}^2 \text{ s}^{-1}$ typically. (This is about the viscosity of glacier ice on Earth.) The microscopic viscosity of a liquid ($\sim 10^{-2} \text{ cm}^2 \text{ s}^{-1}$) or of a gas ($\sim 10^2 \text{ cm}^2 \text{ s}^{-1}$ typically) are small by comparison. The *bulk* viscosity of a foam (bubbly liquid) can be remarkably high [$\sim 10^{11} \text{ cm}^2 \text{ s}^{-1}$; Stevenson 1983b) because of the irreversible entropy production accompanying the induced phase change between gas and liquid, but even this value is not large enough to be important. However, fluid dynamical instabilities could produce viscosities approaching that required. In the Prandtl picture of turbulent viscosity, we can imagine blobs of fluid with relative velocity u and size l ; the resulting viscosity is then $\sim ul$ and could be $\sim 10^{14} \text{ cm}^2 \text{ s}^{-1}$ for $u \sim 10 \text{ km s}^{-1}$ and $l \sim 1000 \text{ km}$. More specifically, Thompson & Stevenson (1983, 1986) point out that a two-phase medium is susceptible to gravitational patch instabilities, even close to the Earth, because the two-phase medium is highly compressible. These instabilities promote turbulence because they cannot evolve all the way to formation of a self-gravitating sphere if they occur within the Roche limit. However, the estimated turbulent viscosity is then $\lesssim l_{\text{crit}}^2\Omega$, where l_{crit} is the wavelength of the instability, probably only about 10^2 km . The resulting turbulent viscosity is then $\lesssim 10^{12} \text{ cm}^2 \text{ s}^{-1}$, less than the constraint give by Equation (9). A viscosity of $10^{12} \text{ cm}^2 \text{ s}^{-1}$ (or even 1000 times less) is still extremely

interesting for lunar formation, however, if the Moon forms from a disk. This is discussed later.

Dynamics

Even with a complete understanding of thermodynamics and rheology, there is a remarkable range of possible fluid dynamical outcomes. These can only be fully understood by numerical simulation. In this section the possibilities are described in very simple terms.

Suppose that non-Keplerian effects (primarily pressure-gradient acceleration) act from the impact site out to a height h above the Earth's surface (radius R_{\oplus}), but that the subsequent motion of an element of material is Keplerian. The initial vertical and horizontal components of velocity at height h are taken to be V_r and V_h , respectively. In the limit $h \ll R_{\oplus}$, the best candidate trajectories for orbital injection are those for which $V_h \gg V_r$, but with negative total energy. We can express the velocity components as

$$\begin{aligned} V_r &= A \left(\frac{GM_{\oplus}}{R_{\oplus}} \right)^{1/2}, \\ V_h &= B \left(\frac{GM_{\oplus}}{R_{\oplus}} \right)^{1/2}. \end{aligned} \quad (10)$$

The requirement that a bound orbit result is that $A^2 + B^2 < 2/(1+x)$, where $x \equiv h/R_{\oplus}$. The requirement that periapse lie above the Earth's surface is

$$B^2 > \frac{A^2 + 2x}{x^2 + 2x}. \quad (11)$$

For each value of total velocity $V_t = (V_r^2 + V_h^2)^{1/2}$, one can define a cone of trajectories for which both criteria are satisfied. The solid angle of this cone, divided by 2π , can be thought of as the "probability" P of orbital injection for each V_t and x . This is shown in Figure 6. Even for quite high starting elevations, represented by x , the value of P is low, since it is truncated at high velocities by escape. In a crude way this calculation gives the requirement that must be satisfied to inject material: Nonballistic processes (the "second burn") must be able to achieve a value of x and V_t so that injection is possible. Notice that other factors being equal, injection into a highly elliptical initial orbit (i.e. total energy only slightly negative) is favored because such an orbit maximizes P .

We turn next to the question of whether pressure-gradient acceleration is capable of acting out to a significant fraction of the Earth's radius, so that the initial conditions in the above calculation could be achieved. For

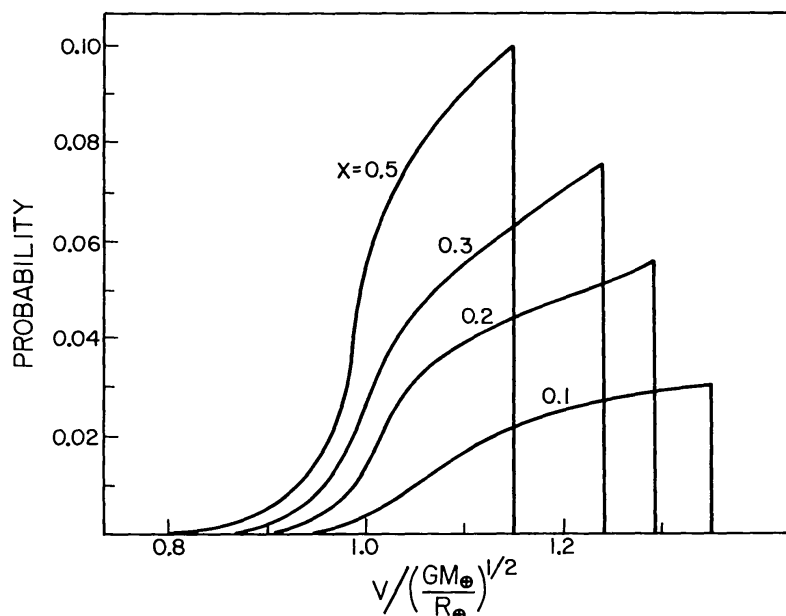


Figure 6 Probability of orbital injection, here defined as the fraction of all possible Keplerian trajectories (equally distributed in solid angle), starting from a nondimensional height $x \equiv h/R_{\oplus}$ with velocity V , for which orbital injection occurs. The truncation at high V occurs because more energetic trajectories lead to escape. The value of x can be thought of as a measure of the regions within which hydrodynamic effects are very important. Hence, large (small) x applies to large (small) impactors.

simplicity, consider steady-state gas flow expanding hemispherically away from a point source. We estimate the total kinetic energy gain between some radius r_0 from the source and infinity from the Euler equations:

$$\begin{aligned}
 \int_{r_0}^{\infty} u \frac{\partial u}{\partial r} dr &= \frac{1}{2} (u^2|_{\infty} - u^2|_{r_0}) \\
 &= \int_{r_0}^{\infty} -\frac{1}{\rho} \frac{\partial p}{\partial r} dr \\
 &= \frac{\gamma p}{(\gamma - 1)\rho} \Big|_{r_0}, \tag{12}
 \end{aligned}$$

where u is the radial velocity and we assume $p \propto \rho^{\gamma}$ as the adiabatic equation of state. We wish to define the hemisphere of radius r_0 as the region outside of which only small hydrodynamic effects occur. Somewhat arbitrarily, let us require that the change in $u^2/2$ is only 10% of the gravitational energy during the expansion from r_0 to ∞ . (The flow need not actually go to infinity; the results are similar if the flow is bound.) We then must define r_0 as the place where $p \approx 10^3$ bar, $\rho \approx 0.1$ g cm $^{-3}$ ($T \approx 5000$ K, $\gamma \approx 1.2$), where we assume most of the mass is in the form

of vapor. This corresponds to a release path at $V_{\text{imp}} \sim 12 \text{ km s}^{-1}$ (Figure 4). For a projectile mass M_{proj} , it follows that

$$\frac{r_o}{R_{\oplus}} \sim 3(M_{\text{proj}}/M_{\oplus})^{1/3}(V_{\text{imp}}/14 \text{ km s}^{-1})^{2/3}. \quad (13)$$

This formula can be loosely translated as the predicted value of x for Figure 6. It suggests that all impacting bodies with mass $\gtrsim 10^{-3} M_{\oplus}$ are capable *in principle* of injecting a significant fraction of their own mass (a mixture of target and projectile) into orbit. Most importantly, Equation (13) suggests a rather weak dependence of injection efficiency on projectile mass. Of course, the actual efficiency can only be assessed by numerical simulation.

Melosh and Sonett (1986) emphasize a particular aspect of impacts that may be very important for orbital injection. They point out that the highest speed ejecta thrown out in the earliest stage of impact cratering is a *jet* of very fast moving material, usually less than 10% of the projectile mass. They find that the fraction of jetted material is greater during oblique impacts and favor this material for orbital injection because it is the most likely material to be vaporized. The existence of a jet is not in dispute (Gault et al 1968, Kieffer 1975, 1977a), but its role in the impact theory of lunar origin is unclear, since (as argued earlier) substantial vaporization is achievable even at 10 km s^{-1} impact velocity, provided that both the projectile and target are hot. Therefore, it is not clear whether the jet is essential or even desirable—since this material is also the most likely to escape.

The final issue we address here concerns large deviations from an r^{-1} gravitational potential. In appropriate circumstances, this could raise the periape of ejecta above the Earth's surface and allow orbital injection. Consider, for example, the situation in which a small body is ejected into orbit about two larger bodies (the protoearth and the remainder of the projectile) that are undergoing merger. There are two effects here: The center of mass is moving toward the center of the protoearth, and there are higher-order terms in the gravitational potential. The former means that even a *closed* orbital trajectory need not reimpact the Earth. In fact, the periape can be “raised” by an order of the distance that the center of mass moves relative to the protoearth (leaving out the small change in Earth radius due to merging). This could easily be $\sim 0.1 R_{\oplus}$. The deviation from r^{-1} in the gravitational potential is more complicated, but it can be thought of as a gravitational torque, analogous to that responsible for the current recession of the lunar orbit, only much larger. Consider a mass anomaly ΔM exerting a torque on an orbiting mass m at distance R (Figure

3d). This torque is of order $Gm\Delta M/R$ times some numerical factors involving the geometry. In a time τ , the angular momentum transferred to m is of order $Gm\Delta M\tau/R$. This value should be compared with $m(GM_{\oplus}R)^{1/2}$. To achieve a 10% change in angular momentum would require

$$\frac{\Delta M}{M_{\oplus}} \sim \frac{0.1}{\Omega\tau}, \quad (14)$$

where Ω is the angular velocity of the orbit of m . Plausibly, we could have $\Omega\tau \sim 0.3$ – 0.6 , and a mass of order two Mars masses would be sufficient. Notice that the result is independent of m , provided $m \ll M_{\oplus}$. Of course, this *only* works if ΔM is transient, since otherwise the torque varies in sign as this bulge rotates beneath the orbiting protomoon. It is desirable that the bulge relax on a dynamical (free-fall) time scale.

NUMERICAL SIMULATIONS

At the Origin of the Moon Conference in 1984, Cameron presented pioneering simulations of an impact origin of the Moon (Cameron 1985b). These early calculations have already been completely superceded, but even so, simulation work is only in its infancy. I concentrate here on the efforts of Benz et al (1986a,b), which are the only reasonably well-documented computations at the time of this writing. Calculations in progress by Kipp & Melosh (1986) are difficult to assess at this stage because these authors' earliest results omitted self-gravity of the ejecta.

Benz et al studied three-dimensional numerical simulations of oblique impacts of Mars-sized bodies on Earth. They chose to explain the present angular momentum of the Earth-Moon system ($3.5 \times 10^{41} \text{ g cm}^2 \text{ s}^{-1}$) entirely by a single impact; this constrains the relationship between projectile mass M_{proj} , impact parameter d , and impact velocity V_{imp} :

$$0.085 \simeq \left(\frac{M_{\text{proj}}}{M_{\oplus}} \right) \left(\frac{d}{R_{\oplus}} \right) \left(\frac{V_{\text{imp}}}{11 \text{ km s}^{-1}} \right). \quad (15)$$

This equation makes the preferred choice of a Mars-sized body ($M_{\text{proj}} \approx 0.1 M_{\oplus}$) self-evident, although one could go to somewhat more massive bodies striking more nearly head-on. The equation also illustrates why at most a small number of projectiles can provide the required angular momentum, since their contributions tend to add incoherently and make the angular momentum constraint harder to satisfy. The constraint assumes that the loss of angular momentum to infinity, carried by fast-moving ejecta, is small. The simulations tend to be consistent with this assumption. The

system is assumed energetically closed in the sense that radiation loss is small on the relevant dynamical time scale of hours. This is already suggested by Equation (6) and is more explicitly demonstrated in the next section. Benz et al use a Tillotson equation of state that agrees with shock-wave data on granite in the high-density limit (Allen 1967) and with an ideal gas in the low-density limit. As briefly discussed earlier, this equation of state does *not* describe a two-phase medium. Moreover, the quoted degree of vaporization in Benz et al (1986a) is incorrect, since it is based only on the sum of the mass elements that are *completely* vaporized. In fact, most of the mass is partially vaporized. The seriousness of these deficiencies is not known at the time of this writing.

Benz et al used a method known as smoothed particle hydrodynamics (SPH) in which the medium is represented by a finite number of mass points whose trajectories are followed in a Lagrangian sense. In this way, the Navier-Stokes equations are translated into an N -body problem, a desirable feature in a problem that has a complicated three-dimensional geometry. The SPH method is a recent development in astrophysical computational fluid dynamics (Lucy 1977, Gingold & Monaghan 1979; see also other references in Benz et al 1986a). Each particle is assumed to have its mass spread out in space according to a given distribution called the kernel. Details of the method can be found in the references given above. One other feature of interest is that an “artificial viscosity” is introduced in order to avoid unacceptably large postshock oscillations. However, this viscosity corresponds to a Reynolds number of $\gtrsim 200$ in the postshock flow (W. Benz, private communication), which is better than what standard finite-difference techniques achieve and large enough to suggest that (unphysical) angular momentum redistribution in the ejecta flow is tolerably low.

Two sets of calculations have been carried out. In Benz et al (1986a), there was no iron core present. For a Mars-sized impactor with grazing incidence, they found that the impactor is not completely destroyed. Instead, a clump of the projectile most distant from the Earth’s center is injected into a highly elliptical orbit. However, this orbit brings the material back to within the Earth’s Roche limit, so they conjecture that the material is sheared out by tidal forces and forms a disk. In the more recent calculations (Benz et al 1986b, W. Benz, private communication), a core is included in the projectile and more massive projectiles are considered. The results were somewhat different and perhaps surprising. Most analysis has been for low-velocity impacts (meaning that the velocity at infinity is small, so $V_{\text{imp}} \sim 11 \text{ km s}^{-1}$). Assuming always that Equation (15) is satisfied, Benz et al find that if $M_{\text{proj}} < 0.12 M_{\oplus}$, then too much iron ends up in orbit. If $M_{\text{proj}} \gtrsim 0.16 M_{\oplus}$, then the impact is closer to head-on and too

little mass ends up in orbit to explain the Moon. If $0.12 M_{\oplus} \lesssim M_{\text{proj}} \lesssim 0.16 M_{\oplus}$, then *three* substantial bodies are present (see Figure 7 for the simulation depicting $M_{\text{proj}} = M_{\oplus}/7$). These bodies are Earth, the core of the projectile, and a small (\sim Moon mass) body in a more distant orbit, which may pick up enough angular momentum by gravitational torque (Equation 14) to remain intact. Alternatively, the latter body may break up tidally to form a disk. Meanwhile, the intermediate-mass (iron-rich) body merges with Earth. If, instead, the impact velocity is higher, then some mass leaves the system completely, but it seems likely that either there is insufficient mass placed in orbit or too much iron in the material placed in orbit.

It would be premature to reach firm conclusions on the basis of these results, for at least three reasons. First, it is not clear whether the thermodynamic code in use is adequate. Second, the final outcome is not determined. (The simulations only cover the first ~ 20 hr.) Third, there have been insufficient simulations to understand the entire range of possibilities. In particular, the requirement that all the Earth-Moon angular momentum is explained by a single event is still a questionable simplification or application of Occam's razor. (All applications of Occam's razor in planetary science should be viewed with skepticism.) Nevertheless, the results yield some interesting insights. The following four suggested implications are most striking:

1. Pressure gradients may not always play an essential role in putting material into orbit.
2. "Clumpiness" of ejecta may happen, rather than well-dispersed clouds.
3. The Moon-forming material may come primarily from the projectile.
4. There are often difficulties in *preventing* the incorporation of metallic iron in Moon-forming material.

I think it is wise to be skeptical about all four "conclusions" at our current premature state.

DISK EVOLUTION

One possible outcome of a giant impact is a gravitationally bound clump of material, of order one Moon mass, that retains its integrity and evolves outward by tidal friction. A more likely outcome is the formation of a disk, either directly from a broadly disseminated "cloud" of two-phase (liquid-gas) ejecta or indirectly by the tidal disruption of a clump of material placed in an orbit so elliptical that it undergoes grazing encounters with Earth. Here, we focus on the evolution of this disk and how the Moon might form from it.

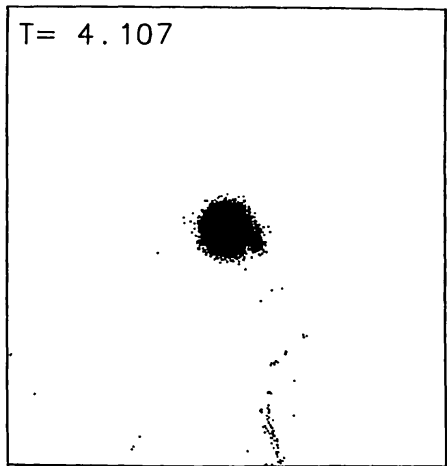
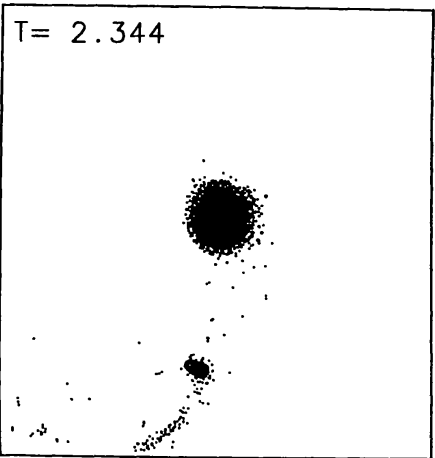
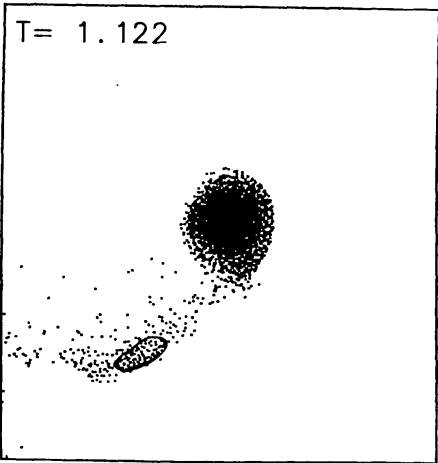
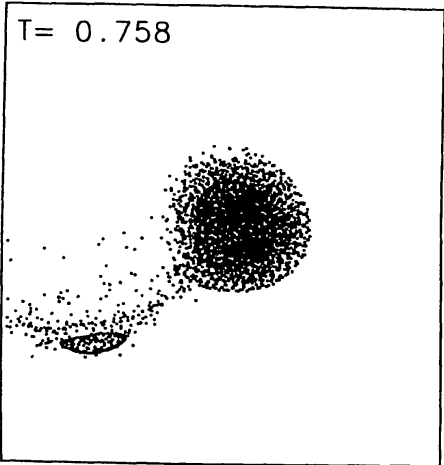
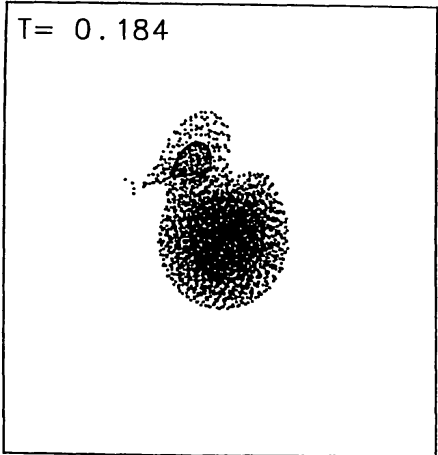
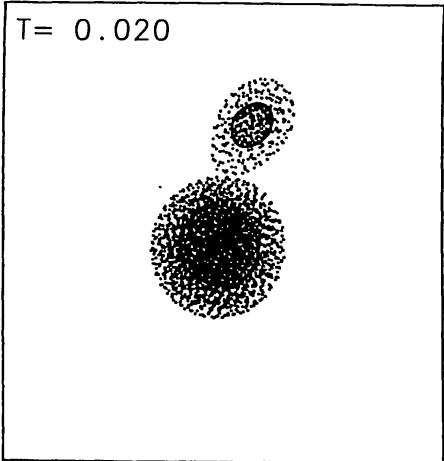
Cameron & Ward (1976) and Ward & Cameron (1978) were first to discuss this disk, but they assumed that it would cool and solidify rapidly to form a massive analogue of the rings of Saturn. Gravitational instabilities within the Roche limit would then provide sufficient dissipation (hence “viscosity”) to spread the disk out beyond the Roche limit, which would allow the Moon to form essentially cold. Thompson & Stevenson (1983, 1986) reconsidered this problem and found that the disk stays hot (liquid and gas) for 10^2 – 10^3 yr, and that it can spread in this time, allowing the Moon to form. The main points are these:

1. The cooling time for a disk of ~ 2 Moon masses placed within the Earth’s Roche limit is $\gtrsim 10^2$ yr.
2. The disk is a two-phase medium that is highly compressible and unstable with respect to gravitational “patch” instabilities, even when it is very hot.
3. The resulting turbulence and eddy viscosity allow the disk to spread in a time comparable to its cooling time. The disk self-regulates, maintaining its two-phase character because of the gravitational energy dissipated as the spreading proceeds.
4. Material spreads beyond the Roche limit, still maintaining its ability to undergo patch instabilities. Progressive cooling allows the instabilities to proceed all the way to the formation of protomoons. These protomoons are probably nearly fully molten and subsequently may coalesce to form the Moon.

We now elaborate on these points, basing our discussion on that of Thompson & Stevenson (1986). Consider 2 Moon masses ($2M_{\text{J}}$) of material spread out in a disk between $1.5R_{\oplus}$ and $3R_{\oplus}$, initially in a molten state with coexisting vapor. The characteristic cooling time of this material is at least

$$\frac{2M_{\text{J}}C_{\text{p}}T}{2\sigma T^4 A} \sim 10 \text{ yr}, \quad (16)$$

where A is the area of the disk ($\sim 7\pi R_{\oplus}^2$), the latent heat of condensation is ignored, and $T \sim 2000$ K is assumed. Actually, cooling turns out to be much slower (10^2 – 10^3 yr) because of the latent heat and because the disk creates its own heat as it spreads, allowing part of the mass to sink down into the Earth’s gravity field while another part climbs farther out in the gravitational potential (so as to conserve angular momentum). For example, suppose all the mass ($2M_{\text{J}}$) were initially at R_{i} . Now suppose half of this mass moves out to $1.5R_{\text{i}}$ and the other half to $[2 - (1.5)^{1/2}]^2 R_{\text{i}}$, thereby conserving angular momentum. The energy release is



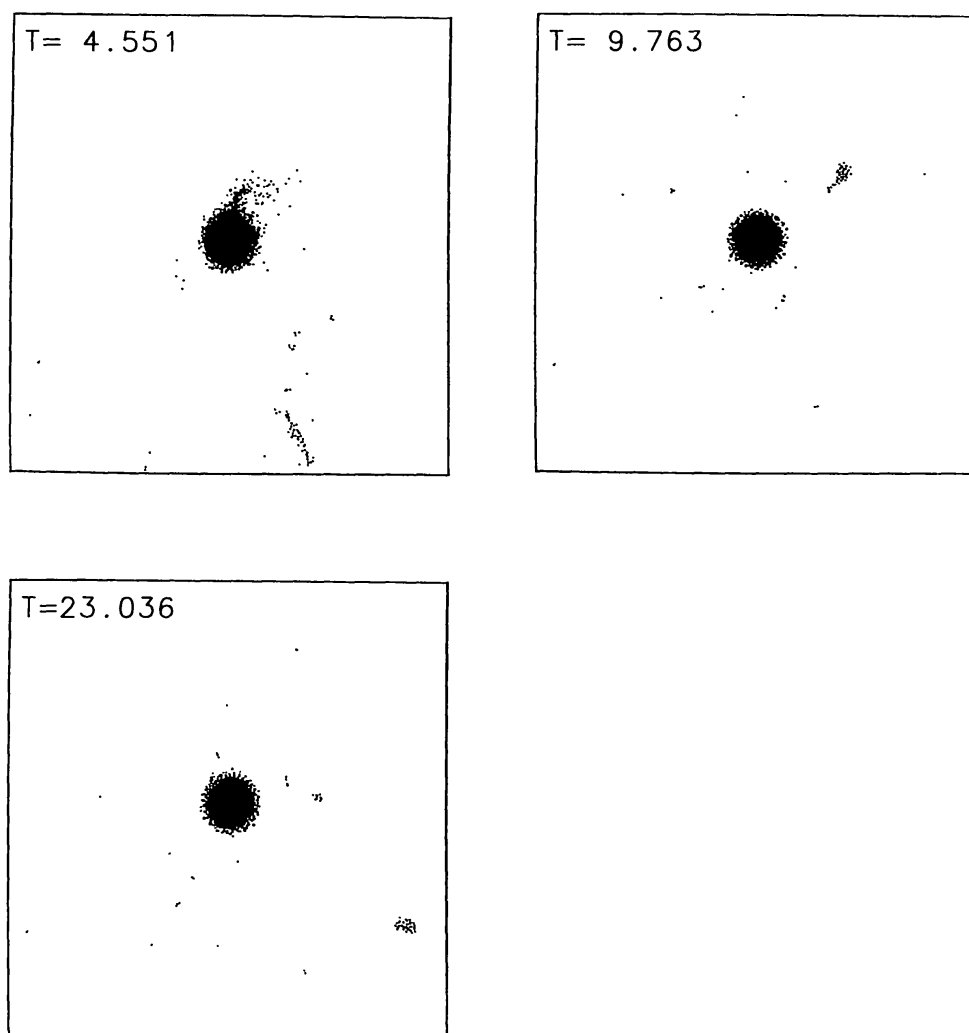


Figure 7 Numerical simulation by Benz et al (1986b) for impact of Earth by a body of mass $M_{\oplus}/7$. The initial conditions are 0 km s^{-1} relative velocity at infinity but the present Earth-Moon system angular momentum. The iron cores are circled. Time is given in hours from impact. The clump in the lower right corner of the last frame has exactly one Moon mass (35 particles). There is another $1/2$ – $1/3$ Moon mass in a disk around Earth. No iron is in orbit. Subsequent evolution after the final frame is not known. Results courtesy of W. Benz (Los Alamos National Laboratory).

$$\frac{GM_{\oplus}M_{\zeta}}{R_i} \left(\frac{1}{0.6} + \frac{1}{1.5} - 2 \right) \simeq \frac{0.33 GM_{\oplus}M_{\zeta}}{R_i}, \quad (17)$$

a value that is larger than $2M_{\zeta}C_pT$ by a factor of roughly five (assuming $R_i \approx 2R_{\oplus}$). We show below that something like this happens, so that much energy is available.

Now, the material emplaced in orbit settles into a disk on a dynamical time scale (\sim hours) after having undergone some adiabatic expansion and cooling since leaving the impact site. According to Figure 4, the cooling does *not* imply a large reduction in the mass fraction of vapor. Moreover, the material has great difficulty in cooling below a temperature ~ 2000 K because of the energy generated by gravitational instabilities. These arise because the medium contains bubbles and is therefore highly compressible. Consider, for example, a thin disk of surface density σ , sound speed c , and (local) orbital angular velocity Ω , assumed Keplerian. The dispersion relationship for waves of the form $\exp(i[kr + \omega t])$ is

$$\omega^2 = k^2c^2 - 2\pi G\sigma|k| + \Omega^2, \quad (18)$$

where r is the radial distance, and $kr \gg 1$ is assumed (e.g. Lin & Shu 1966; see also Goldreich & Ward 1973). Each term on the right-hand side has a simple physical explanation. The first term describes dispersionless sound waves ($\omega = ck$ if k is very large). The second term is negative and describes the possibility of gravitational collapse [analogous to the well-studied Jeans collapse in astrophysics (e.g. Chandrasekhar 1961, Ch. 13)]. The third term is the stabilizing effect of rotation (a combination of the effects of the Coriolis force and Keplerian shear). The important point is that because of the two-phase nature of the medium, the sound speed c is much smaller than the value appropriate to pure liquid or pure gas. This reduces one of the stabilizing terms and makes instability ($\omega^2 < 0$) much more likely. This is best understood graphically (see Figure 8).

The reduction in sound speed is very dramatic and is a well-known effect, for example, in a frothy air-water mixture, where the sound speed can be $\sim 20 \text{ m s}^{-1}$ compared with 1460 m s^{-1} for pure water (Kieffer 1977b). It is even more dramatic when the liquid and vapor are composed of the same material, since much of the compression is then accommodated by gas molecules within bubbles changing phase, which allows the bubbles to shrink. Thompson & Stevenson (1983, 1986) find, for example, that the sound speed of the two-phase medium $\text{SiO}_2(\text{l})\text{--}\text{SiO}(\text{g})\text{--}1/2 \text{ O}_2(\text{g})$ can be three orders of magnitude lower than pure $\text{SiO}_2(\text{l})$. This situation is illustrated in Figure 9. Equation (18) predicts that instability is possible for

$$\sigma > \sigma_{\text{crit}} \equiv \frac{\Omega c}{\pi G}. \quad (19)$$

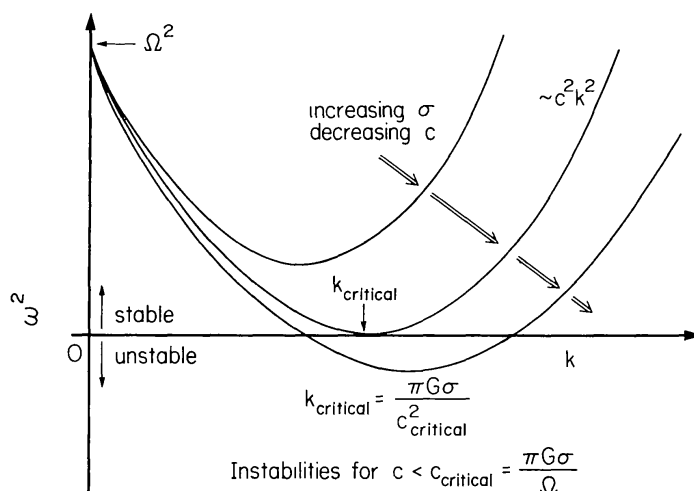


Figure 8 Dispersion relationship [Equation (18)] showing how a reduction in sound speed can cause instability (negative ω^2). This promotes turbulence in the disk.

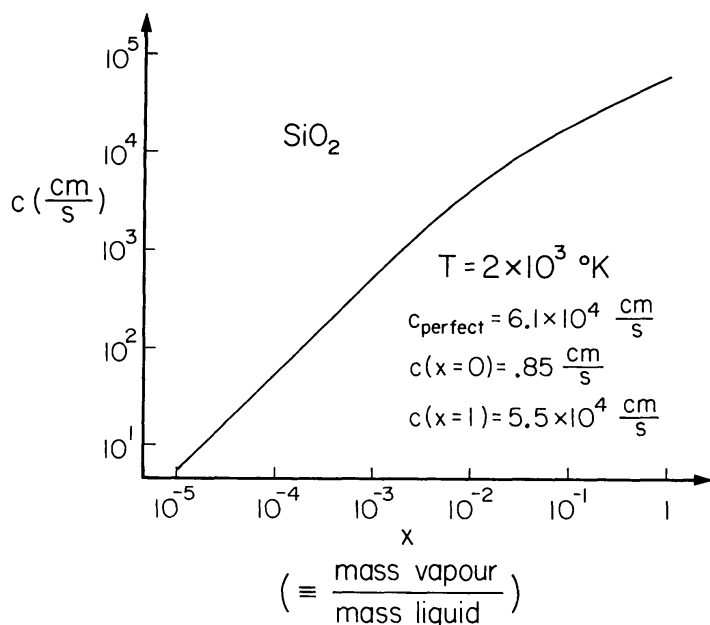


Figure 9 Sound speed of a two-phase medium (SiO_2 liquid and coexisting vapor). The ideal gas value is c_{perfect} . Notice that a medium that is mostly liquid but contains bubbles can have a very low sound speed. The value of c as $x \rightarrow 0$ is low and corresponds to the (artificial) case of an infinitesimal population of bubbles.

For a disk of mass $2M_{\oplus}$ distributed over $\sim 7\pi R_{\oplus}^2$, we have $\sigma \sim 1.5 \times 10^7 \text{ g cm}^{-2}$. Instability occurs for $c \lesssim 8 \times 10^3 \text{ cm s}^{-1}$. This is well below the value of 10^5 cm s^{-1} that roughly characterizes pure liquid or gas, but it is easily achievable in a two-phase medium. Notice that if Equation (19) is satisfied by a substantial factor, then the instabilities grow on a time scale $\lesssim 1$ orbital period, because $2\pi G\sigma|k|$ is comparable to Ω^2 and c^2k^2 for the

fastest-growing instabilities ($k \sim c/\Omega$). These instabilities are patches of collapsing fluid elements that *cannot* collapse all the way to a protomoon because they are within the Roche limit. Instead, they are tidally sheared. The resulting turbulent dissipation is at the expense of the gravitational energy of the disk [cf. Equation (17)] and can be envisaged as a turbulent viscosity that provides the coupling whereby the disk spreads. This idea of Ward & Cameron (1978) does not require a precise identification of the detailed kinematics of the turbulent viscosity, but rather it follows just from the energy argument stated here.

The disk self-regulates for $\sim 10^2$ yr, dissipating just enough energy to maintain itself at a marginally unstable state. It cannot dissipate less rapidly, since that would require cooling to a still more unstable state (one in which the mass in gas is a smaller fraction and the sound speed is even lower—see Figure 9). It cannot dissipate more rapidly, since that would heat the disk up, which makes the sound speed higher and stabilizes the disk. Consequently, the disk spreading time R^2/ν , where ν is the effective eddy viscosity, is of order 10^2 yr. This implies that $\nu \sim 10^9 \text{ cm}^2 \text{ s}^{-1}$.

Material spreads beyond the Roche limit, but because the disk is only marginally unstable, instabilities cannot proceed all the way to self-gravitating spheres immediately. The disk must continue to cool until the Roche limit evolves inward and material stranded beyond this limit can coalesce. This occurs at 10^2 – 10^3 yr after disk formation. The initial protomoons formed have a mass δm comparable to the “patch” instability mass suggested by marginal instability [Equation (18)]:

$$\begin{aligned}\delta m &\sim \pi(\pi/k)^2\sigma, \\ k &\sim \Omega^2/\pi G\sigma.\end{aligned}\tag{20}$$

This implies a body of mass 10^{20} – 10^{22} g (radius 10–100 km). These bodies are rather closely packed and undergo coagulation on a short time scale (\sim years). However, this process has not been modeled quantitatively.

Once even larger bodies are formed, tidal evolution is also rapid. The time that it takes a body of mass m to double its orbital radius from $3R_\oplus$ to $6R_\oplus$, due to tides excited on Earth, is

$$\tau_{\text{tidal}} \approx (500 \text{ yr}) \left(\frac{M_\oplus}{m} \right) \left(\frac{Q}{10} \right),\tag{21}$$

where Q is the tidal quality factor for Earth. This estimate is based on standard tidal theory (e.g. Lambeck 1980). Notice that the time scale is inversely proportional to mass, so smaller bodies can be swept up by larger bodies. Gravitational perturbations (excitation and eccentricity) will play a role also, just as they do in planetary accumulation, but the associated

time scale is much shorter than in terrestrial planet formation because the protolunar disk is very compact and has a short dynamical time scale (orbital period $\ll 1$ year).

The time scale to go from the initial instabilities to a mass of lunar magnitude has an important implication for the thermal state of the initial Moon. Smaller bodies can cool more rapidly, forming transient crusts. As these bodies coalesce, rather little extra energy is dissipated, and the resulting body is only partially molten. This has not been quantified.

CHEMISTRY OF THE IMPACT HYPOTHESIS

If the dynamics are uncertain, then the chemistry is more so, in part because it depends on the dynamics and in part because some aspects of the chemistry involve regimes that have received little attention in the laboratory. The important unresolved dynamical issues are these: Where does the lunar material come from? Is it mostly from the projectile (Benz et al 1986a,b) or from both the projectile and the target (Melosh & Sonett 1986, Kipp & Melosh 1986)? How much iron is put into orbit? In some of the simulations by Benz and coworkers, a substantial amount of iron may be put in orbit. How much material is lost, and is the loss differential (i.e. do some elements or molecules escape preferentially)? It is only on this latter point that a partial answer seems possible at present, but we discuss each issue in order.

The provenance of lunar-forming material is especially important to those who advocate a terrestrial origin (most notably Ringwood 1966, 1979, 1986a,b). In Ringwood's view, the Earth's mantle would be a distinctive reservoir that would differ from the mantle of the projectile unless the latter were sufficiently larger to have undergone the same differentiation and processing as the Earth's mantle. This places a lower bound on the projectile mass that is uncertain but surely larger than the mass of Mars, since models of Mars constructed by Ringwood are significantly different from Earth. As discussed earlier, the similarities and differences between the Moon and the Earth's mantle remain controversial. There seems little doubt that the projectile would be differentiated, since core formation is very effective in bodies with mass \gtrsim Mars (e.g. Stevenson 1981, 1983a), and its mantle might therefore be suitable for Moon-forming material. The relative amounts of projectile and target in the Moon have not been well established yet.

The question of iron is a serious one, given the low iron abundance in the Moon. If a disk that includes iron forms, then it is unlikely that the iron is selectively filtered out of the Moon-forming matter. It is comparably volatile to the silicates (in the sense that its vaporization temperature is

similar), but, more importantly, it would participate in the same patch instabilities that stir the silicate fluid and promote spreading of the disk. In other words, the liquid iron may form a sublayer at the equatorial plane of the disk but spread at the same rate as the rest of the disk. Wood (1986) discusses the data of Hashimoto (1983) indicating that iron (as FeO) is volatile, and he suggests that the iron could be evaporated away. As discussed below, this is insufficiently efficient to be significant. I conclude that if the Moon formed from a molten disk, then it would have the same iron content that at least the outer half of the disk has. The possibility (Benz et al 1986b) that the core of the projectile maintains its integrity and loses angular momentum, merging with the Earth, is a potential alternative.

After the impact, some material may be lost promptly on a hyperbolic trajectory. This material, possibly a jet of the most highly shocked material, would leave hydrodynamically (i.e. preserving composition). But what about “diffusive” loss from a disk emplaced in orbit? At a distance of $5R_{\oplus}$ from Earth center, a molecule of atom of mass μ in orbit needs only $E_{\text{esc}} = GM_{\oplus}\mu/10R_{\oplus}$ energy to escape. (This is only half the gravitational binding energy because of the Keplerian orbital motion.) Comparing this value with the kinetic energy $E_{\text{kin}} \simeq 3/2(kT)$, we have

$$\frac{E_{\text{kin}}}{E_{\text{esc}}} \simeq \left(\frac{4}{\mu}\right) \left(\frac{T}{2000}\right), \quad (22)$$

where μ is expressed in units of the proton mass. Hydrogen, derived from dissociation of H_2O advected to the upper levels of the disk, can clearly escape. The major constituents [MgO ($\mu = 40$), O_2 ($\mu = 32$), and SiO ($\mu = 44$)] cannot escape by a Jeans (thermally activated) process. Actually, the thermal escape would be in a cone of forward velocities, with the axis of the cone defined by the Keplerian orbital motion. In reality, the flow would be more like that of the solar wind because the exobase is at an enormous distance from the disk. It is possible that acoustic waves propagating into the hot-disk atmosphere may drive a wind (Thompson & Stevenson 1986). Since the turbulent velocities in the upper, rarefied layers of the disk are close to the sound speed, a rough estimate of the available energy input is

$$\dot{E} \simeq 2\pi R_d^2 \rho_g C_g^3, \quad (23)$$

where R_d is the disk radius and ρ_g , C_g are the gas density and sound speed, respectively. If all this energy goes into mass loss, then we have

$$\dot{M} \simeq 10^{-5} M_{\oplus} \left(\frac{\rho_g}{10^{-7} \text{ g cm}^{-3}} \right) \left(\frac{C_g}{10^4 \text{ cm s}^{-1}} \right) \text{ yr}^{-1}.$$

The disk photosphere is defined by the pressure level at which $p = g/K$ (g = disk self-gravity, K = opacity). As Thompson & Stevenson (1986) discuss, the opacity is very uncertain but is probably grain dominated, which suggests that $K \sim 1 \text{ cm}^2 \text{ g}^{-1}$. The corresponding density is $\sim 10^{-7} \text{ g cm}^{-3}$, and in view of the turbulent nature of the disk, there may already be enough energy input at this level to excite waves at the rate given in Equation (23). For this upper bound, the mass loss is still only about 0.1% of the mass of the *Moon* and hydrodynamic in nature. It is unlikely that elemental and isotopic fractionation occurs, as suggested by Cameron (1983), except perhaps for a preferential loss of hydrogen, derived from dissociation of water. Water would be present as a molecular species high in the disk atmosphere, advected there by silicate droplets and then exsolved from the silicates in the low-pressure region of grain formation. At $T \sim 2000 \text{ K}$, $P(\text{H}_2\text{O}) \lesssim 10^{-6} \text{ bar}$, plenty of atomic hydrogen would be available by thermal dissociation.

To summarize, the disk is essentially a *closed system*. With the possible exception of hydrogen (and hence water), the loss of material by outflow from the disk is neither large nor capable of differentiation. This partially validates attempts to compare chemically the Moon with the mantle of the Earth or with a large projectile.

COMMENTS AND QUESTIONS

It should be clear that the collision hypothesis of lunar origin is still in a primitive form. This is readily apparent by returning to the 10 propositions listed early in this review and critically assessing their merits and uncertainties:

1. Inadequacies certainly exist in the other “conventional” lunar origin scenarios (fission, capture, binary accretion). But this is not an argument in favor of megaimpact; it may simply be that we have had insufficient time to recognize all the shortcomings of the current bandwagon.
2. The theoretical evidence certainly points to large impacts during planetary formation. However, no simulations exist that cover the entire sequence of events from the formation of small planetesimals all the way to the final planets. Even so, this remains as one of the better-justified propositions.
3. Large impacts should certainly be very different dynamically than small impacts. Quantification of this proposition remains difficult and imperfect, however, and three-dimensional hydrodynamic simulations are still rather primitive and incomplete. Much more work is needed

to characterize the range of outcomes and the efficiency of orbital injection.

4. Estimates of vaporization accompanying large impacts vary widely because they involve extrapolations of existing equation of states or uncertain theoretical modeling (e.g. to evaluate electronic excitations). Current estimates arising from the three-dimensional computer simulations may be incorrect. As a consequence, the role of pressure gradients in the postimpact flow is not yet understood.
5. The impact trauma is very substantial for Earth and almost certainly causes a hot, silicate vapor atmosphere and an underlying magma ocean. But is there a diagnostic consequence of this in Earth history (e.g. the geochemistry of the deep mantle)? More work is needed on the behavior of magma oceans.
6. The material placed in orbit may be clumpy or dispersed. Current calculations do not demonstrate conclusively that a disk will form. Numerical simulations have to be carried out to longer times.
7. The origin of the orbital material is unclear. It may come mainly from the projectile, or it may come partly or wholly from the target. It is also not clear whether this orbital material is devoid of iron from the core of the projectile. Again, more simulations are needed to understand this better.
8. It seems likely that the dynamics of lunar formation are *very fast* irrespective of the details. This seems to be one of the safer propositions. Nevertheless, detailed calculations do not exist.
9. Despite the high energy of a giant impact, the thermal energy of the resulting material is small compared with the escape energy, and loss during the history of the disk is small. Consequently, the system is almost "closed," and this proposition emerges as one of the better-justified ones.
10. Crude estimates of orbital injection suggest that it should work for lower-mass impactors but with gradually diminishing efficiency. Consequently, we must *expect* that the Moon is a combination of several protomoons. The number cannot be too great, since otherwise the angular momentum problem is very serious. (Roughly speaking, the protomoons should have comparable likelihoods of prograde and retrograde orbits.) Formation of the Moon from a *large* number of impacts (e.g. Ringwood 1986a,b) does not seem tenable. We can rationalize the existence of only one Moon by arguing that the largest impact occurred late in Earth accretion. However, more work is needed on the question of whether smaller protomoons are either swept up or lose angular momentum and collide with Earth.

Another way of assessing lunar origin by impact is to ask how well the

scenario explains the data. Such an assessment is presented in Table 2. A quick look confirms that the uncertainties dominate the positive attributes. The best that can be said for the impact hypothesis is that it has no *strongly* negative attributes at present.

Finally, what are the implications of giant impacts for comparative planetology? Here, we see many attractive possibilities. Cameron (1983) has suggested that the differences between Earth and Venus may be partly attributed to the Earth (but not Venus) suffering a very large impact. For example, Cameron speculates that a large impact could “blow away” a massive CO₂ atmosphere. Our discussion here does not provide very strong support for this idea, but it merits further attention. An equally interesting idea concerns Mars. Since the escape velocity from Mars is too low for significant vaporization of incoming impactors with low velocity at infinity, and since fast-moving projectiles would cause ejecta to escape hyperbolically, we conclude tentatively that Mars (or any small planet) should never have a substantial satellite, at least by the impact process. Indeed, Mars has only two small satellites, both plausibly derived by capture (see discussion in Stevenson et al 1986). Last, but certainly not least, the Uranian system is a natural candidate for giant-impact effects. An impact by a body of ~ 2 Earth masses may have stirred the interior (lowering the moment of inertia) and caused the formation of a disk from which the satellites formed. The compositions of these satellites [see Smith et al (1986)

Table 2 Comparison between constraints from lunar data and the impact scenario

Constraints from lunar data	Impact model
1. Mass	Need projectile \gtrsim Mars mass and $\sim 10\%$ efficiency of orbital injection
2. Angular momentum	Implies one or (at most) a few large impacts; many small impacts have canceling angular momenta
3. Low iron content	Iron in projectile must avoid orbital injection; it is not clear whether this is possible
4. Volatile depletion	Uncertain predictions; hydrogen (hence water) probably lost but other volatiles may be partly retained
5. Trace elements and other chemical constraints	Lunar-forming material is an uncertain mix of projectile and Earth mantles
6. High initial temperatures (putative lunar magma ocean)	Readily provided by impact energy release
7. Orbital evolution (tidal theory)	Initial moon may be in equatorial plane. Tidal theory does not clearly preclude this

for Voyager data] are clearly more rock rich than the Saturnian satellites, an observation possibly consistent with an impact origin (Stevenson 1984b, Stevenson & Lunine 1986). The Uranian system is an enticing locale for those who tire with modeling their own backyard.

ACKNOWLEDGMENTS

I thank W. Benz, A. G. W. Cameron, and H. J. Melosh for discussing their unpublished work, A. E. Ringwood and G. W. Wetherill for useful conversations, and the organizers of the Kona Conference (Hartmann et al 1986) for motivating much of this effort. This work was supported by NASA Planetary Geophysics grant NAGW-185 and is contribution number 4348 from the Division of Geological and Planetary Sciences, California Institute of Technology, Pasadena, CA 91125.

Literature Cited

- Ahrens, T. J., O'Keefe, J. D. 1972. Shock melting and vaporization of lunar rocks and minerals. *The Moon* 4: 214–49
- Allen, R. T. 1967. Equation of state of rocks and minerals. *Rep. GAMD-7834*, General Atomic Div., General Dynamics, San Diego, Calif.
- Anderson, D. L. 1972. The origin of the Moon. *Nature* 239: 263–65
- Arrhenius, S. 1908. *Worlds in the Making*. New York/London: Harper & Brothers
- Benz, W., Slattery, W. L., Cameron, A. G. W. 1986a. The origin of the Moon and the single impact hypothesis, I. *Icarus* 66: 515–35
- Benz, W., Slattery, W. L., Cameron, A. G. W. 1986b. The origin of the Moon: 3D numerical simulations of a giant impact. *Lunar Planet. Sci. XVII*, pp. 40–41 (Abstr.)
- Bills, B. G., Ferrari, A. J. 1977. A harmonic analysis of lunar topography. *Icarus* 31: 244–59
- Boss, A. P. 1986. The origin of the Moon. *Science* 231: 341–45
- Cameron, A. G. W. 1972. Orbital eccentricity of Mercury and the origin of the moon. *Nature* 240: 299–300
- Cameron, A. G. W. 1983. Origin of the atmospheres of the terrestrial planets. *Icarus* 56: 195–201
- Cameron, A. G. W. 1985a. Formation and evolution of the primitive solar nebula. In *Protostars and Planets II*, ed. D. C. Black, M. S. Matthews, pp. 1073–99. Tucson: Univ. Ariz. Press
- Cameron, A. G. W. 1985b. Formation of the prelunar accretion disk. *Icarus* 62: 319–27
- Cameron, A. G. W. 1986. The impact theory for origin of the Moon. See Hartmann et al 1986, pp. 609–16
- Cameron, A. G. W., Ward, W. R. 1976. The origin of the Moon. *Lunar Sci. VII*, pp. 120–22 (Abstr.)
- Chandrasekhar, S. 1961. *Hydrodynamic and Hydromagnetic Stability*. Oxford: Oxford Univ. Press
- Darwin, G. H. 1880. On the secular change in elements of the orbit of a satellite revolving around a tidally distorted planet. *Philos. Trans. R. Soc. London* 171: 713–891
- Davies, G. F. 1985. Heat deposition and retention in a solid planet growing by impacts. *Icarus* 63: 45–68
- Drake, M. J. 1983. Geochemical constraints on the origin of the Moon. *Geochim. Cosmochim. Acta* 47: 1759–67
- Durisen, R. H., Scott, E. H. 1984. Implications of recent numerical calculations for the fission theory of the origin of the Moon. *Icarus* 58: 153–58
- Ganapathy, R., Anders, E. 1974. Bulk composition of the Moon and Earth estimated from meteorites. *Proc. Lunar Sci. Conf., 5th*, pp. 1181–1206
- Gault, D. E., Quaide, W. L., Oberbeck, V. R. 1968. Impact cratering mechanics and structures. In *Shock Metamorphism of Natural Materials*, ed. B. M. French, N. M. Short, pp. 87–99. Baltimore: Mono
- Gingold, R. A., Monaghan, J. J. 1979.

- Binary fission in damped rotating polytropes. *Mon. Not. R. Astron. Soc.* 188: 39–44
- Goldreich, P. 1966. History of the lunar orbit. *Rev. Geophys.* 4: 411–39
- Goldreich, P., Ward, W. R. 1973. The formation of planetesimals. *Astrophys. J.* 183: 1051–61
- Greenberg, R. 1982. Planetesimals to planets. In *Formation of Planetary Systems*, ed. A. Brahic, pp. 515–69. Toulouse, Fr: Lepadues Ed.
- Harris, A. W. 1977. An analytical theory of planetary rotation rates. *Icarus* 31: 168–74
- Harris, A. W., Kaula, W. M. 1975. A co-accretional model of satellite formation. *Icarus* 24: 516–24
- Harris, A. W., Ward, W. R. 1982. Dynamical constraints on the formation and evolution of planetary bodies. *Ann. Rev. Earth Planet. Sci.* 10: 61–108
- Hartmann, W. K., Davis, D. R. 1975. Satellite-sized planetesimals. *Icarus* 24: 504–15
- Hartmann, W. K., Phillips, R. J., Taylor, G. J., eds. 1986. *Origin of the Moon*. Houston: Lunar Planet. Sci. Inst. 781 pp.
- Hashimoto, A. 1983. Evaporation metamorphism in the early solar nebula—evaporation experiments on the melt FeO-MgO-SiO₂-CaO-Al₂O₃ and chemical fractionations of primitive materials. *Geochim. J.* 17: 111–45
- Hayashi, C., Nakazawa, K., Nakagawa, Y. 1985. Formation of the solar system. In *Protostars and Planets II*, ed. D. C. Black, M. S. Matthews, pp. 1100–53. Tucson: Univ. Ariz. Press
- Hofmeister, A. M. 1983. Effect of a Hadean terrestrial magma ocean on crust and mantle evolution. *J. Geophys. Res.* 88: 4963–83
- Holmes, N. C., Radousky, H. B., Moss, M. J., Nellis, W. J., Henning, S. 1984. Silica at ultrahigh temperature and expanded volume. *Appl. Phys. Lett.* 45: 626–27
- Horedt, G. P. 1985. Late stages of planetary accretion. *Icarus* 64: 448–70
- JANAF Thermochemical Tables. 1971. *Publ. No. NSRDS-NBS37*. Washington, DC: Natl. Bur. Stand. 1141 pp.
- Jeanloz, R., Thompson, A. B. 1983. Phase transitions and mantle discontinuities. *Rev. Geophys. Space Phys.* 21: 51–75
- Kaula, W. M. 1979. Thermal evolution of Earth and Moon growing by planetesimal impacts. *J. Geophys. Res.* 84: 999–1008
- Kaula, W. M. 1980. The beginning of the Earth's thermal evolution. In *The Continental Crust and its Mineral Deposits*. *Geol. Assoc. Can. Spec. Pap. No. 20*, pp. 25–34
- Kaula, W. M., Beachey, A. E. 1986. Mechanical models of close approaches and collisions of large protoplanets. See Hartmann et al 1986, pp. 567–76
- Kieffer, S. W. 1975. Droplet chondrules. *Science* 189: 333–40
- Kieffer, S. W. 1977a. Impact conditions required for formation of melt by jetting in silicates. See Roddy et al 1977, pp. 751–69
- Kieffer, S. W. 1977b. Sound speed in liquid-gas mixtures: water-air and water-steam. *J. Geophys. Res.* 82: 2895–2904
- Kieffer, S. W. 1979. Thermodynamics and lattice vibrations of minerals. 1. Mineral heat capacities and their relationships to simple lattice vibrational models. *Rev. Geophys. Space Phys.* 17: 1–19
- Kieffer, S. W. 1982. Dynamics and thermodynamics of volcanic eruptions: implications for the plumes on Io. In *Satellites of Jupiter*, ed. D. Morrison, pp. 647–723. Tucson: Univ. Ariz. Press
- Kinslow, R., ed. 1970. *High Velocity Impact Phenomena*. New York: Academic. 579 pp.
- Kipp, M. E., Melosh, H. J. 1986. Origin of the Moon: a preliminary numerical study of colliding planets. *Lunar Planet. Sci. XVII*, pp. 420–21
- Kreutzberger, M. E., Drake, M. J., Jones, J. H. 1986. Origin of the Earth's Moon: constraints from alkali volatile trace elements. *Geochim. Cosmochim. Acta* 50: 91–98
- Krieger, F. J. 1967. The thermodynamics of the magnesium silicate/magnesium-silicon-oxygen vapor system. *Memo. RM-5337-PR*, Rand Corp., Santa Monica, Calif.
- Lambeck, K. 1980. *The Earth's Variable Rotation*. Cambridge: Cambridge Univ. Press
- Lambeck, K. 1986. Banded iron formations. *Nature* 320: 574
- Lecar, M., Aarseth, S. J. 1986. A numerical simulation of the formation of the terrestrial planets. *Astrophys. J.* 305: 564–79
- Lin, C. C., Shu, F. H. 1966. On the spiral arms of disk galaxies. II. Outline of a theory of density waves. *Proc. Natl. Acad. Sci. USA* 55: 229–34
- Lucy, L. B. 1977. A numerical approach to the testing of the fission hypothesis. *Astron. J.* 82: 1013–24
- Lynden-Bell, D., Pringle, J. E. 1974. The evolution of viscous disks and the origin of nebula variables. *Mon. Not. R. Astron. Soc.* 168: 603–37
- Lyzenga, G. A., Ahrens, T. J. 1980. Shock temperature measurements in Mg₂SiO₄ and SiO₂ at high pressure. *Geophys. Res. Lett.* 7: 141–44
- Lyzenga, G. A., Ahrens, T. J., Mitchell, A.

- C. 1983. Shock temperatures of SiO_2 and their geophysical implications. *J. Geophys. Res.* 88: 2431–44
- Melosh, H. J., Sonett, C. P. 1986. When worlds collide: jetted vapor plumes and the Moon's origin. See Hartmann et al 1986, pp. 621–42
- Mizuno, H., Boss, A. P. 1985. Tidal disruption of dissipative planetesimals. *Icarus* 63: 109–33
- Nakazawa, K., Komuro, T., Hayashi, C. 1983. Origin of the Moon: capture by gas drag of the Earth's primordial atmosphere. *The Moon and the Planets* 28: 311–27
- Newsom, H. E. 1984. The lunar core and the origin of the Moon. *Eos, Trans. Am. Geophys. Union* 65: 369–70
- Ohtani, E. 1985. The primordial terrestrial magma ocean and its implications for stratification of the mantle. *Phys. Earth Planet. Inter.* 38: 70–80
- O'Keefe, J. D., Ahrens, T. J. 1977. Impact-induced energy partitioning, melting and vaporization on terrestrial planets. *Proc. Lunar Sci. Conf.*, 8th, pp. 3357–74
- O'Keefe, J. D., Ahrens, T. J. 1982. The interaction of the Cretaceous/Tertiary extinction bolide with the atmosphere, ocean and solid Earth. See Silver & Schultz 1982, pp. 103–20
- Öpik, E. J. 1972. Comments on lunar origin. *Ir. Astron. J.* 10: 190–238
- Orphal, D. L., Borden, W. F., Larson, S. A., Schultz, P. H. 1980. Impact melt generation and transport. *Proc. Lunar Planet. Sci. Conf.*, 11th, pp. 2309–23
- Rigden, S. M., Ahrens, T. J. 1981. Impact vaporization and lunar origin. *Lunar Planet. Sci. XII*, pp. 885–87 (Abstr.)
- Ringwood, A. E. 1966. Chemical evolution of the terrestrial planets. *Geochim. Cosmochim. Acta* 30: 41–104
- Ringwood, A. E. 1979. *Origin of the Earth and Moon*. New York: Springer-Verlag. 295 pp.
- Ringwood, A. E. 1986a. The making of the Moon. *Lunar Planet. Sci. XVII*, pp. 714–15 (Abstr.)
- Ringwood, A. E. 1986b. Composition and origin of the Moon. See Hartmann et al, pp. 673–98
- Roddy, D. J., Pepin, R. O., Merrill, R. B., eds. 1977. *Impact and Explosion Cratering*. Houston: Lunar Sci. Inst. 1299 pp.
- Ruff, O., Schmidt, P. 1921. Die dampfdrucke der oxyde des siliciums, aluminiums, calciums und magnesiums. *Z. Anorg. Allg. Chem.* 117: 172–90
- Ruskol, E. L. 1960. Origin of the Moon. I. *Sov. Astron. AJ* 4: 657–68
- Ruskol, E. L. 1982. Origin of planetary satellites. *Izvestiya Earth Phys.* 18: 425–33
- Safronov, V. S. 1966. Sizes of the largest bodies falling onto the planets during their formation. *Sov. Astron. AJ* 9: 987–91
- Safronov, V. S. 1969. Evolution of the protoplanetary cloud and formation of the Earth and planets. *NASA TT F-677*
- Silver, L. T., Schultz, P. H., eds. 1982. *Geological Implications of Impacts of Large Asteroids and Comets on the Earth*. *GSA Spec. Pap. No. 190*. Boulder, Colo: Geol. Soc. Am.
- Smith, B. A., and 39 others. 1986. Voyager II in the Uranian system: imaging science results. *Science* 233: 43–64
- Stanyukovich, K. P. 1950. Elements of the physical theory of meteors and the formation of meteor craters. *Meteoritika* 7: 39–62
- Stevenson, D. J. 1981. Models of the Earth's core. *Science* 214: 611–19
- Stevenson, D. J. 1983a. The nature of the Earth prior to the oldest known rock record (the Hadean Earth). In *Origin and Evolution of the Earth's Earliest Biosphere*, ed. J. W. Schopf, Ch. 2. Princeton, NJ: Princeton Univ. Press
- Stevenson, D. J. 1983b. Anomalous bulk viscosity of two-phase fluids and implications for planetary interiors. *J. Geophys. Res.* 88: 2445–53
- Stevenson, D. J. 1984a. Lunar origin from impact on the Earth: is it possible? *Conf. Origin of the Moon Abstr. Vol., LPI Contrib. No. 540*, p. 60
- Stevenson, D. J. 1984b. Composition, structure and evolution of Uranian and Neptunian satellites. In *Uranus and Neptune, NASA Conf. Publ. No. 2330*, pp. 405–23
- Stevenson, D. J., Lunine, J. I. 1986. Mobilization of cryogenic ice in outer solar system satellites. *Nature* 323: 46–48
- Stevenson, D. J., Harris, A. W., Lunine, J. I. 1986. Origins of satellites. In *Satellites*, ed. J. Burns, pp. 39–88. Tucson: Univ. Ariz. Press
- Stewart, G. R., Wetherill, G. W. 1986. New formulas for the evolution of planetesimal velocities. *Lunar Planet. Sci. XVII*, pp. 827–28 (Abstr.)
- Taylor, S. R. 1986. The origin of the Moon: geochemical considerations. See Hartmann et al 1986, pp. 125–44
- Taylor, S. R., Jakes, P. 1974. The geochemical evolution of the Moon. *Proc. Lunar Sci. Conf.*, 5th, pp. 1287–1305
- Thompson, A. C., Stevenson, D. J. 1983. Two-phase gravitational instabilities in thin disks with applications to the origin of the Moon. *Lunar Planet. Sci. XIV*, pp. 787–88 (Abstr.)
- Thompson, A. C., Stevenson, D. J. 1986. Two-phase gravitational instabilities in

- thin disks with application to the origin of the Moon. Submitted for publication
- Tillotson, J. H. 1962. Metallic equations of state for hypervelocity impact. *Rep. GA-3216*, General Atomic Div., General Dynamics, San Diego, Calif.
- Walker, J. C. G., Zahnle, K. J. 1986. Lunar nodal tide and distance to the Moon during the Precambrian. *Nature* 320: 600-2
- Walker, J. C. G., Klein, C., Stevenson, D. J., Walter, M. R. 1983. Environmental evolution of the Archean early Proterozoic Earth. In *Origin and Evolution of the Earth's Earliest Biosphere*, ed. J. W. Schopf, pp. 32-40. Princeton, NJ: Princeton Univ. Press
- Wänke, H., Dreibus, G. 1986. Geochemical evidence for formation of the moon by impact induced fission of the proto-Earth. See Hartmann et al 1986, pp. 649-72
- Ward, W. R., Cameron, A. G. W. 1978. Disk evolution within the Roche limit. *Lunar Planet. Sci. IX*, pp. 1205-7 (Abstr.)
- Warren, P. H. 1985. The magma ocean concept and lunar evolution. *Ann. Rev. Earth Planet. Sci.* 13: 201-40
- Weidenschilling, S. J., Greenberg, R., Chapman, C. R., Herbert, F., Davis, D. R., et al. 1986. Origin of the Moon from a circumterrestrial disk. See Hartmann et al 1986, pp. 731-62
- Wetherill, G. W. 1975. Possible slow accretion of the Moon and its thermal and petrological consequences. In *Origins of Mare Basalts*, pp. 184-88. Houston: Lunar Sci. Inst.
- Wetherill, G. W. 1980. Formation of the terrestrial planets. *Ann. Rev. Astron. Astrophys.* 18: 77-113
- Wetherill, G. W. 1985. Occurrence of giant impacts during the growth of the terrestrial planets. *Science* 228: 877-79
- Wetherill, G. W., Cox, L. P. 1984. The range of validity of the two-body approximation in models of terrestrial planet accumulation. I. Gravitational perturbations. *Icarus* 60: 40-55
- Wetherill, G. W., Cox, L. P. 1985. The range of validity of the two-body approximation in models of terrestrial planet accumulation. II. Gravitational cross sections and runaway accretion. *Icarus* 63: 290-303
- Wetherill, G. W., Stewart, G. R. 1986. The early stages of planetesimal accumulation. *Lunar Planet. Sci. XVII*, p. 939 (Abstr.)
- Wood, J. A. 1986. Moon over Mauna Loa: a review of hypotheses of formation of Earth's Moon. See Hartmann et al 1986, pp. 17-55
- Wood, J. A., Mitler, H. E. 1974. Origin of the Moon by a modified capture mechanism, or half a loaf is better than a whole one. *Proc. Lunar Sci. Conf., 5th*, pp. 851-53