

# ISIMA lectures on celestial mechanics. 2

Scott Tremaine, Institute for Advanced Study

July 2014

## 1. Numerical orbit integration

Effective algorithms for numerical orbit integration are among the most important tools for planetary dynamics.

In most planetary problems, the dynamics is that of a Hamiltonian system: with  $N + 1$  particles (planets plus host star) there are  $3(N + 1)$  coordinates that form the components of a vector  $\mathbf{q}(t)$ , and  $3(N + 1)$  components of the corresponding momentum  $\mathbf{p}(t)$ . These vectors satisfy Hamilton's equations,

$$\dot{\mathbf{q}} = \frac{\partial H}{\partial \mathbf{p}} \quad ; \quad \dot{\mathbf{p}} = -\frac{\partial H}{\partial \mathbf{q}}, \quad (1)$$

For simplicity we shall restrict our attention first to the motion of one body of unit mass in a fixed potential, so that the Hamiltonian has the form  $H(\mathbf{q}, \mathbf{p}) = \frac{1}{2}\mathbf{p}^2 + \Phi(\mathbf{q})$ . Given a phase-space position  $(\mathbf{q}, \mathbf{p})$  at time  $t$ , and a **timestep**  $h$ , we require an algorithm—an **integrator**—that generates a new position  $(\mathbf{q}', \mathbf{p}')$  that approximates the true position at time  $t' = t + h$ . The **order** of the integrator is  $k$  if the error after a single timestep varies as  $O(h^{k+1})$ . If the error accumulates from step to step then the expected error after a fixed time interval  $\Delta t = Nh$  should vary as  $O(Nh^{k+1}) = \Delta t O(h^k)$ .

There are many conventional integrators described in numerical analysis books (Runge-Kutta, predictor-corrector, Bulirsch-Stoer, multistep and multivalued methods, etc.). These are ideal for short integrations (up to a few hundred orbits). For longer integrations (up to billions of orbits), conventional integrators are less successful. To see this let us imagine that we want to follow the Earth's orbit (which we take to be circular for simplicity) for  $10^{10}$  years, with a timestep  $h$  of 0.01 years (3.6 days). Suppose that the integrator produces an error in phase of  $\epsilon_\phi$  radians per step, and that this error accumulates from one step to the next. Then the phase error after time  $T$  is  $\Delta\phi = \epsilon_\phi(T/h) = 10^{12}\epsilon_\phi$  at the end of the integration.

Now suppose that in addition the integrator makes a fractional error in energy or semi-major axis of  $\epsilon_a$  per step, which accumulates from step to step (sometimes called **numerical dissipation**). Since the mean motion  $n$  is proportional to  $a^{-3/2}$  this corresponds to

a fractional error in  $n$  of  $\frac{3}{2}\epsilon$  per step. The orbital phase grows in time as  $d\phi/dt = n$  so  $d^2\phi/dt^2 = \frac{3}{2}\epsilon/h$ . Since  $d^2\phi/dt^2 = 0$  in the real system, the phase error grows as  $d^2\Delta\phi/dt^2 = dn/dt = \frac{3}{2}\epsilon n/h$  which is easily integrated (so long as  $\epsilon$  is small enough so that  $n$  is approximately constant) to give  $\Delta\phi = \frac{3}{4}\epsilon n T^2/h$ . For  $T = 10^{10}$  years and  $h = 0.01$  years this yields  $\Delta\phi = 5 \times 10^{22}\epsilon$ , larger by more than  $10^{10}$  if  $\epsilon \sim \epsilon_\phi$ .

As this example shows, energy errors are far more important than phase errors in long integrations, but conventional integrators do not distinguish energy drifts from phase drifts, or oscillatory errors from cumulative ones. Thus, for long integrations it is almost always better to use **geometric integrators**. These are integrators that preserve some of the geometric properties of the phase-space flow exactly. Conventional integrators try to make the position in phase space at the end of one timestep as close as possible to the true position; geometric integrators give up some of this accuracy to ensure that the overall properties of the flow in phase space are the same in the integrated trajectory and the true trajectory. The philosophy is that preserving the phase-space structure of the flow determined by the real dynamical system is more important than minimizing the one-step error.

For example, time-reversible integrators will return an orbit to its exact starting point in phase space if the velocity is reversed. Symplectic integrators ensure that the transformation from initial to final position is symplectic or canonical, a property shared by the original Hamiltonian system. In particular, any integrator that obeys the equations of motion for some Hamiltonian is symplectic.

A good general reference is Hairer et al., *Geometric Numerical Integration*.

### 1.1. Symplectic integrators

**Euler’s method** Consider a Hamiltonian  $H(\mathbf{q}, \mathbf{p}) = \frac{1}{2}\mathbf{p}^2 + \Phi(\mathbf{q})$ . Hamilton’s equations read

$$\dot{\mathbf{q}} = \frac{\partial H_h}{\partial \mathbf{p}} = \mathbf{p} \quad ; \quad \dot{\mathbf{p}} = -\frac{\partial H_h}{\partial \mathbf{q}} = -\nabla\Phi(\mathbf{q}). \quad (2)$$

The simplest integration scheme is **Euler’s method**,

$$\mathbf{p}' = \mathbf{p} - h\nabla\Phi(\mathbf{q}) \quad ; \quad \mathbf{q}' = \mathbf{q} + h\mathbf{p}. \quad (3)$$

This is a first-order method. The error properties of Euler’s method are pretty bad, so we introduce it only as a foil to more sophisticated methods. In particular, it is not symplectic.

**Modified Euler integrator** Replace the original Hamiltonian  $H(\mathbf{q}, \mathbf{p}) = \frac{1}{2}\mathbf{p}^2 + \Phi(\mathbf{q})$  by the time-dependent Hamiltonian

$$H_h(\mathbf{q}, \mathbf{p}, t) = \frac{1}{2}\mathbf{p}^2 + \Phi(\mathbf{q})\delta_h(t), \quad \text{where} \quad \delta_h(t) \equiv h \sum_{j=-\infty}^{\infty} \delta(t - jh) \quad (4)$$

is an infinite series of delta functions. Averaged over a time interval that is long compared to  $h$ ,  $\langle H_h \rangle \simeq H$ , so the trajectories determined by  $H_h$  should approach those determined by  $H$  as  $h \rightarrow 0$ .

Hamilton's equations for  $H_h$  read

$$\dot{\mathbf{q}} = \frac{\partial H_h}{\partial \mathbf{p}} = \mathbf{p} \quad ; \quad \dot{\mathbf{p}} = -\frac{\partial H_h}{\partial \mathbf{q}} = -\nabla\Phi(\mathbf{q})\delta_h(t). \quad (5)$$

We now integrate these equations from  $t = -\epsilon$  to  $t = h - \epsilon$ , where  $0 < \epsilon \ll h$ . Let the system have coordinates  $(\mathbf{q}, \mathbf{p})$  at time  $t = -\epsilon$ , and first ask for its coordinates  $(\bar{\mathbf{q}}, \bar{\mathbf{p}})$  at  $t = +\epsilon$ . During this short interval  $\mathbf{q}$  changes by a negligible amount, and  $\mathbf{p}$  suffers a kick governed by the second of equations (5). Integrating this equation from  $t = -\epsilon$  to  $\epsilon$  is trivial since  $\mathbf{q}$  is fixed, and we find

$$\bar{\mathbf{q}} = \mathbf{q} \quad ; \quad \bar{\mathbf{p}} = \mathbf{p} - h\nabla\Phi(\mathbf{q}); \quad (6)$$

this is called a **kick step** because the momentum changes but the position does not. Equation (6) defines a nonlinear operator, the kick operator  $\mathbf{K}_h$ , which maps  $(\mathbf{q}, \mathbf{p})$  onto  $(\bar{\mathbf{q}}, \bar{\mathbf{p}})$ . Next, between  $t = +\epsilon$  and  $t = h - \epsilon$ , the value of the delta function is zero, so the system has constant momentum, and Hamilton's equations yield for the coordinates at  $t = h - \epsilon$

$$\mathbf{q}' = \bar{\mathbf{q}} + h\bar{\mathbf{p}} \quad ; \quad \mathbf{p}' = \bar{\mathbf{p}}; \quad (7)$$

this is called a **drift step** because the position changes but the momentum does not. Equation (6) defines a linear operator, the drift operator  $\mathbf{D}_h$  which maps  $(\bar{\mathbf{q}}, \bar{\mathbf{p}})$  onto  $(\mathbf{q}', \mathbf{p}')$ , Combining these results, we find that over a timestep  $h$  starting at  $t = -\epsilon$  the Hamiltonian  $H_h$  generates a map  $(\mathbf{q}, \mathbf{p}) \rightarrow (\mathbf{q}', \mathbf{p}')$  given by

$$\mathbf{p}' = \mathbf{p} - h\nabla\Phi(\mathbf{q}) \quad ; \quad \mathbf{q}' = \mathbf{q} + h\mathbf{p}', \quad (8)$$

and called the **kick-drift** or **modified Euler** method (compare this to Euler's method, eq. 3). This method corresponds to the sequence of operators (from right to left)  $\mathbf{D}_h\mathbf{K}_h$ . Similarly, starting at  $t = +\epsilon$  yields the map

$$\mathbf{q}' = \mathbf{q} + h\mathbf{p} \quad ; \quad \mathbf{p}' = \mathbf{p} - h\nabla\Phi(\mathbf{q}'), \quad (9)$$

called the **drift-kick** or (also) the modified Euler method, which is represented as the operator  $\mathbf{K}_h \mathbf{D}_h$ . Since both integrators are derived from the Hamiltonian (4) they are both symplectic.

There is a slightly simpler version of this argument. Hamiltonian flows conserve the volume element in phase space (Liouville’s theorem) so another desirable geometric feature of an integrator for Hamiltonian systems is volume conservation. Examination of equation (6) shows that the Jacobian  $\partial(\bar{\mathbf{q}}, \bar{\mathbf{p}})/\partial(\mathbf{q}, \mathbf{p}) = 1$  so the kick operator  $\mathbf{K}_h$  conserves phase-space volume. Similarly, equation (7) shows that  $\partial(\mathbf{q}', \mathbf{p}')/\partial(\bar{\mathbf{q}}, \bar{\mathbf{p}}) = 1$  so the drift operator  $\mathbf{D}_h$  also conserves phase-space volume. Thus any composition of these operators conserves phase-space volume. Volume conservation is a less powerful constraint than symplecticity, since symplectic transformations conserve phase-space volume but not all volume-conserving flows are symplectic.

According to the kick-drift modified Euler integrator, the position after timestep  $h$  is

$$\mathbf{q}' = \mathbf{q} + h\mathbf{p}' = \mathbf{q} + h\mathbf{p} - h^2 \nabla \Phi(\mathbf{q}), \quad (10)$$

while the exact result may be written as a Taylor series,

$$\mathbf{q}' = \mathbf{q} + h\dot{\mathbf{q}}(t=0) + \frac{1}{2}h^2\ddot{\mathbf{q}}(t=0) + O(h^3) = \mathbf{q} + h\mathbf{p} - \frac{1}{2}h^2 \nabla \Phi(\mathbf{q}) + O(h^3). \quad (11)$$

The error after a single step of the modified Euler integrator is seen to be  $O(h^2)$ , so it is a first-order integrator.

*Exercise:* What are the drift and kick operators for cosmological comoving coordinates?

**Leapfrog integrator** By alternating kick and drift steps in more elaborate sequences, we can construct higher-order integrators; these are automatically symplectic since they are the composition of maps (the kick and drift steps) that are themselves symplectic. The simplest and most widely used of these is the **leapfrog** or **Verlet** integrator in which we drift for  $\frac{1}{2}h$ , kick for  $h$  and then drift for  $\frac{1}{2}h$ , namely:

$$\mathbf{q}_{1/2} = \mathbf{q} + \frac{1}{2}h\mathbf{p} ; \mathbf{p}' = \mathbf{p} - h\nabla\Phi(\mathbf{q}_{1/2}) ; \mathbf{q}' = \mathbf{q}_{1/2} + \frac{1}{2}h\mathbf{p}'. \quad (12)$$

This algorithm is sometimes called “drift-kick-drift” leapfrog; an equally good form is “kick-drift-kick” leapfrog:

$$\mathbf{p}_{1/2} = \mathbf{p} - \frac{1}{2}h\nabla\Phi(\mathbf{q}) ; \mathbf{q}' = \mathbf{q} + h\mathbf{p}_{1/2} ; \mathbf{p}' = \mathbf{p}_{1/2} - \frac{1}{2}h\nabla\Phi(\mathbf{q}'). \quad (13)$$

Drift-kick-drift leapfrog can also be derived by considering motion in the Hamiltonian (4) from  $t = -\frac{1}{2}h$  to  $t = \frac{1}{2}h$ . These can be written in operator form as  $\mathbf{D}_{h/2}\mathbf{K}_h\mathbf{D}_{h/2}$  and  $\mathbf{K}_{h/2}\mathbf{D}_h\mathbf{K}_{h/2}$  respectively.

The leapfrog integrator has many appealing features: (i) In contrast to the modified Euler integrator, it is second- rather than first-order accurate, in that the error in phase-space position after a single timestep is  $O(h^3)$ . (ii) Leapfrog is **time reversible** in the sense that if leapfrog advances the system from  $(\mathbf{q}, \mathbf{p})$  to  $(\mathbf{q}', \mathbf{p}')$  in a given time, it will also advance it from  $(\mathbf{q}', -\mathbf{p}')$  to  $(\mathbf{q}, -\mathbf{p})$  in the same time. Time-reversibility is a constraint on the phase-space flow that, like symplecticity, suppresses numerical dissipation, since dissipation is not time-reversible. (iii) A sequence of  $n$  leapfrog steps can be regarded as a drift step for  $\frac{1}{2}h$ , then  $N$  kick-drift steps of the modified Euler integrator, then a drift step for  $-\frac{1}{2}h$ ; thus if  $N \gg 1$  the leapfrog integrator requires negligibly more work than the same number of steps of the modified Euler integrator. In operator form this sequence can be written  $\mathbf{D}_{-h/2}(\mathbf{D}_h\mathbf{K}_h)^N\mathbf{D}_{h/2}$ . (iv) Leapfrog requires no additional memory for storing intermediate results. This is usually not a concern for planetary integrations but is a huge advantage for large N-body simulations.

The performance of Euler’s method, the modified Euler method, leapfrog, and a fourth-order Runge–Kutta method are compared in Figure 1. Even though modified Euler and leapfrog are only first- and second-order schemes they outperform Runge–Kutta—in the sense that the energy error is smaller—at large times, because Runge–Kutta exhibits numerical dissipation (linear growth in the energy error), while the symplectic schemes do not.

**The Forest fourth-order integrator** Over timestep  $h$  this integrator has the form

$$\mathbf{D}_{ah}\mathbf{K}_{ah}\mathbf{D}_{bh}\mathbf{K}_{ch}\mathbf{D}_{bh}\mathbf{K}_{ah}\mathbf{D}_{ah}; \tag{14}$$

the symmetry of this formula shows that the operator is time-reversible, and with the particular choice  $a = 1.35120$ ,  $b = -0.3512$ ,  $c = -1.7024$  the operator becomes fourth-order rather than second-order (all time-reversible operators must have even order). There are also higher-order generalizations (Yoshida 1993, *Cel. Mech.* 56, 27).

*Exercise:* Let  $\mathbf{P}_h$  be an integration algorithm of order  $2k$  for timestep  $h$ . Prove that  $\mathbf{P}_{ah}\mathbf{P}_{bh}\mathbf{P}_{ah}$  is of order  $2k + 2$  if  $a^{-1} = 2 - 2^{1/(2k+1)}$  and  $b = 1 - 2a$ .

**Variable timesteps** One serious limitation of symplectic integrators is that they work well only with fixed timesteps, as the following argument shows. Suppose the timestep depends

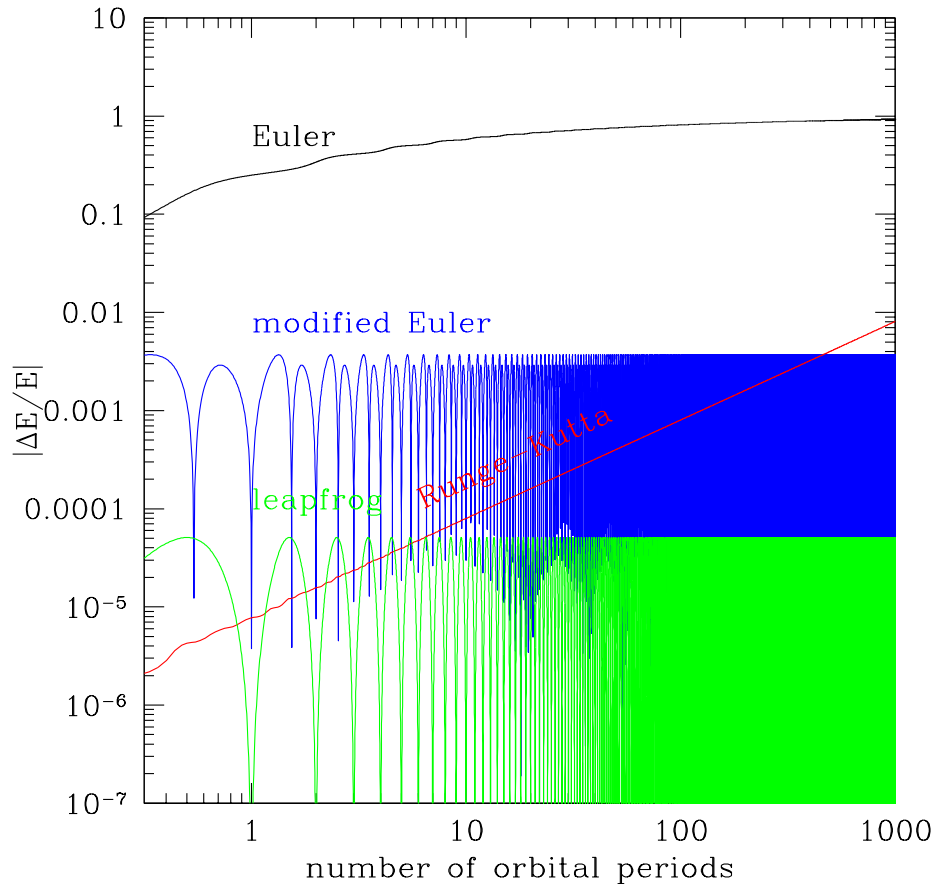


Fig. 1.— Fractional energy errors in the integration of a Kepler orbit with eccentricity  $e = 0.1$ . Each integrator is allowed 200 force evaluations per orbit. Euler’s method behaves very poorly, showing energy errors of order unity after only 100 orbits. In contrast, the modified Euler method shows an oscillatory energy error with no long-term growth; leapfrog behaves similarly although with maximum energy errors that are smaller by a factor of 30 or so (as expected since it is a second-order method, while modified Euler is first-order). Runge–Kutta, a fourth-order method, behaves much better than the others over short times but exhibits a linear drift in energy that eventually makes its performance worse than modified Euler or leapfrog.

on phase-space position,  $h(\mathbf{q}, \mathbf{p})$ . The Hamiltonian (4) becomes

$$H_h(\mathbf{q}, \mathbf{p}, t) = \frac{1}{2}\mathbf{p}^2 + \Phi(\mathbf{q})\delta_{h(\mathbf{q}, \mathbf{p})}(t), \quad \text{where} \quad \delta_h(t) \equiv h \sum_{j=-\infty}^{\infty} \delta(t - jh) \quad (15)$$

Since Hamilton's equations (1) require derivatives of  $H_h$  they now involve derivatives of delta functions, which means that the equations of motion are not the usual kick and drift operators with position-dependent timestep, nor are there *any* simple analytic operators corresponding to this Hamiltonian. In words, a symplectic integrator with fixed timestep is generally no longer symplectic once the timestep is varied<sup>1</sup>.

Fortunately, the geometric constraints on the phase-space flow imposed by time-reversibility are also strong, so the leapfrog integrator retains its good behavior if the timestep is adjusted in a time-reversible manner, even though the resulting integrator is no longer symplectic. Here is one way to do this: suppose that the appropriate timestep  $h$  is given by some function  $\tau(\mathbf{q}, \mathbf{p})$  of the phase-space coordinates. Then we modify equations (12) to

$$\begin{aligned} \mathbf{q}_{1/2} &= \mathbf{q} + \frac{1}{2}h\mathbf{p} \quad ; \quad \mathbf{p}_{1/2} = \mathbf{p} - \frac{1}{2}h\nabla\Phi(\mathbf{q}_{1/2}), \\ t' &= t + \frac{1}{2}(h + h'), \\ \mathbf{p}' &= \mathbf{p}_{1/2} - \frac{1}{2}h'\nabla\Phi(\mathbf{q}_{1/2}) \quad ; \quad \mathbf{q}' = \mathbf{q}_{1/2} + \frac{1}{2}h'\mathbf{p}'. \end{aligned} \quad (16)$$

Here  $h'$  is determined from  $h$  by solving the equation  $u(h, h') = \tau(\mathbf{q}_{1/2}, \mathbf{p}_{1/2})$ , where  $u(h, h')$  is any symmetric function of  $h$  and  $h'$  such that  $u(h, h) = h$ ; for example,  $u(h, h') = \frac{1}{2}(h + h')$  or  $u(h, h') = 2hh'/(h + h')$ .

Another more general way to do this is a clever algorithm by Mikkola & Merritt (2006, MNRAS 372, 219) which converts *any* integrator into a time-reversible one. Let  $\mathbf{y}(t)$  be the trajectory we are trying to integrate and let  $\mathbf{G}_h(\mathbf{y})$  be any integrator. Then set  $\mathbf{z} = \mathbf{y}$  and write

$$\begin{aligned} \bar{\mathbf{y}} &= \mathbf{y} + \mathbf{G}_{h/2}(\mathbf{z}) \quad ; \quad \bar{\mathbf{z}} = \mathbf{z} - \mathbf{G}_{-h/2}(\bar{\mathbf{y}}), \\ u(h, h') &= \tau(\bar{\mathbf{y}}), \\ \mathbf{z}' &= \bar{\mathbf{z}} + \mathbf{G}_{h/2}(\bar{\mathbf{y}}) \quad ; \quad \mathbf{y}' = \bar{\mathbf{y}} - \mathbf{G}_{-h/2}(\mathbf{z}'), \end{aligned} \quad (17)$$

which is time-reversible, that is, if  $(\mathbf{y}, \mathbf{z}) \rightarrow (\mathbf{y}', \mathbf{z}')$  in a timestep  $h$ , then  $(\mathbf{y}', \mathbf{z}') \rightarrow (\mathbf{y}, \mathbf{z})$  in a timestep  $-h$ .

---

<sup>1</sup>A symplectic integrator *does* remain symplectic if the timesteps are varied in some fixed pattern that does not depend on the phase-space coordinates.

## 1.2. Mixed-variable symplectic integrators

There is a more general way to describe the derivation of the modified Euler and leapfrog integrators. Suppose that a Hamiltonian  $H$  can be written as the sum of two terms,

$$H(\mathbf{q}, \mathbf{p}) = H_A(\mathbf{q}, \mathbf{p}) + H_B(\mathbf{q}, \mathbf{p}) \quad (18)$$

where  $H_A$  and  $H_B$  are separately integrable, that is, the phase-space position  $(\mathbf{q}, \mathbf{p})$  evolves under the action of  $H_A$  over a time  $h$  to

$$(\mathbf{q}', \mathbf{p}') = \mathbf{A}_h(\mathbf{q}, \mathbf{p}) \quad (19)$$

where  $\mathbf{A}_h$  is a (usually) nonlinear operator on the argument  $(\mathbf{q}, \mathbf{p})$ , with a similar operator  $\mathbf{B}_h$  for  $H_B$ . In the example we gave earlier,  $H_A = \frac{1}{2}\mathbf{p}^2$  and  $H_B = \Phi(\mathbf{q})$  so

$$\mathbf{A}_h(\mathbf{q}, \mathbf{p}) = (\mathbf{q} + h\mathbf{p}, \mathbf{p}), \quad \mathbf{B}_h(\mathbf{q}, \mathbf{p}) = (\mathbf{q}, \mathbf{p} - h\nabla\Phi(\mathbf{q})). \quad (20)$$

These will be recognized as the drift and kick operators so the drift-kick and kick-drift modified Euler integrators are (operators are applied from right to left)

$$\mathbf{B}_h\mathbf{A}_h, \quad \mathbf{A}_h\mathbf{B}_h, \quad (21)$$

and the drift-kick-drift and kick-drift-kick leapfrog integrators are

$$\mathbf{A}_{h/2}\mathbf{B}_h\mathbf{A}_{h/2}, \quad \mathbf{B}_{h/2}\mathbf{A}_h\mathbf{B}_{h/2}. \quad (22)$$

Now specialize to a Hamiltonian of the form  $\frac{1}{2}\mathbf{p}^2 - \mathfrak{G}M/|\mathbf{q}| + \epsilon\phi(\mathbf{q})$ , i.e., the Kepler problem with a small additional perturbing potential,  $\epsilon \ll 1$ . For leapfrog we would set  $H_A = \frac{1}{2}\mathbf{p}^2$ ,  $H_B = -\mathfrak{G}M/|\mathbf{q}| + \epsilon\phi(\mathbf{q}, t)$ . So far there is nothing new. However, this is not the only split we could make. The only requirement is that  $H_A$  and  $H_B$  be integrable, so we could equally well split as  $H_A = \frac{1}{2}\mathbf{p}^2 - \mathfrak{G}M/|\mathbf{q}|$  (the Kepler problem, which is integrable) and  $H_B = \epsilon\phi(\mathbf{q}, t)$ . Then the operator  $\mathbf{A}_h$  corresponds to advancing the particle on a Kepler orbit for a time  $h$ —most easily done using the Gauss  $f$  and  $g$  functions—and  $\mathbf{B}_h$  is the usual kick operator due to the potential  $\epsilon\phi(\mathbf{q})$ . The advantage of this approach is that the integration errors after one timestep are of order  $\epsilon h^3$  rather than  $h^3$  as in standard leapfrog (Wisdom & Holman 1991, AJ 102 1528; Kinoshita et al. 1991, Cel. Mech. 50, 59). The term “mixed variable” arises because the integrator consists of alternating steps which are trivial in Cartesian coordinates and orbital elements: thus all of the work is done in converting back and forth from one set of canonical variables to the other.

There are many refinements to this basic scheme. (i) If the perturbing potential  $\phi(\mathbf{q})$  is due to other planets whose orbits must also be followed, the integration is best done in



Jacobi coordinates rather than center-of-mass or heliocentric coordinates. (ii) The dominant contributions from general relativity are straightforward to include. (iii) Individual timesteps for each planet can speed up the calculation. (iv) The errors in long-term integrations can be reduced by an additional factor of  $\epsilon$  by methods that add negligible computational cost, including symplectic correctors and “warmup”; the latter term refers to slowly increasing the timestep from nearly zero to its final value over the first few thousand orbits (Wisdom & Holman 1991; Saha & Tremaine 1992, AJ 104, 1633; Saha & Tremaine 1994, AJ 108, 1962; Wisdom, Holman, & Touma 1996, Fields Institute Communications 10, 217).

Mixed-variable symplectic methods have become the workhorse for long planetary system integrations, at least if the planets are on approximately circular and coplanar orbits. Public software packages that implement this method include MERCURY (<http://star.arm.ac.uk/~jec/home.html>) and SWIFTER (<http://www.boulder.swri.edu/swifter/>).

### 1.3. Regularization

Highly eccentric orbits are difficult for most integrators to handle, because the acceleration is very high for a small fraction of the orbit. This problem can be circumvented by transforming to a coordinate system in which the two-body problem has no singularity—this procedure is called **regularization**. Standard integrators, symplectic or not, can then be used to solve the equations of motion in the regularized coordinates.

**Burdet–Heggie regularization** The simplest approach to regularization is time transformation. We write the equations of motion for the two-body problem as

$$\ddot{\mathbf{r}} = -\mathfrak{G} M \frac{\mathbf{r}}{r^3} + \mathbf{g}, \quad (23)$$

where  $\mathbf{g}$  is the gravitational field from the other  $N - 2$  bodies in the simulation (or other sources), and change to a fictitious time  $\tau$  that is defined by

$$dt = r d\tau. \quad (24)$$

Denoting derivatives with respect to  $\tau$  by a prime we find

$$\dot{\mathbf{r}} = \frac{d\tau}{dt} \frac{d\mathbf{r}}{d\tau} = \frac{1}{r} \mathbf{r}' \quad ; \quad \ddot{\mathbf{r}} = \frac{d\tau}{dt} \frac{d}{d\tau} \frac{1}{r} \mathbf{r}' = \frac{1}{r^2} \mathbf{r}'' - \frac{r'}{r^3} \mathbf{r}'. \quad (25)$$

Substituting these results into the equation of motion, we obtain

$$\mathbf{r}'' = \frac{r'}{r} \mathbf{r}' - \mathfrak{G} M \frac{\mathbf{r}}{r} + r^2 \mathbf{g}. \quad (26)$$

The eccentricity vector  $\mathbf{e}$  helps us to simplify this equation. We have

$$\begin{aligned}\mathbf{e} &= \frac{\mathbf{v} \times (\mathbf{r} \times \mathbf{v})}{\mathfrak{G} M} - \hat{\mathbf{r}} \\ &= \frac{1}{\mathfrak{G} M} \left( |\mathbf{r}'|^2 \frac{\mathbf{r}}{r^2} - \frac{r'}{r} \mathbf{r}' \right) - \frac{\mathbf{r}}{r},\end{aligned}\quad (27)$$

where we have used  $\mathbf{v} = \dot{\mathbf{r}} = \mathbf{r}'/r$ . Thus equation (26) can be written

$$\mathbf{r}'' = |\mathbf{r}'|^2 \frac{\mathbf{r}}{r^2} - 2\mathfrak{G} M \frac{\mathbf{r}}{r} - \mathfrak{G} M \mathbf{e} + r^2 \mathbf{g}.\quad (28)$$

The energy of the two-body orbit is

$$E_2 = \frac{1}{2}v^2 - \frac{\mathfrak{G} M}{r} = \frac{|\mathbf{r}'|^2}{2r^2} - \frac{\mathfrak{G} M}{r},\quad (29)$$

so we arrive at the regularized equation of motion

$$\mathbf{r}'' - 2E_2 \mathbf{r} = -\mathfrak{G} M \mathbf{e} + r^2 \mathbf{g},\quad (30)$$

in which the singularity at the origin has disappeared. This must be supplemented by equations for the rates of change of  $E_2$ ,  $\mathbf{e}$ , and  $t$  with fictitious time  $\tau$ ,

$$E_2' = \mathbf{g} \cdot \mathbf{r}' \quad ; \quad \mathbf{e}' = 2\mathbf{r}(\mathbf{r}' \cdot \mathbf{g}) - \mathbf{r}'(\mathbf{r} \cdot \mathbf{g}) - \mathbf{g}(\mathbf{r} \cdot \mathbf{r}') \quad ; \quad t' = r.\quad (31)$$

When the external field  $\mathbf{g}$  vanishes, the energy  $E_2$  and eccentricity vector  $\mathbf{e}$  are constants, and the equation of motion (30) is that of a harmonic oscillator that is subject to a constant force  $-\mathbf{e}$ . The singularity at  $V - R = 0$  has disappeared.

*Exercise:* In the absence of external forces, prove that the fictitious time  $\tau$  is proportional to the eccentric anomaly.

**Kustaanheimo–Stiefel (KS) regularization** An alternative regularization procedure, which involves the transformation of the coordinates in addition to time, can be derived using the symmetry group of the Kepler problem, the theory of quaternions and spinors, or several other methods (Stiefel & Schiefele, *Linear and Regular Celestial Mechanics*; Heggie & Hut, *The Gravitational Million-Body Problem*; Saha 2009, MNRAS 400, 228). Once again we use the fictitious time  $\tau$  defined by equation (24). We also define a four-vector  $\mathbf{u} = (u_1, u_2, u_3, u_4)$

that is related to the position  $\mathbf{r} = (x, y, z)$  by

$$\begin{aligned} u_1^2 &= \frac{1}{2}(x+r)\cos^2\psi \\ u_4^2 &= \frac{1}{2}(x+r)\sin^2\psi \\ u_2 &= \frac{yu_1 + zu_4}{x+r} \\ u_3 &= \frac{zu_1 - yu_4}{x+r}, \end{aligned} \tag{32}$$

where  $\psi$  is an arbitrary parameter. The inverse relations are

$$x = u_1^2 - u_2^2 - u_3^2 + u_4^2; \quad y = 2(u_1u_2 - u_3u_4); \quad z = 2(u_1u_3 + u_2u_4). \tag{33}$$

Note that  $r = u_1^2 + u_2^2 + u_3^2 + u_4^2$ . Let  $\Phi_{\text{ext}}$  be the potential that generates the external field  $\mathbf{g} = -\nabla\Phi_{\text{ext}}$ . Then in terms of the new variables the equation of motion (23) reads

$$\begin{aligned} \mathbf{u}'' - \frac{1}{2}E\mathbf{u} &= -\frac{1}{4}\frac{\partial}{\partial\mathbf{u}}(|\mathbf{u}|^2\Phi_{\text{ext}}), \\ E = \frac{1}{2}v^2 - \frac{\mathfrak{G}M}{r} + \Phi_{\text{ext}} &= 2\frac{|\mathbf{u}'|^2}{|\mathbf{u}|^2} - \frac{\mathfrak{G}M}{|\mathbf{u}|^2} + \Phi_{\text{ext}}, \\ E' = |\mathbf{u}|^2\frac{\partial\Phi_{\text{ext}}}{\partial t} \quad ; \quad t' = |\mathbf{u}|^2, \end{aligned} \tag{34}$$

When the external force vanishes, the first of these is the equation of motion for a four-dimensional harmonic oscillator.

Figure 2 shows the fractional energy error that arises in the integration of an orbit with eccentricity  $e = 0.99$  using KS regularization. Using the same integrator, the energy error is more than an order of magnitude smaller than the error using Burdet–Heggie regularization.

## 2. The gravitational N-body problem

Now consider the case when  $N$  planets are present, with masses  $m_j$  and positions  $\mathbf{r}_j$ ,  $j = 1, \dots, N$ . In addition there is a host star, with mass and position  $m_0, \mathbf{r}_0$ . The equations of motion for the two-body problem are replaced by

$$\frac{d^2\mathbf{r}_j}{dt^2} = \sum_{k=0, k \neq j}^N \frac{\mathfrak{G}m_k(\mathbf{r}_k - \mathbf{r}_j)}{|\mathbf{r}_k - \mathbf{r}_j|^3}, \quad j = 0, \dots, N. \tag{35}$$

Here  $\mathbf{r}_j$  is the position in an inertial frame. In some cases it is more useful to work in a frame centered on the host star. Let  $\mathbf{x}_j = \mathbf{r}_j - \mathbf{r}_0$  be the position of planet  $j$  relative to the

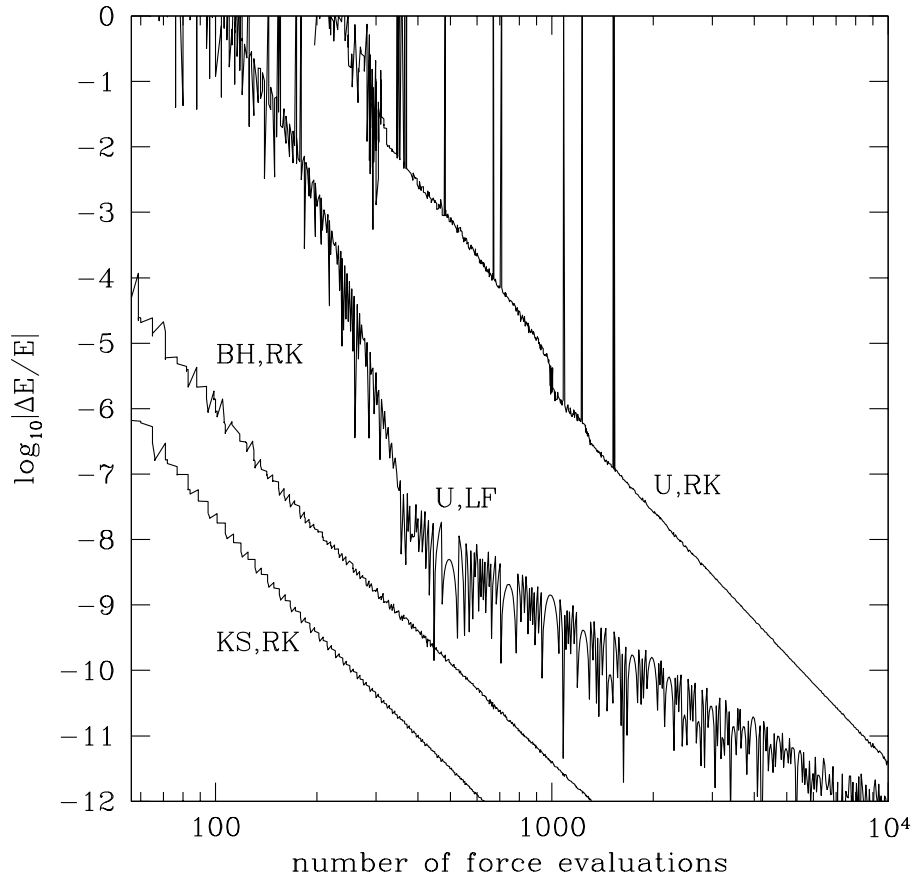


Fig. 2.— Fractional energy errors in the integration of a Kepler orbit with eccentricity  $e = 0.99$  for one orbital period. The  $x$ -axis is the number of force evaluations. Curves labeled by “RK” are followed using a fourth-order Runge–Kutta integrator with adaptive timestep control. The curve labeled “U” for “unregularized” is integrated in Cartesian coordinates, the curve “BH” uses Burdet–Heggie regularization, and the curve “KS” uses Kustaanheimo–Stiefel regularization. The curve labeled “U,LF” is followed in Cartesian coordinates using a leapfrog integrator with timestep proportional to radius (eq. 16). From Binney & Tremaine, *Galactic Dynamics*, 2nd ed.

host star. Then

$$\frac{d^2 \mathbf{x}_j}{dt^2} = -\frac{\mathfrak{G}(m_0 + m_j)}{|\mathbf{x}_j|^3} \mathbf{x}_j + \sum_{k=1, k \neq j}^N \mathfrak{G} m_k \left[ \frac{(\mathbf{x}_k - \mathbf{x}_j)}{|\mathbf{x}_k - \mathbf{x}_j|^3} - \frac{\mathbf{x}_k}{|\mathbf{x}_k|^3} \right], \quad j = 1, \dots, N. \quad (36)$$

We can write  $d^2 \mathbf{x}_j / dt^2 = -\nabla \Phi_j$  where

$$\Phi_j = -\frac{\mathfrak{G}(m_0 + m_j)}{|\mathbf{x}_j|} - \sum_{k=1, k \neq j}^N \frac{\mathfrak{G} m_k}{|\mathbf{x}_k - \mathbf{x}_j|} + \sum_{k=1, k \neq j}^N \frac{\mathfrak{G} m_k \mathbf{x}_j \cdot \mathbf{x}_k}{|\mathbf{x}_k|^3}. \quad (37)$$

The last term is called the **indirect potential**. For most calculations the frame centered on the host star is more convenient, but there are some cases in which the inertial frame is far better (e.g., following the motion of a comet or other body far outside the planetary system).

There are other useful coordinate systems, in particular **Jacobi coordinates**, which are needed for mixed-variable symplectic integrations of the N-body equations of motion.

Of course, the equations we have derived also govern the motion of stars in globular clusters and at the centers of galaxies, the topics of other lectures in this course. Celestial mechanics is distinguished from these subjects by its focus on gravitational N-body systems in which (i) there is a dominant central mass; (ii) the eccentricities and inclinations are small; (iii) the system must be followed for large numbers of orbits, typically up to  $10^{10}$ – $10^{11}$ .

## 2.1. The stability of the solar system

This is one of the oldest problems in theoretical physics. From Newton’s laws of motion and his law of universal gravitation, we know that each planet causes small oscillatory variations in the orbits of the other planets. Although the fractional variations in the orbital elements are small (typically less than  $10^{-3}$ – $10^{-4}$ ) the age of the system is very large ( $10^8$ – $10^{10}$  orbital periods). Over these vary large times, do the variations in the orbits remain oscillatory or do they gradually grow, leading eventually to ejection of a planet or collisions between two planets or a planet and the Sun?

Many famous mathematicians and physicists have thought about this problem: Newton, Leibnitz, Laplace, Lagrange, Gauss, Poincaré, Kolmogorov, Arnold, etc. These investigations have yielded remarkable insights into the behavior of planetary systems and into nonlinear dynamics in general. However, the analytic results are all restricted to planetary systems with unrealistically small masses, eccentricities, and inclinations.

To investigate the stability of the actual solar system we must use numerical integrations. It is only in the last decade or so that we have had the ability to follow the motions of the eight planets accurately for timescales comparable to the age of the solar system (4.5 Gyr) or the time remaining until the inner planets are swallowed by the Sun (7.7 Gyr). See for example Ito & Tanikawa (2002), MNRAS 336, 483; Laskar & Gastineau (2009), Nature 459, 817.

These integrations show that in most cases the solar system is stable over timescales of a few Gyr. All of the planets survive, and most remain in orbits very similar to their present ones. However, all the planetary orbits are chaotic, with a Liapunov time  $t_L \sim 10^7$  years. In more detail: (i) on timescales  $\lesssim t_L$  the planetary orbits are regular, that is, small changes in position and velocity grow only linearly with time, so the planetary orbits are highly predictable—this of course is why we can accurately predict eclipses, send spacecraft to other planets, etc. (ii) On timescales  $\gg t_L$  small changes in the orbits grow exponentially, as  $\exp(t/t_L)$ . This means that small changes now in the Earth’s orbit—for example from the difference in the gravitational attraction by Jupiter on your coffee cup when you lift it to drink—will be amplified by a factor of  $\sim \exp(7.7 \text{ Gyr}/t_L) \sim 10^{330}$  before the death of the Sun. (iii) Most of this chaotic behavior is restricted to narrow regions in phase space associated with high-order resonances, and this means in turn that the chaotic behavior is mostly in the planetary phase rather than other orbital elements such as semi-major axis or eccentricity. (iv) Most, but not all: the eccentricity of Mercury’s orbit undergoes a random walk and there is about a 1% probability that it will experience some catastrophic event—a collision with the Sun, Venus, or even Earth—before the death of the Sun.

Of course, long after such an event there would be no obvious sign that Mercury had ever been present in the solar system. Thus it is a very plausible speculation that the solar system had more planets early in its history, and that one or more of these has been lost. We conclude that the answer to the question “is the solar system stable?” is not a simple “yes” or “no”: in the future the solar system is probably stable (at about the 99% confidence level) up to the time when the Sun dies in about 8 Gyr; on the other hand, in the past the solar system probably was *not* stable but we will probably never know for certain.