

# Geometric numerical integration illustrated by the Störmer–Verlet method

Ernst Hairer

*Section de Mathématiques,  
Université de Genève, Switzerland  
E-mail: Ernst.Hairer@math.unige.ch*

Christian Lubich

*Mathematisches Institut,  
Universität Tübingen, Germany  
E-mail: Lubich@na.uni-tuebingen.de*

Gerhard Wanner

*Section de Mathématiques,  
Université de Genève, Switzerland  
E-mail: Gerhard.Wanner@math.unige.ch*

The subject of geometric numerical integration deals with numerical integrators that preserve geometric properties of the flow of a differential equation, and it explains how structure preservation leads to improved long-time behaviour. This article illustrates concepts and results of geometric numerical integration on the important example of the Störmer–Verlet method. It thus presents a cross-section of the recent monograph by the authors, enriched by some additional material.

After an introduction to the Newton–Störmer–Verlet–leapfrog method and its various interpretations, there follows a discussion of geometric properties: reversibility, symplecticity, volume preservation, and conservation of first integrals. The extension to Hamiltonian systems on manifolds is also described. The theoretical foundation relies on a backward error analysis, which translates the geometric properties of the method into the structure of a modified differential equation, whose flow is nearly identical to the numerical method. Combined with results from perturbation theory, this explains the excellent long-time behaviour of the method: long-time energy conservation, linear error growth and preservation of invariant tori in near-integrable systems, a discrete virial theorem, and preservation of adiabatic invariants.

## CONTENTS

1	The Newton–Störmer–Verlet–leapfrog method	400
2	Geometric properties	408
3	Conservation of first integrals	416
4	Backward error analysis	419
5	Long-time behaviour of numerical solutions	425
6	Constrained Hamiltonian systems	440
7	Geometric integration beyond Störmer–Verlet	447
	References	447

### 1. The Newton–Störmer–Verlet–leapfrog method

We start by considering systems of second-order differential equations

$$\ddot{q} = f(q), \quad (1.1)$$

where the right-hand side  $f(q)$  does not depend on  $\dot{q}$ . Many problems in astronomy, molecular dynamics, and other areas of physics are of this form.

#### 1.1. Two-step formulation

If we choose a step size  $h$  and grid points  $t_n = t_0 + nh$ , the most natural discretization of (1.1) is

$$q_{n+1} - 2q_n + q_{n-1} = h^2 f(q_n), \quad (1.2)$$

which determines  $q_{n+1}$  whenever  $q_{n-1}$  and  $q_n$  are known. Geometrically, this amounts to determining an interpolating parabola which, in the mid-point, assumes the second derivative prescribed by equation (1.1); see Figure 1.1, left.

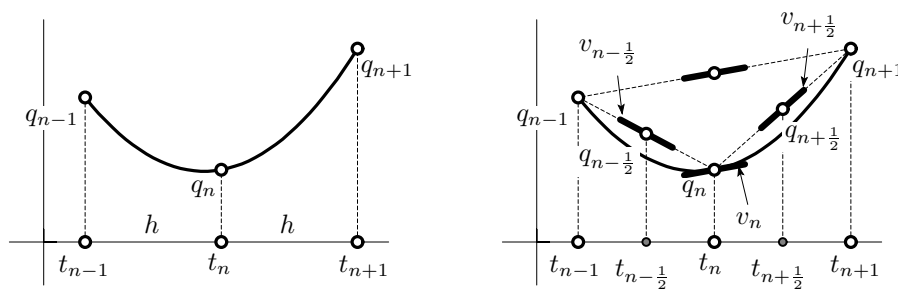


Figure 1.1. Method (1.2): two-step formulation (left); one-step formulations (right).

1.2. One-step formulations

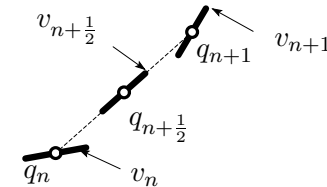
Introducing the velocity  $\dot{q} = v$  turns equation (1.1) into a first-order system of doubled dimension

$$\dot{q} = v, \quad \dot{v} = f(q), \tag{1.3}$$

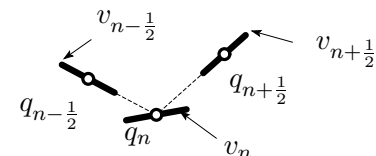
an equation in the so-called *phase space*. In analogy to this, we introduce discrete approximations of  $v$  and  $q$  as follows:

$$v_n = \frac{q_{n+1} - q_{n-1}}{2h}, \quad v_{n-\frac{1}{2}} = \frac{q_n - q_{n-1}}{h}, \quad q_{n-\frac{1}{2}} = \frac{q_n + q_{n-1}}{2}, \tag{1.4}$$

where some derivatives, in order to preserve second-order and symmetry, are evaluated on the *staggered grid*  $t_{n-\frac{1}{2}}, t_{n+\frac{1}{2}}, \dots$ ; see Figure 1.1, right. Inserting these expressions into the method (or simply looking at the picture) we see that method (1.2) can now be interpreted as a *one-step method*  $\Phi_h$ :  $(q_n, v_n) \mapsto (q_{n+1}, v_{n+1})$ , given by

$$\begin{aligned} (A) \quad & v_{n+\frac{1}{2}} = v_n + \frac{h}{2} f(q_n), \\ & q_{n+1} = q_n + h v_{n+\frac{1}{2}}, \\ & v_{n+1} = v_{n+\frac{1}{2}} + \frac{h}{2} f(q_{n+1}). \end{aligned} \tag{1.5}$$


There is a *dual* variant of the method on the staggered grid  $(v_{n-\frac{1}{2}}, q_{n-\frac{1}{2}}) \mapsto (v_{n+\frac{1}{2}}, q_{n+\frac{1}{2}})$  as follows:

$$\begin{aligned} (B) \quad & q_n = q_{n-\frac{1}{2}} + \frac{h}{2} v_{n-\frac{1}{2}}, \\ & v_{n+\frac{1}{2}} = v_{n-\frac{1}{2}} + h f(q_n), \\ & q_{n+\frac{1}{2}} = q_n + \frac{h}{2} v_{n+\frac{1}{2}}. \end{aligned} \tag{1.6}$$


For both arrays (A) and (B), we can concatenate, in the actual step-by-step procedure, the last line of the previous step with the first line of the subsequent step. Both schemes then turn into the same method, where the  $q$ -values are evaluated on the original grid, and the  $v$ -values are evaluated on the staggered grid:

$$\begin{aligned} v_{n+\frac{1}{2}} &= v_{n-\frac{1}{2}} + h f(q_n), \\ q_{n+1} &= q_n + h v_{n+\frac{1}{2}}. \end{aligned} \tag{1.7}$$

This is the computationally most economic implementation, and numerically more stable than (1.2); see Hairer, Nørsett and Wanner (1993, p. 472).

### 1.3. Historical remarks

Isn't that ingenious? I borrowed it straight from Newton. It comes right out of the *Principia*, diagram and all. R. Feynman (1965, p. 43)

The above schemes are known in the literature under various names. In particular, in molecular dynamics they are often called the *Verlet method* (Verlet 1967) and have become by far the most widely used integration scheme in this field.

Another name for this method is the *Störmer method*, since C. Störmer, in 1907, used higher-order variants of it for his computations of the motion of ionized particles in the earth's magnetic field (aurora borealis); see, *e.g.*, Hairer *et al.* (1993, Section III.10). Sometimes it is also called the *Encke method*, because J. F. Encke, around 1860, did extensive calculations for the perturbation terms of planetary orbits, which obey systems of second-order differential equations of precisely the form (1.1). Mainly in the context of partial differential equations of wave propagation, this method is called the *leapfrog method*. In yet another context, this formula is the basic method for the GBS extrapolation scheme, as it was proposed, for the case of equation (1.1), by Gragg in 1965; see Hairer *et al.* (1993, p. 294*f.*). Furthermore, the scheme (1.7) is equivalent to Nyström's method of order 2; see Hairer *et al.* (1993, p. 362, formula (III.1.13')).

A curious fact is that Professor Loup Verlet, who later became interested in the history of science, discovered precisely 'his' method in several places in the classical literature, for example, in the calculations of logarithms and astronomical tables by J. B. Delambre in 1792: this paper was translated and discussed in McLachlan and Quispel (2002, Appendix C). Even more spectacular is the finding that the 'Verlet method' was used in Newton's *Principia* from 1687 to prove Kepler's second law. An especially clear account can be found in Feynman's *Messenger Lecture* from 1964; see Feynman (1965, p. 41), from which we reproduce with pleasure<sup>1</sup> two of Feynman's original hand drawings.

The argument is as follows: if there are no forces, the body advances with uniform speed, and the radius vector covers equal areas in equal times, simply because the two triangles Sun-1-2 and Sun-2-3 have the same base and common altitudes (see Figure 1.2). If the gravitational force acts at the midpoint, the planet is deviated in such a way that the top of the second triangle moves parallel to the sun ray (see Figure 1.3). Hence, the triangle Sun-2-4 also has the same area. The whole procedure (uniform motion on half the interval, then a 'kick' to the velocity in the direction of the Sun, and another uniform motion on the second half) is precisely variant (B) of the Störmer–Verlet scheme.

<sup>1</sup> ... and with permission of the publisher

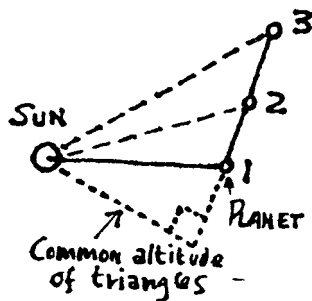


Figure 1.2. Uniform motion of a planet (drawing by R. Feynman).

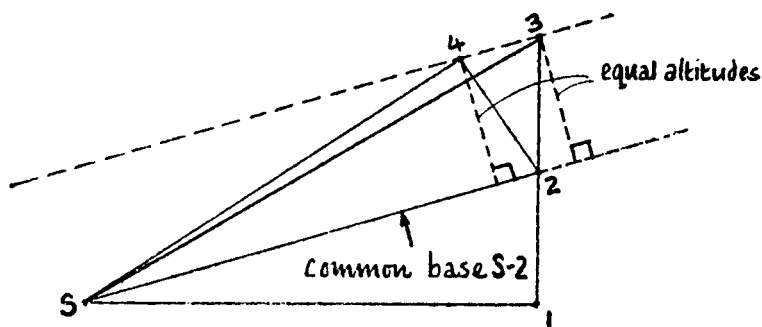


Figure 1.3. Gravitation acting at midpoint (drawing by R. Feynman).

1.4. Interpretation as composition method (symplectic Euler)

We can go a step further and split the formulae in the middle of the schemes (1.5) and (1.6). We then arrive at the schemes  $(v_n, q_n) \mapsto (v_{n+\frac{1}{2}}, q_{n+\frac{1}{2}})$  given by

$$(SE1) \quad \begin{aligned} v_{n+\frac{1}{2}} &= v_n + \frac{h}{2} f(q_n), \\ q_{n+\frac{1}{2}} &= q_n + \frac{h}{2} v_{n+\frac{1}{2}}, \end{aligned} \quad \begin{array}{c} \nearrow v_{n+\frac{1}{2}} \\ q_{n+\frac{1}{2}} \\ \nwarrow v_n \\ q_n \end{array} \quad (1.8)$$

as well as the *adjoint* scheme  $(v_{n+\frac{1}{2}}, q_{n+\frac{1}{2}}) \mapsto (v_{n+1}, q_{n+1})$  obtained by formally replacing the subscript  $n$  by  $n+1$  and  $h$  by  $-h$ ,

$$(SE2) \quad \begin{aligned} q_{n+1} &= q_{n+\frac{1}{2}} + \frac{h}{2} v_{n+\frac{1}{2}}, \\ v_{n+1} &= v_{n+\frac{1}{2}} + \frac{h}{2} f(q_{n+1}). \end{aligned} \quad \begin{array}{c} q_{n+1} \\ \nwarrow v_{n+1} \\ q_{n+\frac{1}{2}} \\ \nearrow v_{n+\frac{1}{2}} \\ q_n \\ \nwarrow v_n \\ q_{n-\frac{1}{2}} \\ \nearrow v_{n-\frac{1}{2}} \end{array} \quad (1.9)$$

Both these schemes, in which one variable is used at the old value and the other at the new value, are called the *symplectic Euler method*.

We thus see that the above scheme (A) is the composition of the symplectic Euler schemes (SE1) with (SE2), while the scheme (B) is the composition of method (SE2) followed by (SE1).

1.5. Interpretation as splitting method

We consider the vector field  $(v, f(q))$  of (1.3) ‘split’ as the sum of two vector fields  $(v, 0)$  and  $(0, f(q))$ , as indicated in Figure 1.4. The *exact* flows  $\varphi_t^{[1]}$  and  $\varphi_t^{[2]}$  of these two vector fields, which both have a constant time derivative, are easily obtained:

$$\varphi_t^{[1]} : \begin{cases} q_1 = q_0 + t \cdot v_0 \\ v_1 = v_0 \end{cases} \quad \text{and} \quad \varphi_t^{[2]} : \begin{cases} q_1 = q_0 \\ v_1 = v_0 + t \cdot f(q_0). \end{cases} \quad (1.10)$$

These formulae are precisely those which build up the formulae (SE1) and (SE2) above:

$$\begin{aligned} \text{(SE2)} &= \varphi_{h/2}^{[2]} \circ \varphi_{h/2}^{[1]}, & \text{(SE2)} & \text{diagram} & \varphi_{h/2}^{[1]} & \text{(SE1)} & \text{(1.11)} \\ \text{(SE1)} &= \varphi_{h/2}^{[1]} \circ \varphi_{h/2}^{[2]}. \end{aligned}$$

For the two versions of the Störmer–Verlet method we thus obtain the diagrams

$$\begin{aligned} \text{(A)} &= \text{(SE2)} \circ \text{(SE1)}, & \text{(A)} & \text{diagram} & \varphi_{h/2}^{[1]} & \text{(B)} & \text{(1.12)} \\ \text{(B)} &= \text{(SE1)} \circ \text{(SE2)}, \end{aligned}$$

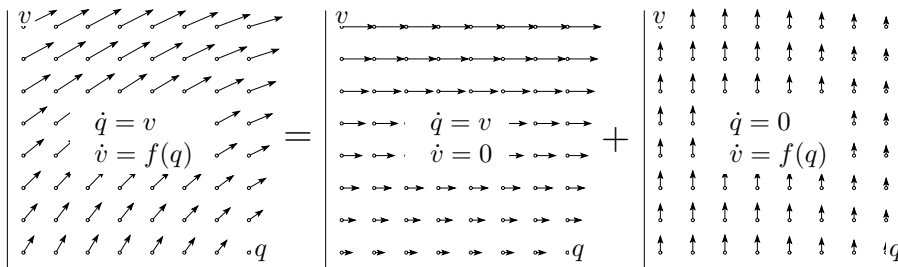


Figure 1.4. The phase space vector field split into two fields.

or, in explicit formulae,

$$\begin{aligned}\Phi_h^{(A)} &= \varphi_{h/2}^{[2]} \circ \varphi_h^{[1]} \circ \varphi_{h/2}^{[2]}, \\ \Phi_h^{(B)} &= \varphi_{h/2}^{[1]} \circ \varphi_h^{[2]} \circ \varphi_{h/2}^{[1]}.\end{aligned}\tag{1.13}$$

This way of composing the flows of split vector fields is often referred to as *Strang splitting*, after Strang (1968). For a careful survey of splitting methods we refer to McLachlan and Quispel (2002).

1.6. Interpretation as variational integrator

A further approach to the Störmer–Verlet method is obtained by discretizing Hamilton’s principle. This variational principle states that the motion of a mechanical system between any two positions  $q(t_0) = q_0$  and  $q(t_N) = q_N$  is such that the action integral

$$\int_{t_0}^{t_N} L(q(t), \dot{q}(t)) dt \quad \text{is minimized,}\tag{1.14}$$

where  $L(q, v)$  is the *Lagrangian* of the system. Typically, it is the difference between the kinetic and the potential energy, that is,

$$L(q, v) = \frac{1}{2}v^T Mv - U(q),\tag{1.15}$$

with a symmetric positive definite mass matrix  $M$ . When  $M$  does not depend on  $q$ , the Euler–Lagrange equations of this variational problem,  $\frac{d}{dt} \frac{\partial L}{\partial v} = \frac{\partial L}{\partial q}$ , reduce to the second-order differential equation  $M\ddot{q} = -\nabla U(q)$ .

We now approximate  $q(t)$  by a piecewise linear function, interpolating grid values  $(t_n, q_n)$  for  $n = 0, 1, \dots, N$ , and the action integral by the trapezoidal rule. We then require that  $q_1, \dots, q_{N-1}$  be such that, instead of (1.14),

$$\sum_{n=0}^{N-1} S_h(q_n, q_{n+1}) \quad \text{is minimized,}\tag{1.16}$$

where

$$S_h(q_n, q_{n+1}) = \frac{h}{2}L\left(q_n, \frac{q_{n+1} - q_n}{h}\right) + \frac{h}{2}L\left(q_{n+1}, \frac{q_{n+1} - q_n}{h}\right).\tag{1.17}$$

The requirement that the gradient with respect to  $q_n$  be zero, yields the discrete Euler–Lagrange equations

$$\nabla_Q S_h(q_{n-1}, q_n) + \nabla_q S_h(q_n, q_{n+1}) = 0$$

for  $n = 1, \dots, N - 1$ , where the partial gradients  $\nabla_q, \nabla_Q$  refer to  $S_h = S_h(q, Q)$ . In the case of the Lagrangian (1.15) these equations reduce to

$$M(q_{n+1} - 2q_n + q_{n-1}) + h^2 \nabla U(q_n) = 0,\tag{1.18}$$

which is just the two-step formulation (1.2) of the Störmer–Verlet method, with  $f(q) = -M^{-1}\nabla U(q)$ .

This variational interpretation of the Störmer–Verlet method was given by MacKay (1992). A comprehensive survey of variational integrators can be found in Marsden and West (2001).

### 1.7. Numerical example

We choose the Kepler problem

$$\ddot{q}_1 = -\frac{q_1}{(q_1^2 + q_2^2)^{3/2}}, \quad \ddot{q}_2 = -\frac{q_2}{(q_1^2 + q_2^2)^{3/2}}. \quad (1.19)$$

As initial values we take

$$q_1(0) = 1 - e, \quad q_2(0) = 0, \quad \dot{q}_1(0) = 0, \quad \dot{q}_2(0) = \sqrt{\frac{1+e}{1-e}}, \quad (1.20)$$

with  $e = 0.6$ . The period of the exact solution is  $2\pi$ . Figure 1.5 presents the numerical values of the Störmer–Verlet method for two different step sizes. These solutions are compared to those of the explicit midpoint rule in Runge’s one-step formulation; see Hairer, Lubich and Wanner (2002, p. 24, Figure 1.2. and equation (1.3)). This second method is of the same order and for the first steps it behaves very similarly to the Störmer–Verlet scheme (the first step is even identical!), but it deteriorates significantly as the integration interval increases. The explanation of this strange difference is the subject of the theories below.

### 1.8. Extension to general partitioned problems

For the extension of the above formulae to the more general system

$$\dot{q} = g(q, v), \quad \dot{v} = f(q, v), \quad (1.21)$$

we follow the ideas of De Vogelaere (1956). This is a marvellous paper, short, clear, elegant, written in one week, submitted for publication – and never published. We first extend the formulae (1.8) and (1.9), by taking over the missing arguments from one equation to the other. This gives

$$\begin{aligned} \text{(SE1)} \quad v_{n+\frac{1}{2}} &= v_n + \frac{h}{2} f(q_n, v_{n+\frac{1}{2}}), \\ q_{n+\frac{1}{2}} &= q_n + \frac{h}{2} g(q_n, v_{n+\frac{1}{2}}), \end{aligned} \quad (1.22)$$

and

$$\begin{aligned} \text{(SE2)} \quad q_{n+1} &= q_{n+\frac{1}{2}} + \frac{h}{2} g(q_{n+1}, v_{n+\frac{1}{2}}), \\ v_{n+1} &= v_{n+\frac{1}{2}} + \frac{h}{2} f(q_{n+1}, v_{n+\frac{1}{2}}). \end{aligned} \quad (1.23)$$



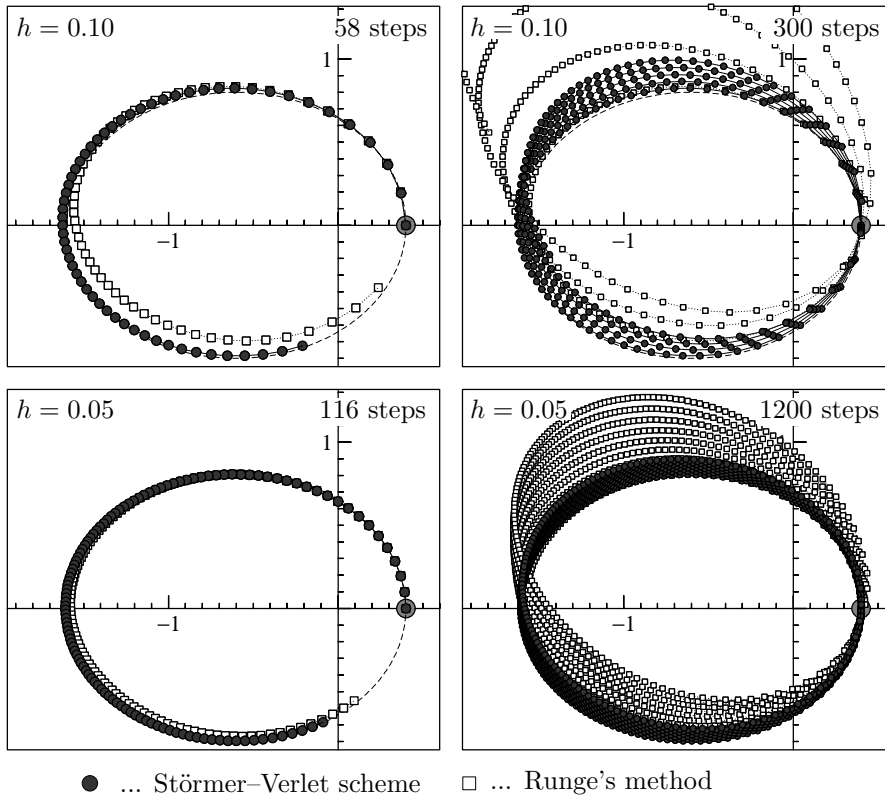


Figure 1.5. Kepler problem: the dashed line is the exact solution.

In each of these algorithms the derivative evaluations of both formulae are taken at the same point. The extensions of the Störmer-Verlet schemes are now obtained by composition, in the same way as in Section 1.4:

$$\begin{aligned}
 v_{n+\frac{1}{2}} &= v_n + \frac{h}{2} f(q_n, v_{n+\frac{1}{2}}), \\
 \text{(A) = (SE2) } \circ \text{ (SE1)} \quad q_{n+1} &= q_n + \frac{h}{2} \left( g(q_n, v_{n+\frac{1}{2}}) + g(q_{n+1}, v_{n+\frac{1}{2}}) \right), \quad (1.24) \\
 v_{n+1} &= v_{n+\frac{1}{2}} + \frac{h}{2} f(q_{n+1}, v_{n+\frac{1}{2}}),
 \end{aligned}$$

and, for the dual version,

$$\begin{aligned}
 q_n &= q_{n-\frac{1}{2}} + \frac{h}{2} g(q_n, v_{n-\frac{1}{2}}), \\
 \text{(B) = (SE1) } \circ \text{ (SE2)} \quad v_{n+\frac{1}{2}} &= v_{n-\frac{1}{2}} + \frac{h}{2} \left( f(q_n, v_{n-\frac{1}{2}}) + f(q_n, v_{n+\frac{1}{2}}) \right), \quad (1.25) \\
 q_{n+\frac{1}{2}} &= q_n + \frac{h}{2} g(q_n, v_{n+\frac{1}{2}}).
 \end{aligned}$$

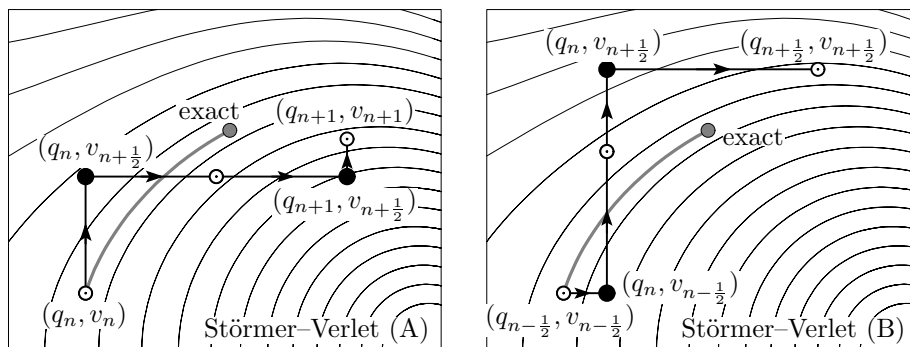


Figure 1.6. Störmer-Verlet methods for  $\dot{q} = v$ ,  $\dot{v} = -\sin q - v^2/5$ , initial values  $(-1.8, 0.3)$ , step size  $h = 1.5$ . Black points indicate where the vector field is evaluated.

For illustrations see Figure 1.6. The first equation of (1.24) is now an *implicit* formula for  $v_{n+\frac{1}{2}}$ , the second one for  $q_{n+1}$ , while only the last one is explicit. Such implicit methods were not common in the 1950s and might then not have delighted journal editors – or programmers:

No detailed example or discussion is given. This will best be done by those working on these problems in the Brookhaven, Harwell, MURA or CERN group.  
De Vogelaere (1956)

## 2. Geometric properties

We study geometric properties of the flow of differential equations which are preserved by the Störmer-Verlet method. The properties discussed are reversibility, symplecticity, and volume preservation.

### 2.1. Symmetry and reversibility

The Störmer-Verlet method is *symmetric* with respect to changing the direction of time: in its one-step formulation (1.5), replacing  $h$  by  $-h$  and exchanging the subscripts  $n \leftrightarrow n + 1$  (*i.e.*, reflecting time at the centre  $t_{n+1/2}$ ) gives the same method again. Similarly, the replacements  $h \leftrightarrow -h$  and  $n - \frac{1}{2} \leftrightarrow n + \frac{1}{2}$  leave the formulation (1.6) unchanged. In terms of the numerical one-step map  $\Phi_h : (q_n, v_n) \mapsto (q_{n+1}, v_{n+1})$ , this symmetry can be stated more formally as

$$\Phi_h = \Phi_{-h}^{-1}. \quad (2.1)$$

Such a relation does not hold for the symplectic Euler methods (1.8) and (1.9), where the above time-reflection transforms (SE1) to (SE2) and *vice versa*.

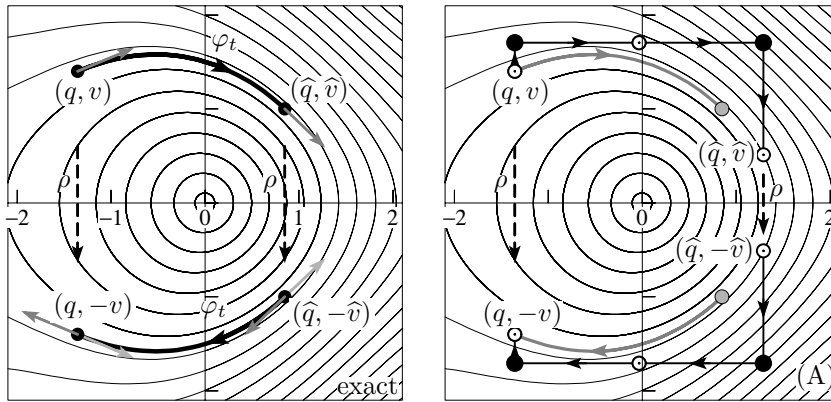


Figure 2.1. A reversible system (left) and the symmetric Störmer–Verlet method (right); the same equation as in Figure 1.6.

The time-symmetry of the Störmer–Verlet method implies an important geometric property of the numerical map in the phase space, namely *reversibility*, to which we turn next. The importance of this property in numerical analysis was first emphasized by Stoffer (1988).

The system (1.3) has the property that inverting the direction of the initial velocity does not change the solution trajectory, it just inverts the direction of motion. The flow  $\varphi_t$  thus satisfies that

$$\varphi_t(q, v) = (\hat{q}, \hat{v}) \quad \text{implies} \quad \varphi_t(\hat{q}, -\hat{v}) = (q, -v), \tag{2.2}$$

and we call it *reversible* with respect to the reflection  $\rho : (q, v) \mapsto (q, -v)$ . This property is illustrated in Figure 2.1, left. The numerical one-step map  $\Phi_h$  of the Störmer–Verlet method satisfies similarly

$$\Phi_h(q, v) = (\hat{q}, \hat{v}) \quad \text{implies} \quad \Phi_h(\hat{q}, -\hat{v}) = (q, -v), \tag{2.3}$$

for all  $q, v$  and all  $h$ ; see Figure 2.1, right. This holds because practically all numerical methods for (1.3), and in particular the Störmer–Verlet method and the symplectic Euler methods, are such that

$$\Phi_h(q, v) = (\hat{q}, \hat{v}) \quad \text{implies} \quad \Phi_{-h}(q, -v) = (\hat{q}, -\hat{v}), \tag{2.4}$$

as is readily seen from the defining formulae such as (1.5). The symmetry (2.1) of the Störmer–Verlet method is therefore equivalent to the reversibility (2.3). Let us summarize these considerations.

**Theorem 2.1.** The Störmer–Verlet method applied to the second-order differential equation (1.1) is both symmetric and reversible, *i.e.*, its one-step map satisfies (2.1) and (2.3).

In some situations, the flow is  $\rho$ -reversible with respect to involutions  $\rho$  other than  $(q, v) \mapsto (q, -v)$ , that is, it satisfies

$$\rho \circ \varphi_t = \varphi_t^{-1} \circ \rho. \quad (2.5)$$

For example, the flow of the Kepler problem (1.19) is  $\rho$ -reversible also with respect to  $\rho : (q_1, q_2, v_1, v_2) \mapsto (q_1, -q_2, -v_1, v_2)$ . In general, the flow of a differential equation  $\dot{y} = F(y)$  is  $\rho$ -reversible if and only if the vector field satisfies  $\rho \circ F = -F \circ \rho$ . We then call the differential equation  $\rho$ -reversible.

By the same argument as above, the Störmer–Verlet method is then also  $\rho$ -reversible for  $\rho$  of the form  $\rho(q, v) = (\rho_1(q), \rho_2(v))$ , that is,

$$\rho \circ \Phi_h = \Phi_h^{-1} \circ \rho. \quad (2.6)$$

## 2.2. Hamiltonian systems and symplecticity

We now turn to the important class of *Hamiltonian systems*

$$\dot{p} = -\nabla_q H(p, q), \quad \dot{q} = \nabla_p H(p, q), \quad (2.7)$$

where  $H(p, q)$  is an arbitrary scalar function of the variables  $(p, q)$ . When the Hamiltonian is of the form

$$H(p, q) = \frac{1}{2} p^T M^{-1} p + U(q), \quad (2.8)$$

with a positive definite mass matrix  $M$  and a potential  $U(q)$ , then the system (2.7) turns into the second-order differential equation (1.3) upon expressing the momenta  $p = Mv$  in terms of the velocities and setting  $f(q) = -M^{-1}\nabla U(q)$ . Equation (2.8) expresses the total energy  $H$  as the sum of kinetic and potential energy.

A characteristic geometric property of Hamiltonian systems is that the flow  $\varphi_t$  is *symplectic*, that is, the derivative  $\varphi'_t = \partial\varphi_t/\partial(p, q)$  of the flow satisfies, for all  $(p, q)$  and  $t$  where  $\varphi_t(p, q)$  is defined,

$$\varphi'_t(p, q)^T J \varphi'_t(p, q) = J \quad \text{with} \quad J = \begin{pmatrix} 0 & I \\ -I & 0 \end{pmatrix}, \quad (2.9)$$

where  $I$  is the identity matrix of the dimension of  $p$  or  $q$ ; see, *e.g.*, Arnold (1989, p. 204) or Hairer *et al.* (2002, p. 172).

The relation (2.9) is formally similar to orthogonality (which it would be if  $J$  were replaced by the identity matrix), but, unlike orthogonality, it is not related to the conservation of lengths but of *areas* in phase space. In fact, for systems with one degree of freedom (*i.e.*,  $p, q \in \mathbb{R}$ ), equation (2.9) expresses that the flow preserves the area of sets of initial values in the  $(p, q)$ -plane; see Figure 2.2, left. For higher-dimensional systems, symplecticity (2.9) means that the flow preserves the sum of the oriented areas of the projections of  $\varphi_t(A)$  onto the  $(p_i, q_i)$ -coordinate planes, for any two-dimensional bounded

manifold of initial values  $A$ ; see, *e.g.*, Hairer *et al.* (2002, p. 171*f.*) for a justification of this interpretation.

The Störmer–Verlet method (1.24) applied to the Hamiltonian system (2.7) reads

$$\begin{aligned}
 p_{n+\frac{1}{2}} &= p_n - \frac{h}{2} \nabla_q H(p_{n+\frac{1}{2}}, q_n), \\
 \text{(A)} \quad q_{n+1} &= q_n + \frac{h}{2} \left( \nabla_p H(p_{n+\frac{1}{2}}, q_n) + \nabla_p H(p_{n+\frac{1}{2}}, q_{n+1}) \right), \\
 p_{n+1} &= p_{n+\frac{1}{2}} - \frac{h}{2} \nabla_q H(p_{n+\frac{1}{2}}, q_{n+1}),
 \end{aligned}
 \tag{2.10}$$

and a similar formula for variant (B). In the particular case of the Hamiltonian (2.8), the method reduces to the Störmer–Verlet method (1.5) with  $f(q) = -M^{-1} \nabla U(q)$ , upon setting  $p_n = Mv_n$ .

A numerical method is called *symplectic* if, for Hamiltonian systems (2.7), the Jacobian of the numerical flow  $\Phi_h : (p_n, q_n) \mapsto (p_{n+1}, q_{n+1})$  satisfies condition (2.9), that is, if

$$\Phi'_h(p, q)^T J \Phi'_h(p, q) = J \tag{2.11}$$

for all  $(p, q)$  and all step sizes  $h$ .

Symplecticity of numerical methods was first considered by De Vogelaere (1956), but was not followed up until Ruth (1983) and Feng (1985). In the late 1980s, the results of Lasagni (1988), Sanz-Serna (1988), and Suris (1988) started off an avalanche of papers on symplectic numerical methods. Sanz-Serna and Calvo (1994) was the first book dealing with this subject.

**Theorem 2.2.** The Störmer–Verlet method applied to a Hamiltonian system is symplectic.

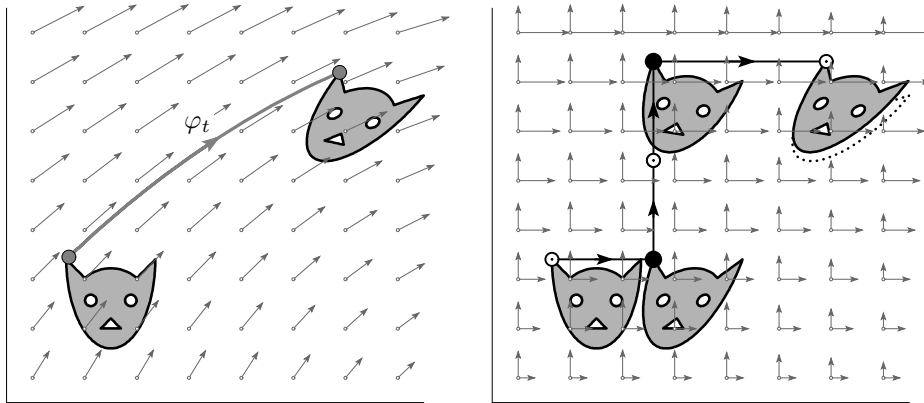


Figure 2.2. Symplecticity of the Störmer–Verlet method for a separable Hamiltonian.

We give four different proofs of this result, which all correspond to different interpretations of the method: as a composition method, as a splitting method, as a variational integrator, and using generating functions. Each of these interpretations lends itself to generalizations to other symplectic integrators, of higher order and/or for constrained Hamiltonian systems. Yet another proof is based on the preservation of quadratic invariants and will be mentioned in Section 3 below. The second proof applies only to Hamiltonians of the special form (2.8), the third proof is formulated for such Hamiltonians for convenience.

The historically *first proof*, due to De Vogelaere (1956), uses the interpretation of the Störmer–Verlet method as the composition of the symplectic Euler method

$$(SE1) \quad \begin{aligned} p_{n+\frac{1}{2}} &= p_n - \frac{h}{2} \nabla_q H(p_{n+\frac{1}{2}}, q_n), \\ q_{n+\frac{1}{2}} &= q_n + \frac{h}{2} \nabla_p H(p_{n+\frac{1}{2}}, q_n), \end{aligned} \quad (2.12)$$

and its adjoint

$$(SE2) \quad \begin{aligned} q_{n+1} &= q_{n+\frac{1}{2}} + \frac{h}{2} \nabla_p H(p_{n+\frac{1}{2}}, q_{n+1}), \\ p_{n+1} &= p_{n+\frac{1}{2}} - \frac{h}{2} \nabla_q H(p_{n+\frac{1}{2}}, q_{n+1}). \end{aligned} \quad (2.13)$$

The method (SE1) is indeed symplectic, as is seen by direct verification of the symplecticity condition

$$\left( \frac{\partial(p_{n+1/2}, q_{n+1/2})}{\partial(p_n, q_n)} \right)^T J \left( \frac{\partial(p_{n+1/2}, q_{n+1/2})}{\partial(p_n, q_n)} \right) = J.$$

The matrix of partial derivatives is obtained from differentiating equation (2.12):

$$\begin{pmatrix} I + hH_{qp}^T & 0 \\ -hH_{pp} & I \end{pmatrix} \begin{pmatrix} \partial(p_{n+1/2}, q_{n+1/2}) \\ \partial(p_n, q_n) \end{pmatrix} = \begin{pmatrix} I & -hH_{qq} \\ 0 & I + hH_{qp} \end{pmatrix},$$

where all the submatrices of the Hessian,  $H_{qp}, H_{pp}$ , etc., are evaluated at  $(p_{n+1/2}, q_n)$ . In the same way, (SE2) is seen to be symplectic. Hence their composition (2.10) is also symplectic.

The *second proof* is the most elegant one, but it applies only to the case of separable Hamiltonians  $H(p, q) = T(p) + U(q)$ . It is based on the interpretation of the Störmer–Verlet method as a splitting method. As in (1.13), we have for variant (A)

$$\Phi_h = \varphi_{h/2}^U \circ \varphi_h^T \circ \varphi_{h/2}^U, \quad (2.14)$$

where  $\varphi_t^T$  and  $\varphi_t^U$  are the exact flows of the Hamiltonian systems with Hamiltonian  $T(p) = \frac{1}{2} p^T M^{-1} p$  and  $U(q)$ , i.e.,  $\dot{p} = 0$ ,  $\dot{q} = M^{-1} p$  and  $\dot{p} = -\nabla U(q)$ ,  $\dot{q} = 0$ , respectively, corresponding to the splitting  $H(p, q) =$

$T(p)+U(q)$  of the Hamiltonian (2.8) into kinetic and potential energy. Since the flows of Hamiltonian systems are symplectic, so is their composition (2.14). This is illustrated in Figure 2.2, right. Variant (B) has the flows of  $T$  and  $U$  interchanged in (2.14), and is thus likewise symplectic.

The *third proof* uses the interpretation of the Störmer–Verlet method as a variational integrator (see Section 1.6). The symplecticity of variational integrators derives from non-numerical work by Maeda (1980) and Veselov (1991). Using (1.4) and the first line of (1.5), we have for  $S_h(q, Q)$  of (1.17), in the case of the Lagrangian (1.15) which corresponds to the Hamiltonian (2.8),

$$-\nabla_q S_h(q_n, q_{n+1}) = M \frac{q_{n+1} - q_n}{h} + \frac{h}{2} \nabla U(q_n) = M v_n = p_n \quad (2.15)$$

and similarly

$$\nabla_Q S_h(q_n, q_{n+1}) = M \frac{q_{n+1} - q_n}{h} - \frac{h}{2} \nabla U(q_{n+1}) = M v_{n+1} = p_{n+1}. \quad (2.16)$$

Given  $(p_n, q_n)$ , the first of the above two equations determines  $q_{n+1}$ , and the second one  $p_{n+1}$ . The one-step map  $\Phi_h : (p_n, q_n) \mapsto (p_{n+1}, q_{n+1})$  of the Störmer–Verlet method is thus generated by the scalar-valued function  $S_h$  via (2.15) and (2.16). The desired result then follows from the fact that a map  $(p, q) \mapsto (P, Q)$  generated by

$$-\nabla_q S(q, Q) = p, \quad \nabla_Q S(q, Q) = P,$$

is symplectic for *any* function  $S$ . This is verified by directly checking the symplecticity condition. Differentiation of the above equations gives the following relations for the matrices of partial derivatives  $P_p, P_q, Q_p, Q_q$ :

$$\begin{aligned} S_{qq} + S_{qQ}Q_q &= 0, & S_{qQ}Q_p &= I, \\ S_{Qq} + S_{QQ}Q_q &= P_q, & S_{QQ}Q_p &= P_p. \end{aligned}$$

These equations yield

$$\begin{pmatrix} P_p & P_q \\ Q_p & Q_q \end{pmatrix}^T \begin{pmatrix} 0 & I \\ -I & 0 \end{pmatrix} \begin{pmatrix} P_p & P_q \\ Q_p & Q_q \end{pmatrix} = \begin{pmatrix} 0 & I \\ -I & 0 \end{pmatrix}$$

after multiplying out, as is required for symplecticity. This completes the third proof of symplecticity of the Störmer–Verlet method.

A *fourth proof* of the symplecticity is based on ideas of Lasagni (1988). A step of the Störmer–Verlet method can be generated by a function  $\widehat{S}_h(p_1, q_0)$  in the same way as the symplectic Euler method:

$$\begin{aligned} p_1 &= p_0 - \nabla_q \widehat{S}_h(p_1, q_0), \\ q_1 &= q_0 + \nabla_p \widehat{S}_h(p_1, q_0). \end{aligned} \quad (2.17)$$

As we have seen in the first proof, such maps are symplectic. The generating

function is simply  $\widehat{S}_h = hH$  for the symplectic Euler method. For the Störmer–Verlet method  $\widehat{S}_h$  is obtained as

$$\begin{aligned} \widehat{S}_h(p_1, q_0) &= \frac{h}{2} \left( H(p_{1/2}, q_0) + H(p_{1/2}, q_1) \right) \\ &\quad - \frac{h^2}{4} \nabla_q H(p_{1/2}, q_1)^T \left( \nabla_p H(p_{1/2}, q_0) + \nabla_p H(p_{1/2}, q_1) \right), \end{aligned} \quad (2.18)$$

where  $q_1$  and  $p_{1/2}$  are defined by the Störmer–Verlet formulae and are now considered as functions of  $(p_1, q_0)$ . We do not give the computational details, which can be found in Hairer *et al.* (2002, Section VI.5) for a more general class of symplectic integrators.

### 2.3. Volume preservation

The flow  $\varphi_t$  of a system of differential equations  $\dot{y} = F(y)$  with divergence-free vector field ( $\operatorname{div} F(y) = 0$  for all  $y$ ) satisfies  $\det \varphi'_t(y) = 1$  for all  $y$ . It therefore preserves volume in phase space: for every bounded open set  $\Omega$ , and for every  $t$  for which  $\varphi_t(y)$  exists for all  $y \in \Omega$ ,

$$\operatorname{vol}(\varphi_t(\Omega)) = \operatorname{vol}(\Omega).$$

The vector field  $(v, f(q))$  of a second-order differential equation (1.3), written as a first-order system, is divergence-free. The same is true for Hamiltonian vector fields  $(-\nabla_q H(p, q), \nabla_p H(p, q))$ .

The Störmer–Verlet method preserves volume,

$$\operatorname{vol}(\Phi_h(\Omega)) = \operatorname{vol}(\Omega),$$

in the following two situations.

For the method (2.10), applied to a Hamiltonian system (2.7), this follows from its symplecticity (2.11), which implies  $\det \Phi'_h(p, q) = 1$  for all  $(p, q)$ .

For partitioned differential equations of the form

$$\dot{q} = g(v), \quad \dot{v} = f(q), \quad (2.19)$$

the method (1.24) can be interpreted as the splitting (1.13), where  $\varphi_t^{[1]}$  and  $\varphi_t^{[2]}$  are the exact flows of  $\dot{q} = g(v)$ ,  $\dot{v} = 0$  and  $\dot{q} = 0$ ,  $\dot{v} = f(q)$ , respectively. Since the vector fields of these flows are divergence-free, they are volume-preserving and so is their composition.

The same idea allows us to extend the Störmer–Verlet method to a volume-preserving algorithm for systems partitioned into the *three* equations

$$\dot{x} = a(y, z), \quad \dot{y} = b(x, z), \quad \dot{z} = c(x, y), \quad (2.20)$$

for which the diagonal blocks of the Jacobian are zero. We split them symmetrically, giving

$$\varphi_{h/2}^{[1]} \circ \varphi_{h/2}^{[2]} \circ \varphi_h^{[3]} \circ \varphi_{h/2}^{[2]} \circ \varphi_{h/2}^{[1]}, \quad (2.21)$$



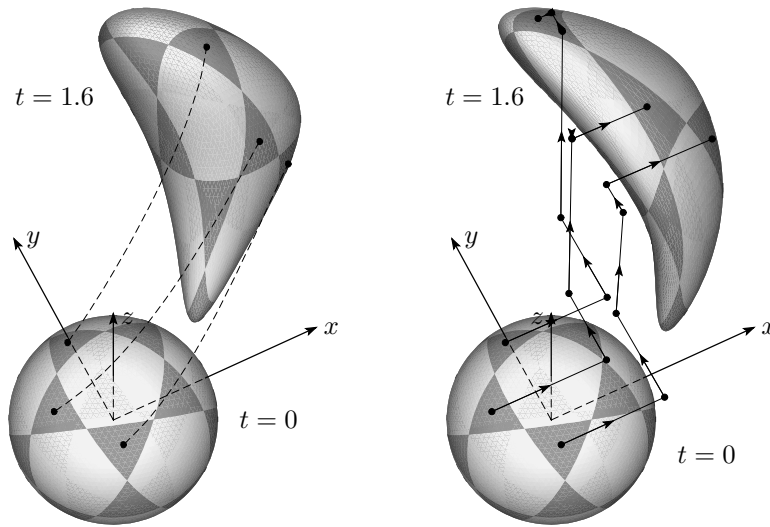


Figure 2.3. Volume-preserving deformation of the ball of radius 0.9, centred at the origin, by the ABC flow (left) and by method (2.22) (right).

where  $\varphi_t^{[1]}$  is the (volume-preserving) flow of  $\dot{x} = a(y, z)$ ,  $\dot{y} = 0$ ,  $\dot{z} = 0$  and similarly for  $\varphi_t^{[2]}$  and  $\varphi_t^{[3]}$ . Written out, this becomes

$$\begin{aligned}
 x_{n+\frac{1}{2}} &= x_n + \frac{h}{2} a(y_n, z_n), \\
 y_{n+\frac{1}{2}} &= y_n + \frac{h}{2} b(x_{n+\frac{1}{2}}, z_n), \\
 z_{n+1} &= z_n + h c(x_{n+\frac{1}{2}}, y_{n+\frac{1}{2}}), \\
 y_{n+1} &= y_{n+\frac{1}{2}} + \frac{h}{2} b(x_{n+\frac{1}{2}}, z_{n+1}), \\
 x_{n+1} &= x_{n+\frac{1}{2}} + \frac{h}{2} a(y_{n+1}, z_{n+1}).
 \end{aligned}
 \tag{2.22}$$

An illustration of this algorithm, applied to the ABC-flow

$$\begin{aligned}
 \dot{x} &= A \sin z + C \cos y, \\
 \dot{y} &= B \sin x + A \cos z, \\
 \dot{z} &= C \sin y + B \cos x,
 \end{aligned}$$

is presented in Figure 2.3 for  $A = 1/2$ ,  $B = C = 1$ .

More ingenuity is necessary if the system is divergence-free with nonzero elements on the diagonal of the Jacobian. Feng and Shang (1995) give a volume-preserving extension of the above scheme to the general case.

### 3. Conservation of first integrals

A non-constant function  $I(y)$  is a *first integral* (or conserved quantity, or constant of motion, or invariant) of the differential equation  $\dot{y} = F(y)$  if  $I(y(t))$  is constant along every solution, or equivalently, if

$$I'(y)F(y) = 0 \quad \text{for all } y. \quad (3.1)$$

The latter condition says that the gradient  $\nabla I(y)$  is orthogonal to the vector field  $F(y)$  in every point of the phase space.

The foremost example is the Hamiltonian  $H(p, q)$  of a Hamiltonian system (2.7): since  $H' = (\nabla_p H^T, \nabla_q H^T)$  and  $\nabla_p H^T(-\nabla_q H) + \nabla_q H^T \nabla_p H = 0$ , the total energy  $H$  is a first integral. Apart from very exceptional cases,  $H$  is *not* constant along numerical solutions computed with the Störmer–Verlet method. Later we will see, however, that  $H$  is conserved up to  $\mathcal{O}(h^2)$  over extremely long time intervals.

**Example 3.1.** The Kepler problem (1.19) is Hamiltonian with  $H(p, q) = \frac{1}{2}(p_1^2 + p_2^2) - 1/\sqrt{q_1^2 + q_2^2}$ . In addition to the Hamiltonian, this system has the following conserved quantities, as can be easily checked: the angular momentum  $L = q_1 p_2 - q_2 p_1$ , and the nonzero components of the Runge–Lenz–Pauli vector

$$\begin{pmatrix} A_1 \\ A_2 \\ 0 \end{pmatrix} = \begin{pmatrix} p_1 \\ p_2 \\ 0 \end{pmatrix} \times \begin{pmatrix} 0 \\ 0 \\ q_1 p_2 - q_2 p_1 \end{pmatrix} - \frac{1}{\sqrt{q_1^2 + q_2^2}} \begin{pmatrix} q_1 \\ q_2 \\ 0 \end{pmatrix}.$$

Figure 3.1 shows the behaviour of these quantities along a numerical solution of the Störmer–Verlet method. The method preserves the angular momentum exactly (see Section 1.3), and there are only small errors in the Hamiltonian along the numerical solution, but no drift. There is, however, a linear drift in the Runge–Lenz–Pauli vector. In contrast, for explicit Runge–Kutta methods, none of the first integrals is preserved, and there is a drift away from the constant value for all of them.

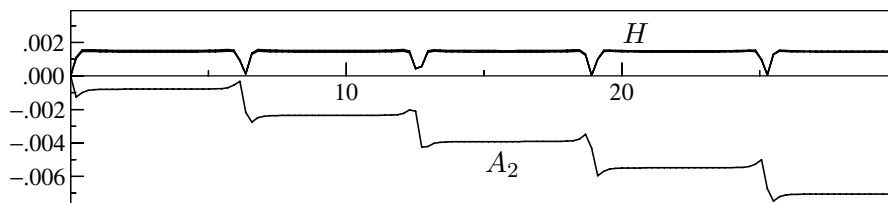


Figure 3.1. The Hamiltonian and the second component of the Runge–Lenz–Pauli vector along the numerical solution of the Störmer–Verlet method with step size  $h = 0.02$ .

**Example 3.2. (Conservation of total linear and angular momentum of  $N$ -body systems)** A system of  $N$  particles interacting pairwise, with potential forces depending on the distances of the particles, is formulated as a Hamiltonian system with total energy

$$H(p, q) = \frac{1}{2} \sum_{i=1}^N \frac{1}{m_i} p_i^T p_i + \sum_{i=2}^N \sum_{j=1}^{i-1} V_{ij}(\|q_i - q_j\|). \quad (3.2)$$

Here  $q_i, p_i \in \mathbb{R}^3$  represent the position and momentum of the  $i$ th particle of mass  $m_i$ , and  $V_{ij}(r)$  ( $i > j$ ) is the interaction potential between the  $i$ th and  $j$ th particle. The equations of motion read

$$\dot{q}_i = \frac{1}{m_i} p_i, \quad \dot{p}_i = \sum_{j=1}^N \nu_{ij} (q_i - q_j)$$

where, for  $i > j$ , we have  $\nu_{ij} = \nu_{ji} = -V'_{ij}(r_{ij})/r_{ij}$  with  $r_{ij} = \|q_i - q_j\|$ , and  $\nu_{ii}$  is arbitrary, say  $\nu_{ii} = 0$ . The conservation of the total *linear momentum*  $P = \sum_{i=1}^N p_i$  and the total *angular momentum*  $L = \sum_{i=1}^N q_i \times p_i$  is a consequence of the symmetry relation  $\nu_{ij} = \nu_{ji}$ :

$$\begin{aligned} \frac{d}{dt} \sum_{i=1}^N p_i &= \sum_{i=1}^N \sum_{j=1}^N \nu_{ij} (q_i - q_j) = 0, \\ \frac{d}{dt} \sum_{i=1}^N q_i \times p_i &= \sum_{i=1}^N \frac{1}{m_i} p_i \times p_i + \sum_{i=1}^N \sum_{j=1}^N q_i \times \nu_{ij} (q_i - q_j) = 0. \end{aligned}$$

The exact preservation of *linear* first integrals, such as the total linear momentum, is common to most numerical integrators.

**Theorem 3.3.** The Störmer–Verlet method preserves linear first integrals.

*Proof.* Let the linear first integral be  $I(q, v) = b^T q + c^T v$ , so that  $b^T v + c^T f(q) = 0$  for all  $q, v$ . Necessarily then,  $c^T f(q) = 0$  for all  $q$ , and  $b = 0$ . Multiplying the formulae for  $v$  in (1.5) by  $c^T$  thus yields  $c^T v_1 = c^T v_0$ .  $\square$

*Quadratic first integrals* are not generally preserved by the Störmer–Verlet method, as the following example shows.

**Example 3.4.** Consider the harmonic oscillator, which has the Hamiltonian  $H(p, q) = \frac{1}{2}p^2 + \frac{1}{2}\omega^2 q^2$  ( $p, q \in \mathbb{R}$ ). Applying the Störmer–Verlet method gives

$$\begin{pmatrix} p_{n+1} \\ \omega q_{n+1} \end{pmatrix} = A(h\omega) \begin{pmatrix} p_n \\ \omega q_n \end{pmatrix} \quad (3.3)$$

with the propagation matrix

$$A(h\omega) = \begin{pmatrix} 1 - \frac{h^2\omega^2}{2} & -\frac{h\omega}{2} \left(1 - \frac{h^2\omega^2}{4}\right) \\ \frac{h\omega}{2} & 1 - \frac{h^2\omega^2}{2} \end{pmatrix}. \quad (3.4)$$

Since  $A(h\omega)$  is not an orthogonal matrix,  $H(p, q)$  is not preserved along numerical solutions. Notice, however, that the characteristic polynomial is  $\lambda^2 - (2 - h^2\omega^2)\lambda + 1$ , so that the eigenvalues are of modulus one if (and only if)  $|h\omega| \leq 2$ . The matrix  $V$  of eigenvectors is close to the identity for small  $h\omega$ , and the norm of  $V^{-1}(p_n, \omega q_n)^T$  is conserved.

The Störmer–Verlet method does, however, preserve an important subclass of quadratic first integrals, and in particular the total angular momentum of  $N$ -body systems. As we have seen in Section 1.3, Newton was aware that the method preserves angular momentum in the Kepler problem and used this fact to prove Kepler’s second law. In the following result  $C$  is a constant square matrix and  $c$  a constant vector.

**Theorem 3.5.** The Störmer–Verlet method preserves quadratic first integrals of the form  $I(q, v) = v^T(Cq + c)$  (or  $I(p, q) = p^T(Bq + b)$  in the Hamiltonian case).

*Proof.* By (3.1),  $f(q)^T(Cq + c) + v^T C v = 0$  for all  $q, v$ . Writing the Störmer–Verlet method as the composition of the two symplectic Euler methods (1.8) and (1.9), we obtain for the first half-step

$$\begin{aligned} v_{n+1/2}^T(Cq_{n+1/2} + c) &= v_n^T(Cq_n + c) \\ &\quad + \frac{h}{2} (f(q_n)^T(Cq_n + c) + v_{n+1/2}^T C v_{n+1/2}), \end{aligned}$$

where we notice that the term in the second line vanishes. For the second half-step we obtain in the same way  $v_{n+1}^T(Cq_{n+1} + c) = v_{n+1/2}^T(Cq_{n+1/2} + c)$ , and the result follows.  $\square$

The most important source of first integrals of Hamiltonian systems is *Noether’s theorem*, which states that continuous symmetries yield first integrals: if the associated Lagrangian is invariant under the flow  $\alpha_s$  of the vector field  $a(q)$ , that is,  $L(\alpha_s(q), \alpha_s'(q)v) = L(q, v)$  for all real  $s$  near 0 and all  $(q, v)$ , then  $I(p, q) = p^T a(q)$  is a first integral; see, *e.g.*, Arnold (1989, p. 88). For Hamiltonian systems of the form (2.8), where the associated Lagrangian is (1.15), it can be shown that  $a(q)$  must be linear:  $a(q) = Bq + b$ , with  $MB$  skew-symmetric. Hence, for Hamiltonian systems (2.7) with a Hamiltonian of the form (2.8), all first integrals originating from Noether’s theorem are preserved by the Störmer–Verlet method.

Theorem 3.5 yields yet another proof (and further insight) of the symplecticity of the Störmer–Verlet method, following an argument by Bochev and

Scovel (1994): consider the Hamiltonian system  $\dot{p} = -\nabla U(q)$ ,  $\dot{q} = M^{-1}p$  together with its variational equation

$$\dot{Y} = \begin{pmatrix} 0 & -\nabla^2 U(q) \\ M^{-1} & 0 \end{pmatrix} Y \quad \text{with } Y = \begin{pmatrix} P_p & P_q \\ Q_p & Q_q \end{pmatrix}.$$

The derivative of the flow is then  $\varphi'_t(p, q) = Y(t)$  corresponding to the initial conditions  $p, q$  and  $Y(0) = I$ . The derivative  $\Phi'_h(p, q)$  of the numerical solution with respect to the initial values equals the result  $Y_1$  obtained by applying the method to the combined system of the Hamiltonian system together with its variational equation, partitioned into  $(p, P_p, P_q)$  and  $(q, Q_p, Q_q)$ . Symplecticity means that the components of  $Y^T J Y$  are first integrals. Since they are of the mixed quadratic type considered above, Theorem 3.5 shows that they are preserved by the Störmer–Verlet method:  $Y_1^T J Y_1 = Y_0^T J Y_0$ , which is just the symplecticity  $\Phi_h'^T J \Phi_h' = J$ .

#### 4. Backward error analysis

The theoretical foundation of geometric integrators is mainly based on a backward interpretation which considers the numerical approximation as the exact solution of a modified problem. Such an interpretation has been intuitively used in the physics literature, *e.g.*, Ruth (1983). A rigorous formulation evolved around 1990, beginning with the papers by Feng (1991), McLachlan and Atela (1992), Sanz-Serna (1992) and Yoshida (1993). Exponentially small error bounds and applications of backward error analysis to explain the long-time behaviour of numerical integrators were subsequently given by Benettin and Giorgilli (1994), Hairer and Lubich (1997), and Reich (1999*a*). We explain the essential ideas and we illustrate them for the Störmer–Verlet method.

##### 4.1. Construction of the modified equation

The idea of backward error analysis applies to general ordinary differential equations and to general numerical integrators, and a restriction to special methods for second-order problems would hide the essentials. We therefore consider the differential equation

$$\dot{y} = F(y), \tag{4.1}$$

and a numerical one-step method  $y_{n+1} = \Phi_h(y_n)$ . The idea consists in searching and studying a *modified differential equation*

$$\dot{y} = F(y) + hF_2(y) + h^2F_3(y) + \dots, \tag{4.2}$$

such that the exact time- $h$  flow  $\tilde{\varphi}_h(y)$  of (4.2) is equal to the numerical flow  $\Phi_h(y)$ . Unfortunately, the series in (4.2) cannot be expected to converge in general, and the precise statement has to be formulated as follows.

**Theorem 4.1.** Consider (4.1) with an infinitely differentiable vector field  $F(y)$ , and assume that the numerical method admits a Taylor series expansion of the form

$$\Phi_h(y) = y + hF(y) + h^2D_2(y) + h^3D_3(y) + \dots \quad (4.3)$$

with smooth  $D_j(y)$ . Then there exist unique vector fields  $F_j(y)$  such that, for any  $N \geq 1$ ,

$$\Phi_h(y) = \tilde{\varphi}_{h,N}(y) + \mathcal{O}(h^{N+1}),$$

where  $\tilde{\varphi}_{t,N}$  is the exact flow of the truncated modified equation

$$\dot{y} = F(y) + hF_2(y) + \dots + h^{N-1}F_N(y). \quad (4.4)$$

*Proof.* Disregarding convergence issues, we expand the exact flow of (4.2) into a Taylor series (using the notation  $\tilde{y}(t) = \tilde{\varphi}_t(y)$ )

$$\begin{aligned} \tilde{\varphi}_h(y) &= y + h\dot{\tilde{y}}(0) + \frac{h^2}{2!}\ddot{\tilde{y}}(0) + \frac{h^3}{3!}\tilde{y}^{(3)}(0) + \dots \\ &= y + h(F(y) + hF_2(y) + h^2F_3(y) + \dots) \\ &\quad + \frac{h^2}{2!}(F'(y) + hF_2'(y) + \dots)(F(y) + hF_2(y) + \dots) + \dots \end{aligned} \quad (4.5)$$

and we compare like powers of  $h$  in the expressions (4.5) and (4.3). This yields recurrence relations for the functions  $F_j(y)$ , namely,

$$\begin{aligned} F_2(y) &= D_2(y) - \frac{1}{2!}F'F(y), \\ F_3(y) &= D_3(y) - \frac{1}{3!}(F''(F, F)(y) + F'F'F(y)) - \frac{1}{2!}(F'F_2(y) + F_2'F(y)), \end{aligned} \quad (4.6)$$

and uniquely defines the functions  $F_j(y)$  in a constructive manner.  $\square$

#### 4.2. Modified equation of the Störmer–Verlet method

Putting  $y = (q, v)^T$  and  $F(y) = (v, f(q))^T$ , the differential equation (1.3) is of the form (4.1). For the Störmer–Verlet scheme (1.5) we have

$$\Phi_h(q, v) = \begin{pmatrix} q + hv + \frac{h^2}{2}f(q) \\ v + \frac{h}{2}f(q) + \frac{h}{2}f(q + hv + \frac{h^2}{2}f(q)) \end{pmatrix}. \quad (4.7)$$

Expanding this function into a Taylor series we get (4.3) with

$$D_2(q, v) = \frac{1}{2} \begin{pmatrix} f(q) \\ f'(q)v \end{pmatrix}, \quad D_3(q, v) = \frac{1}{4} \begin{pmatrix} 0 \\ f'(q)f(q) + f''(q)(v, v) \end{pmatrix}, \quad \dots$$

and the functions  $F_j(q, v)$  can be computed as in the proof of Theorem 4.1. Since the Störmer–Verlet method is of second order, the function  $D_2(q, v)$

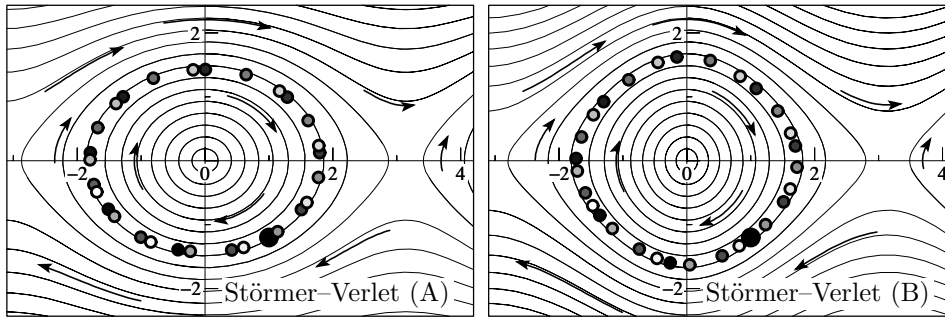


Figure 4.1. Numerical solution with step size  $h = 0.9$  for the two versions of the Störmer–Verlet method compared to the exact flow of their modified differential equations truncated after the  $\mathcal{O}(h^2)$  term.

has to coincide with the  $h^2$ -coefficient of the exact solution and we have  $F_2(q, v) = 0$ . We then get

$$F_3(q, v) = \frac{1}{12} \begin{pmatrix} -2 f'(q)v \\ f'(q)f(q) + f''(q)(v, v) \end{pmatrix}, \quad (4.8)$$

and for the next function we obtain  $F_4(q, v) = 0$ . The vanishing of this function follows from the symmetry of the method (*cf.* Section 4.3). For larger (odd)  $j$  the functions  $F_j(q, v)$  become more and more complicated, and higher derivatives of  $f(q)$  are involved. The explicit formula for  $F_3(q, v)$  also shows that the modified differential equation (4.2) is no longer a second-order equation like (1.3).

A similar computation for the version (B) of the Störmer–Verlet method (see (1.6)) gives

$$F_3(q, v) = \frac{1}{24} \begin{pmatrix} 2 f'(q)v \\ -4 f'(q)f(q) - f''(q)(v, v) \end{pmatrix}, \quad (4.9)$$

and, obviously, also  $F_2(q, v) = F_4(q, v) = 0$ .

As a concrete example consider the pendulum equation for which  $f(q) = -\sin q$ . The two pictures of Figure 4.1 show the exact flow of the modified differential equations (truncated after the  $\mathcal{O}(h^2)$  term) corresponding to the two versions (1.5) and (1.6) of the Störmer–Verlet scheme together with the numerical solution for the initial value  $(p_0, q_0) = (1.0, -1.2)$ . The shade of the numerical approximations (dark grey to light grey) indicates the increasing time. We observe a surprisingly good agreement.

In both cases the solutions of the modified equation are periodic, and the numerical approximation lies near a closed curve, so that correct qualitative behaviour is obtained. This is explained by the fact that for  $f(q) = -\nabla U(q)$

the vector fields (4.8) and (4.9) are Hamiltonian with

$$\begin{aligned} H_3(p, q) &= \frac{1}{12} \nabla^2 U(q)(p, p) + \frac{1}{24} \nabla U(q)^T \nabla U(q) \quad \text{and} \\ H_3(p, q) &= -\frac{1}{24} \nabla^2 U(q)(p, p) - \frac{1}{6} \nabla U(q)^T \nabla U(q), \end{aligned}$$

respectively. Consequently, the exact solutions of the truncated modified equation stay on the level curves of  $\tilde{H}(p, q) = H(p, q) + h^2 H_3(p, q)$  which are drawn in Figure 4.1.

#### 4.3. Properties of the modified differential equation

In Section 4.4 we shall see that the numerical solution is extremely close to the exact solution of a truncated modified equation. To study properties of the numerical solution it is therefore justified to investigate instead the corresponding properties of the modified differential equation.

It follows from the definition of the modified equation that for methods of order  $r$ , that is,  $\Phi_h(y) = \varphi_h(y) + \mathcal{O}(h^{r+1})$ , we have

$$F_j(y) = 0 \quad \text{for } j = 2, \dots, r.$$

Furthermore, if the leading term of the *local truncation error* is  $E_{r+1}(y)$ , that is,  $\Phi_h(y) = \varphi_h(y) + h^{r+1} E_{r+1}(y) + \mathcal{O}(h^{r+2})$ , then

$$F_{r+1}(y) = E_{r+1}(y).$$

By Theorem 2.1 the Störmer–Verlet method is *symmetric*. For such methods the modified equation has an expansion in even powers of  $h$ , that is,

$$F_{2j}(y) = 0 \quad \text{for } j = 1, 2, \dots \quad (4.10)$$

This can be proved as follows: to indicate the  $h$ -dependence of the vector field (4.2), we let  $\tilde{\varphi}_{t,h}(y)$  denote the (formal) flow of (4.2). Backward error analysis tells us that  $\Phi_h(y) = \tilde{\varphi}_{h,h}(y)$ . We thus have  $\Phi_{-h}(y) = \tilde{\varphi}_{-h,-h}(y)$  and, by the group property of the exact flow,  $\Phi_{-h}^{-1}(y) = \tilde{\varphi}_{h,-h}(y)$ . The symmetry condition (2.1) thus implies that  $\tilde{\varphi}_{t,h}(y) = \tilde{\varphi}_{t,-h}(y)$  for  $t = h$ , and the computation of (4.5) shows that this is only possible if (4.10) holds.

Geometric properties of a numerical method have their counterparts in the modified equation. Let us explain this for the properties discussed in Sections 2 and 3.

**Theorem 4.2. (Reversible systems)** If the Störmer–Verlet method (1.5) is applied to a differential equation (1.3), then every truncation of the modified differential equation is reversible with respect to the reflection  $\rho(q, v) = (q, -v)$ .

**Theorem 4.3. (Hamiltonian systems)** If the Störmer–Verlet method (2.10) is applied to a Hamiltonian system, then every truncation of the modified differential equation is Hamiltonian.



**Theorem 4.4. (Divergence-free systems)** If the Störmer–Verlet method (1.24) is applied to a divergence-free system of the form (2.19), then every truncation of the modified differential equation is divergence-free.

**Theorem 4.5. (First integrals)** If the Störmer–Verlet method (1.24) is applied to a differential equation with a first integral of the form  $I(q, v) = v^T(Cq + c)$ , then every truncation of the modified differential equation has  $I(q, v)$  as a first integral.

The proofs are based on an induction argument. Since they are all very similar (see Hairer *et al.* (2002, Chapter IX)), we only present the proof of Theorem 4.3, for the case where the Hamiltonian  $H(p, q)$  is defined on a simply connected domain. This proof was first given by Benettin and Giorgilli (1994) and Tang (1994), and its ideas can be traced back to Moser (1968).

*Proof.* With  $y = (p, q)$ , the Hamiltonian system (2.7) is written more compactly as  $\dot{y} = J^{-1}\nabla H(y)$  with  $J$  of (2.9). We will show that all the coefficient functions of the modified equation can be written as

$$F_j(y) = J^{-1}\nabla H_j(y). \tag{4.11}$$

Assume, by induction, that (4.11) holds for  $j = 1, 2, \dots, N$  (this is satisfied for  $N = 1$ , because  $F_1(y) = F(y) = J^{-1}\nabla H(y)$ ). We have to prove the existence of a Hamiltonian  $H_{N+1}(y)$ . The idea is to consider the truncated modified equation (4.4), which is then a Hamiltonian system with Hamiltonian  $H(y) + hH_2(y) + \dots + h^{N-1}H_N(y)$ . Its flow  $\varphi_{N,t}(y)$ , compared to that of (4.2) and thus to the one-step map  $\Phi_h$  of the Störmer–Verlet method, satisfies

$$\Phi_h(y) = \varphi_{N,h}(y) + h^{N+1}F_{N+1}(y) + \mathcal{O}(h^{N+2}),$$

and also

$$\Phi'_h(y) = \varphi'_{N,h}(y) + h^{N+1}F'_{N+1}(y) + \mathcal{O}(h^{N+2}).$$

By Theorem 2.2 and by the induction hypothesis,  $\Phi_h$  and  $\varphi_{N,h}$  are symplectic transformations. This, together with  $\varphi'_{N,h}(y) = I + \mathcal{O}(h)$ , therefore implies

$$J = \Phi'_h(y)^T J \Phi'_h(y) = J + h^{N+1}(F'_{N+1}(y)^T J + JF'_{N+1}(y)) + \mathcal{O}(h^{N+2}).$$

Consequently, the matrix  $JF'_{N+1}(y)$  is symmetric. The function  $JF'_{N+1}(y)$  is therefore the gradient of some scalar function  $H_{N+1}(y)$ , which proves (4.11) for  $j = N + 1$ .  $\square$

The last argument of this proof requires that the domain be simply connected. For general domains, we have to use the representation (2.17) with the help of the generating function (2.18). We refer to Section IX.3.2 of Hairer *et al.* (2002) for details of the proof.

#### 4.4. Exponentially small error estimates

Theorem 4.1 proves a statement that is valid for all  $N \geq 1$ , and it is natural to ask which choice of  $N$  gives the best estimate.

**Example 4.6.** Consider the simple differential equation  $\ddot{q} = f(t)$  (which becomes autonomous after adding  $\dot{t} = 0$ ). If we try to compute the modified equation for the Störmer–Verlet method, we are readily convinced that its  $q$ -component is of the form

$$\ddot{q}(t) = f(t) + h^2 b_2 \ddot{f}(t) + h^4 b_4 f^{(4)}(t) + h^6 b_6 f^{(6)}(t) + \dots \quad (4.12)$$

Putting  $f(t) = e^t$ , the solution of this modified equation is

$$\tilde{q}(t) = C_1 + tC_2 + (1 + b_2 h^2 + b_4 h^4 + b_6 h^6 + \dots) e^t,$$

and inserted into (1.2) we obtain

$$(1 + b_2 h^2 + b_4 h^4 + b_6 h^6 + \dots)(e^{-h} - 2 + e^h) = h^2. \quad (4.13)$$

This shows that  $1 + b_2 h^2 + b_4 h^4 + \dots$  is analytic in a disc of radius  $2\pi$  centred at the origin. Consequently, the coefficients behave like  $b_{2k} \approx \text{Const} (2\pi)^{-2k}$  for  $k \rightarrow \infty$ .

Now consider functions  $f(t)$  whose derivatives grow like  $f^{(k)}(t) \approx k! M R^{-k}$ . This is the case for analytic  $f(t)$  with finite poles. The individual terms of the modified equation (4.12) then behave like

$$h^{2k} b_{2k} f^{(2k)}(t) \approx \text{Const} \frac{h^{2k} (2k)!}{(R \cdot 2\pi)^{2k}} \approx \text{Const} \sqrt{4\pi k} \left( \frac{h \cdot 2k}{R \cdot 2\pi e} \right)^{2k} \quad (4.14)$$

(using Stirling's formula). Even for very small step sizes  $h$  this expression is unbounded for  $k \rightarrow \infty$ , so that the series (4.12) cannot converge. However, formula (4.14) tells us that the terms of the series decrease until  $2k$  approaches the value  $2\pi R/h$ , and then they tend rapidly to  $\infty$ . It is therefore natural to truncate the modified equation after  $N$  terms, where  $N \approx 2\pi R/h$ .

To find a reasonably good truncation index  $N$  for general differential equations, we have to know estimates for all derivatives of  $F(y)$  and of the coefficient functions  $D_j(y)$  of the Taylor expansion of the numerical flow. One convenient way of doing this is to assume analyticity of these functions.

Exponentially small error bounds were first derived by Benettin and Giorgilli (1994). The following estimates are from Hairer *et al.* (2002, p. 306).

**Theorem 4.7.** Let  $F(y)$  be analytic in  $B_{2R}(y_0)$ , let the coefficients  $D_j(y)$  of the method (4.3) be analytic in  $B_R(y_0)$ , and assume that

$$\|F(y)\| \leq M \quad \text{and} \quad \|D_j(y)\| \leq \mu M \left( \frac{2\kappa M}{R} \right)^{j-1} \quad (4.15)$$

hold for  $y \in B_{2R}(y_0)$  and  $y \in B_R(y_0)$ , respectively. If  $h \leq h_0/4$  with

$h_0 = R/(e\eta M)$  and  $\eta = 2 \max(\kappa, \mu/(2 \ln 2 - 1))$ , then there exists  $N = N(h)$  (namely  $N$  equal to the largest integer satisfying  $hN \leq h_0$ ) such that the difference between the numerical solution  $y_1 = \Phi_h(y_0)$  and the exact solution  $\tilde{\varphi}_{N,t}(y_0)$  of the truncated modified equation (4.4) satisfies

$$\|\Phi_h(y_0) - \tilde{\varphi}_{N,h}(y_0)\| \leq h\gamma M e^{-h_0/h},$$

where  $\gamma = e(2 + 1.65\eta + \mu)$  depends only on the method.

The proof of this theorem is technical and long; see Hairer *et al.* (2002, Section IX.7) for details. We just explain how the assumptions can be checked for the Störmer–Verlet method (4.7).

We let  $y = (q, v)^T$ ,  $F(y) = (v, f(q))^T$ , and we consider the scaled norm  $\|y\| = \|q\| + h\|v\|$ . The quantities  $R$  and  $M$  are then given by the problem. The computation of the beginning of Section 4.2 shows that the functions  $D_j(q, v)$  are composed of derivatives of  $f(q)$  so that they are analytic on the same domain as  $f(q)$ . To find the constants  $\mu$  and  $\kappa$  in (4.15), we use  $\|F(y)\| = \|v\| + h\|f(q)\| \leq M$  for  $\|q - q_0\| + h\|v - v_0\| \leq 2R$ , and we estimate

$$\left\| \Phi_h(q, v) - \begin{pmatrix} q + hv \\ v + \frac{h}{2} f(q) \end{pmatrix} \right\| = \left\| \begin{pmatrix} \frac{h^2}{2} f(q) \\ \frac{h}{2} f(q + hv + \frac{h^2}{2} f(q)) \end{pmatrix} \right\| \leq hM$$

for  $\|q - q_0\| + h\|v - v_0\| \leq R$  and for  $hM \leq R$ . This follows from the fact that the argument of  $f$  satisfies  $\|q + hv + \frac{h^2}{2} f(q) - q_0\| \leq R + hM \leq 2R$ . Considered as a function of  $h$ ,  $\Phi_h(q, v)$  is analytic in the complex disc  $|h| \leq R/M$ . Cauchy’s estimate therefore yields

$$\|D_j(q, v)\| = \frac{1}{j!} \left\| \frac{d^j}{dh^j} \left( \Phi_h(q, v) - \begin{pmatrix} q + hv \\ v + \frac{h}{2} f(q) \end{pmatrix} \right) \Big|_{h=0} \right\| \leq M \left( \frac{M}{R} \right)^{j-1}$$

for  $j \geq 2$ . This proves the estimates (4.15) with  $\mu = 1$  and  $\kappa = 1/2$ .

### 5. Long-time behaviour of numerical solutions

In this section we show how the geometric properties of Section 2 turn into favourable long-term behaviour. Most of the results are obtained with the help of backward error analysis.

#### 5.1. Energy conservation

We have seen in Example 3.4 that the total energy  $H(p, q)$  of a Hamiltonian system is not preserved exactly by the Störmer–Verlet method. In that example it is, however, approximately preserved. Also for the Kepler problem, Figure 3.1 indicates no drift in the energy. As the following theorem shows, the Hamiltonian is in fact approximately preserved over very long times for general Hamiltonian systems.

**Theorem 5.1.** The total energy along a numerical solution  $(p_n, q_n)$  of the Störmer–Verlet method satisfies

$$|H(p_n, q_n) - H(p_0, q_0)| \leq Ch^2 + C_N h^N t \quad \text{for } 0 \leq t = nh \leq h^{-N}$$

for arbitrary positive integer  $N$ . The constants  $C$  and  $C_N$  are independent of  $t$  and  $h$ .  $C_N$  depends on bounds of derivatives of  $H$  up to  $(N + 1)$ th order in a region that contains the numerical solution values  $(p_n, q_n)$ .

We give two different proofs of this result, the first one based on the symplecticity, the second one on the symmetry of the method. When the Hamiltonian is analytic, both proofs can be refined to yield an estimate  $Ch^2 + C_0 e^{-c/h} t$  over exponentially long times  $t \leq e^{c/h}$ , with  $c$  proportional to  $1/\Omega$ , where  $\Omega$  is an upper bound of  $\|M^{-1/2} \nabla^2 U(q) M^{-1/2}\|^{1/2}$ , *i.e.*, of the highest frequency in the linearized system.

The *first proof* uses the symplecticity of the Störmer–Verlet method via backward error analysis, in an argument due to Benettin and Giorgilli (1994). It applies to general symplectic methods for general (smooth) Hamiltonian systems (2.7). We know from Theorem 4.3 that the modified differential equation, truncated after  $N$  terms, is again Hamiltonian, with a modified Hamiltonian  $\tilde{H}$  that is  $\mathcal{O}(h^2)$  close to the original Hamiltonian  $H$  in a neighbourhood of the numerical solution values. Consider now  $\tilde{H}$  along the numerical solution. We write the deviation of  $\tilde{H}$  as a telescoping sum

$$\tilde{H}(p_n, q_n) - \tilde{H}(p_0, q_0) = \sum_{j=0}^{n-1} (\tilde{H}(p_{j+1}, q_{j+1}) - \tilde{H}(p_j, q_j)).$$

By construction of the modified equation, we have for its flow  $\tilde{\varphi}_h(p_j, q_j) = (p_{j+1}, q_{j+1}) + \mathcal{O}(h^{N+1})$ . On the other hand, the flow  $\tilde{\varphi}_t$  preserves the modified Hamiltonian, and hence

$$\tilde{H}(p_{j+1}, q_{j+1}) - \tilde{H}(p_j, q_j) = \tilde{H}(p_{j+1}, q_{j+1}) - \tilde{H}(\tilde{\varphi}_h(p_j, q_j)) = \mathcal{O}(h^{N+1}).$$

Inserting this estimate in the above sum yields the result.

The *second proof* uses only the symmetry of the Störmer–Verlet method. It was given in Hairer and Lubich (2000b) because its arguments extend to numerical energy conservation in oscillatory systems when the product of the step size with the highest frequencies is bounded away from 0 (see Section 5.4). Backward error analysis, or the asymptotic  $h^2$ -expansion of the numerical solution, shows that there exists, for every  $n$ , a function  $q^n(t)$  with  $q^n(0) = q_n$  and  $q^n(-h) = q_{n-1} + \mathcal{O}(h^{N+1})$  satisfying

$$q^n(t+h) - 2q^n(t) + q^n(t-h) = h^2 f(q^n(t)) + \mathcal{O}(h^{N+2}) \quad (5.1)$$

for  $t$  in some fixed interval around 0. The functions  $q^n(t+h)$  and  $q^{n+1}(t)$  agree up to  $\mathcal{O}(h^{N+1})$ , as do their  $k$ th derivatives multiplied with  $h^k$ ,

for  $k \leq N$ . By Taylor expansion in (5.1),

$$\sum_{l=1}^{N/2} \frac{2}{(2l)!} \frac{d^{2l}q^n}{dt^{2l}}(t) h^{2l-2} = f(q^n(t)) + \mathcal{O}(h^N). \quad (5.2)$$

Because of the symmetry of the method, only even-order derivatives of  $q^n(t)$  (and even powers of the step size) are present in (5.2).

We multiply (5.2) with  $\dot{q}^n(t)^T M$  and integrate over  $t$ . The key observation is now that the product of  $\dot{q}^n(t)$  with an *even*-order derivative of  $q^n(t)$  is a total differential (we omit the superscript  $n$  in the following formula):

$$\dot{q}^T M q^{(2l)} = \frac{d}{dt} \mathcal{A}_l[q]$$

with

$$\mathcal{A}_l[q] = \left( \dot{q}^T M q^{(2l-1)} - \ddot{q}^T M q^{(2l-2)} + \dots \mp (q^{(l-1)})^T M q^{(l+1)} \pm \frac{1}{2} (q^{(l)})^T M q^{(l)} \right).$$

In particular,  $\mathcal{A}_1[q] = \frac{1}{2} \dot{q}^T M \dot{q}$ . Moreover, for  $f(q) = -M^{-1} \nabla U(q)$  we clearly have  $\dot{q}^T M f(q) = -(\mathrm{d}/\mathrm{d}t)U(q)$ . For the energy functional

$$\mathcal{H}[q](t) = \sum_{l=1}^{N/2} \frac{2}{(2l)!} \mathcal{A}_l[q](t) h^{2l-2} + U(q(t))$$

we thus obtain  $(\mathrm{d}/\mathrm{d}t)\mathcal{H}[q^n](t) = \mathcal{O}(h^N)$ , and hence

$$\mathcal{H}[q^n](h) - \mathcal{H}[q^n](0) = \mathcal{O}(h^{N+1}). \quad (5.3)$$

Since the functions  $q^n(t+h)$  and  $q^{n+1}(t)$ , together with their  $k$ th derivatives scaled by  $h^k$  ( $k \leq N$ ), are equal up to  $\mathcal{O}(h^{N+1})$ , we further have

$$\mathcal{H}[q^{n+1}](0) - \mathcal{H}[q^n](h) = \mathcal{O}(h^{N+1}). \quad (5.4)$$

Moreover, with  $p^n(t) = M\dot{q}^n(t)$  we have

$$\mathcal{H}[q^n](0) = H(p^n(0), q^n(0)) + \mathcal{O}(h^2) = H(p_n, q_n) + \mathcal{O}(h^2), \quad (5.5)$$

where the last equation follows by noting

$$p_n = M \frac{q_{n+1} - q_{n-1}}{2h} = p^n(0) + \mathcal{O}(h^2).$$

Hence, from (5.3)–(5.5),

$$\begin{aligned} H(p_n, q_n) - H(p_0, q_0) &= \mathcal{H}[q^n](0) - \mathcal{H}[q^0](0) + \mathcal{O}(h^2) \\ &= \mathcal{O}(nh^{N+1}) + \mathcal{O}(h^2), \end{aligned}$$

which completes the proof.

### 5.2. Linear error growth for integrable systems

General Hamiltonian systems may have extremely complicated dynamics, and little can be said about the long-time behaviour of their discretizations apart from the long-time near-conservation of the total energy considered above. At the other end, the simplest conceivable dynamics – uniform motion on a Cartesian product of circles – appears in integrable Hamiltonian systems. Their practical interest lies in the fact that many physical systems are perturbations of integrable systems, with planetary motion as the classical example and historical driving force.

A Hamiltonian system (2.7) is *integrable* if there exists a symplectic transformation

$$(p, q) = \psi(a, \theta) \tag{5.6}$$

to *action-angle variables*  $(a, \theta)$ , defined for actions  $a = (a_1, \dots, a_d)$  in some open set of  $\mathbb{R}^d$  and for angles  $\theta$  on the whole  $d$ -dimensional torus

$$\mathbb{T}^d = \mathbb{R}^d / (2\pi\mathbb{Z}^d) = \{(\theta_1, \dots, \theta_d); \theta_i \in \mathbb{R} \bmod 2\pi\},$$

such that the Hamiltonian in these variables depends only on the actions

$$H(p, q) = H(\psi(a, \theta)) = K(a). \tag{5.7}$$

In the action-angle variables, the equations of motion are simply

$$\dot{a} = 0, \quad \dot{\theta} = \omega(a), \tag{5.8}$$

with the *frequencies*  $\omega = (\omega_1, \dots, \omega_d)^T = \nabla_a K$  (note that  $\nabla_\theta K = 0$ ). This has a quasi-periodic (or possibly periodic) flow:

$$\varphi_t(a, \theta) = (a, \theta + \omega(a)t). \tag{5.9}$$

For every  $a$ , the torus  $\{(a, \theta) : \theta \in \mathbb{T}^d\}$  is thus invariant under the flow. We express the actions and angles in terms of the original variables  $(p, q)$  via the inverse transform of (5.6) as

$$(a, \theta) = (I(p, q), \Theta(p, q)),$$

and note that the components of  $I = (I_1, \dots, I_d)$  are first integrals of the integrable system.

Integrability of a Hamiltonian system is an exceptional property: the system has  $d$  independent first integrals  $I_1, \dots, I_d$  whose Poisson brackets vanish pairwise, that is,

$$\{I_i, I_j\} = \nabla_q I_i^T \nabla_p I_j - \nabla_p I_i^T \nabla_q I_j = 0 \quad \text{for all } i, j.$$

The solution trajectories of the Hamiltonian systems with Hamiltonian  $I_i$  exist for all time (in the action-angle variables, their flow is simply  $\varphi_t^{[i]}(a, \theta) =$

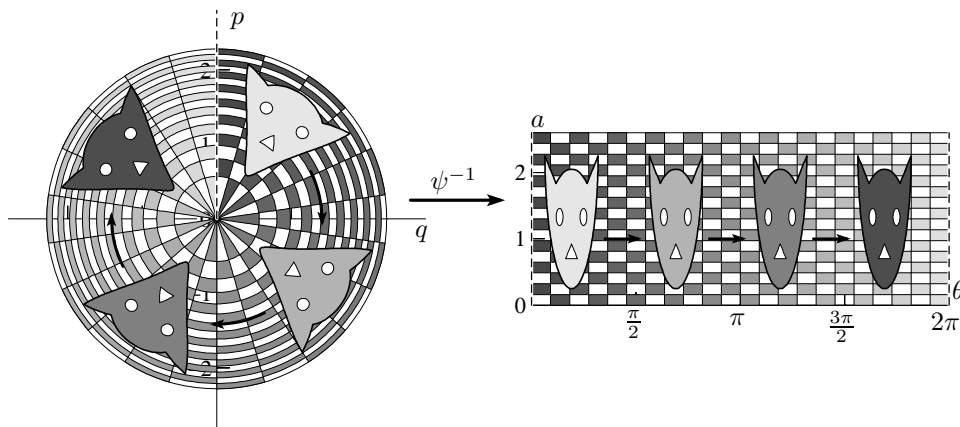


Figure 5.1. Transformation to action-angle variables.

$(a, \theta + te_i)$  with  $e_i$  denoting the  $i$ th unit vector of  $\mathbb{R}^d$ , and the level sets of  $I$  are compact (the invariant tori  $\{a = \text{Const } \theta \in \mathbb{T}^d\}$ ). Conversely, the *Arnold–Liouville theorem* (Arnold 1963) states that every Hamiltonian system having  $d$  first integrals with the above properties can be transformed to action-angle variables with a Hamiltonian depending only on the actions.

**Example 5.2.** The harmonic oscillator  $H(p, q) = \frac{1}{2}p^2 + \frac{1}{2}q^2$  is integrable, with the transformation to action-angle coordinates given by

$$\begin{pmatrix} p \\ q \end{pmatrix} = \begin{pmatrix} \sqrt{2a} \cos \theta \\ \sqrt{2a} \sin \theta \end{pmatrix},$$

with  $a = H(p, q)$ : see Figure 5.1. Here, the action-angle coordinates are symplectic polar coordinates.

**Example 5.3.** The Kepler problem, with  $H(p, q) = \frac{1}{2}(p_1^2 + p_2^2) - (q_1^2 + q_2^2)^{-\frac{1}{2}}$  in the range  $H < 0$ , is integrable with actions  $a_1 = 1/\sqrt{-2H}$  and  $a_2 = L$  (the angular momentum,  $L = q_1p_2 - q_2p_1$ ). The frequencies are  $\omega_1 = \omega_2 = 2\pi/T$ , where  $T = 2\pi/(-2H)^{3/2}$  is the period of a trajectory with total energy  $H$ .

**Example 5.4.** A further celebrated example of an integrable system is the Toda lattice (Toda 1970, Flaschka 1974), which describes a system of particles on a line interacting with exponential forces. The Hamiltonian is

$$H(p, q) = \sum_{k=1}^d \left( \frac{1}{2}p_k^2 + \exp(q_k - q_{k+1}) \right)$$

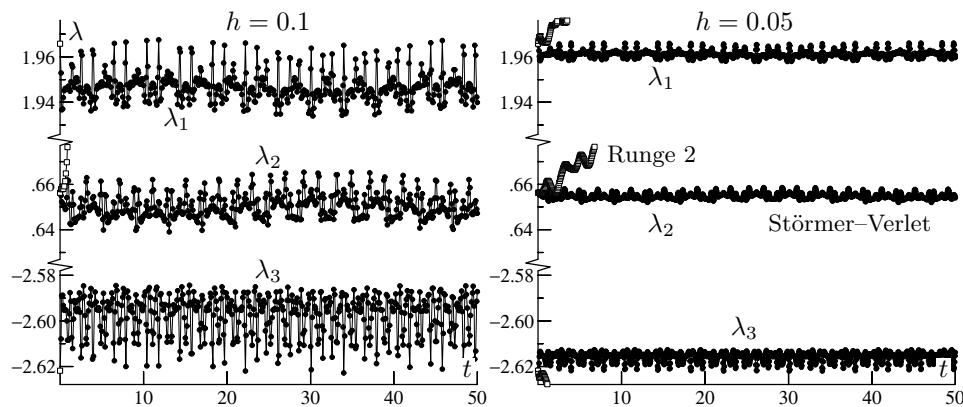


Figure 5.2. Toda eigenvalues along the numerical solution.

with periodic extension  $q_{d+1} = q_1$ . The eigenvalues of the matrix

$$L = \begin{pmatrix} a_1 & b_1 & & & b_d \\ b_1 & a_2 & b_2 & 0 & \\ & b_2 & \ddots & \ddots & \\ & 0 & \ddots & a_{d-1} & b_{d-1} \\ b_d & & & b_{d-1} & a_d \end{pmatrix}, \quad \begin{aligned} a_k &= -\frac{1}{2}p_k, \\ b_k &= \frac{1}{2} \exp\left(\frac{1}{2}(q_k - q_{k+1})\right) \end{aligned}$$

are first integrals whose Poisson brackets vanish pairwise.

We consider the case  $d = 3$  and choose initial values  $q_0 = (1, 2, -1)^T$  and  $p_0 = (-1.5, 1, 0.5)^T$ . Figure 5.2 shows the eigenvalues of  $L$  along the numerical solution of the Störmer–Verlet and the second-order Runge method obtained with step sizes  $h = 0.1$  (left) and  $h = 0.05$  (right) on the interval  $0 \leq t \leq 50$ . Not only the Hamiltonian (Theorem 5.1), but all  $d$  first integrals of the integrable system are well approximated over long times with an error of size  $\mathcal{O}(h^2)$ . This is explained by Theorem 5.5 below.

The global error in  $(p, q)$  is plotted in Figure 5.3. We observe a linear error growth for the Störmer–Verlet method, in contrast to a quadratic error growth for the second-order Runge method.

The study of the error behaviour of the numerical method combines *backward error analysis*, by means of which the numerical map is interpreted as being essentially the time- $h$  flow of a modified Hamiltonian system, and the *perturbation theory* of integrable systems, a rich mathematical theory originally developed for problems of celestial mechanics (Poincaré 1892/1893/1899, Siegel and Moser 1971). The effect of a small perturbation



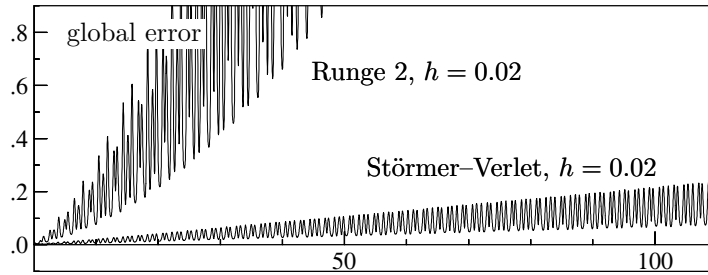


Figure 5.3. Global error of the Störmer-Verlet and the second-order Runge method on the Toda lattice.

of an integrable system is well under control in subsets of the phase space where the frequencies  $\omega$  satisfy Siegel’s *diophantine condition*:

$$|k \cdot \omega| \geq \gamma |k|^{-\nu} \quad \text{for all } k \in \mathbb{Z}^d \tag{5.10}$$

for some positive constants  $\gamma$  and  $\nu$ , with  $|k| = \sum_i k_i$ . For  $\nu > d - 1$ , almost all frequencies (in the sense of the Lebesgue measure) satisfy (5.10) for some  $\gamma > 0$ . For any choice of  $\gamma$  and  $\nu$  the complementary set is, however, open and dense in  $\mathbb{R}^d$ .

For general numerical integrators applied to integrable systems (or perturbations thereof) the error grows quadratically with time, and there is a linear drift away from the first integrals  $I_i$ . For symplectic methods such as the Störmer-Verlet method there is linear error growth and long-time near-preservation of the first integrals  $I_i$ , as is shown by the following result from Hairer *et al.* (2002, Section X.3).

**Theorem 5.5.** Consider applying the Störmer-Verlet method to an integrable system (2.7) with real-analytic Hamiltonian. Suppose that  $\omega^* \in \mathbb{R}^d$  satisfies the diophantine condition (5.10). Then there exist positive constants  $C, c$  and  $h_0$  such that the following holds for all step sizes  $h \leq h_0$ : every numerical solution  $(p_n, q_n)$  starting with frequencies  $\omega_0 = \omega(I(p_0, q_0))$  such that  $\|\omega_0 - \omega^*\| \leq c |\log h|^{-\nu-1}$ , satisfies

$$\begin{aligned} \|(p_n, q_n) - (p(t), q(t))\| &\leq C t h^2, \\ \|I(p_n, q_n) - I(p_0, q_0)\| &\leq C h^2, \end{aligned} \quad \text{for } t = nh \leq h^{-2}.$$

The constants  $h_0, c, C$  depend on  $d, \gamma, \nu$  and on bounds of the Hamiltonian.

The basic steps of the proof are summarized in Figure 5.4. By backward error analysis, the numerical method coincides, up to arbitrary order in  $h$ , with the flow of the modified differential equation, which is a Hamiltonian

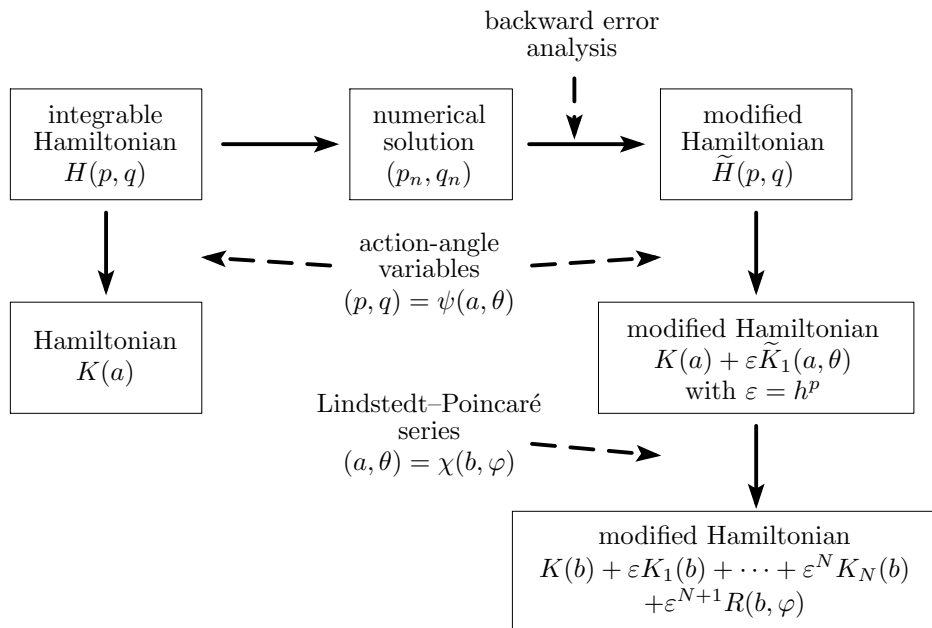


Figure 5.4. Transformations in the proof of Theorem 5.5.

perturbation of size  $\varepsilon = h^2$  of the original, integrable system. We are thus in the realm of classical perturbation theory. In addition to the transformation to the action-angle variables  $(a, \theta)$ , which gives the modified Hamiltonian in the form  $K(a) + \varepsilon \tilde{K}_1(a, \theta)$ , we use a further symplectic coordinate transformation  $(a, \theta) = \chi(b, \varphi)$  which eliminates, up to high-order terms in  $\varepsilon$ , the dependence on the angles in the modified Hamiltonian. This transformation is  $\mathcal{O}(\varepsilon)$  close to the identity. It is constructed as

$$b = a - \nabla_{\theta} S(b, \theta), \quad \varphi = \theta + \nabla_b S(b, \theta),$$

where the generating function  $S(b, \theta)$  is given by a Lindstedt–Poincaré perturbation series,

$$S(b, \theta) = \varepsilon S_1(b, \theta) + \varepsilon^2 S_2(b, \theta) + \dots + \varepsilon^N S_N(b, \theta).$$

The error propagation is then studied in the  $(b, \varphi)$ -variables, with the result

$$\begin{aligned} \|b(t) - b_0\| &\leq C t \varepsilon^{N+1}, \\ \|\varphi(t) - \varphi_0 - \omega_{\varepsilon}(b_0) t\| &\leq C (t + t^2) \varepsilon^N, \end{aligned} \quad \text{for } t^2 \leq 1/\varepsilon^N,$$

with  $\omega_{\varepsilon}(b) = \omega(b) + \mathcal{O}(\varepsilon)$ . Transforming back to the original variables  $(p, q)$  finally yields the stated result.

Theorem 5.5 admits extensions in several directions. It is just one of a series of results on the long-time behaviour of geometric integrators.

- The theorem does not apply directly to the Kepler problem, which has two identical frequencies  $\omega_1 = \omega_2 = (-2H)^{3/2}$ . However, since the angular momentum  $a_2 = L$  is preserved exactly by the Störmer–Verlet method, it turns out that the modified Hamiltonian written in the action-angle variables of the Kepler problem is independent of the angle  $\theta_2$ . Only the angle  $\theta_1$  must therefore be eliminated via the perturbation series, and this involves only the single frequency  $\omega_1$  for which the diophantine condition is trivially satisfied. The proof and result of Theorem 5.5 thus extend to the Kepler problem.
- The linear error growth remains intact when the method is applied to *perturbed integrable systems*  $H(p, q) + \varepsilon G(p, q)$  with a perturbation parameter of size  $\varepsilon = \mathcal{O}(h^\alpha)$  for some positive exponent  $\alpha$ .
- Under stronger conditions on the initial values or on the system, the near-preservation of the action variables along the numerical solution holds for times that are exponentially long in a negative power of the step size (Hairer and Lubich 1997, Moan 2002). For a Cantor set of initial values and a Cantor set of step sizes this holds even perpetually, in view of the existence of invariant tori of the numerical integrator close to the invariant tori of the integrable system (Shang 1999, 2000).
- Perturbed integrable systems have KAM tori, *i.e.*, deformations of the invariant tori of the integrable system corresponding to diophantine frequencies  $\omega$ , which are invariant under the flow of the perturbed system. If the method is applied to such a perturbed integrable system, then the numerical method has tori which are near-invariant over exponentially long times (Hairer and Lubich 1997). For a Cantor set of non-resonant step sizes there are even truly invariant tori on which the numerical one-step map reduces to rotation by  $h\omega$  in suitable coordinates (Hairer *et al.* 2002, p. 371).
- There is a completely analogous theory for *integrable reversible systems* (Hairer *et al.* 2002, Chapter XI). These are differential equations with reversible flow (2.2), which are transformed to the form (5.8) by a transformation  $(q, v) = (\mu(a, \theta), \nu(a, \theta))$  that preserves reversibility, *i.e.*,  $\mu$  is odd in  $\theta$  and  $\nu$  is even in  $\theta$ . In that theory, only the reversibility of the numerical method comes into play, not the symplecticity. There is again linear error growth, long-time near-preservation of the action variables, and an abundance of invariant tori.
- For dissipatively perturbed integrable systems, where only one torus survives the perturbation and becomes weakly attractive, the existence of a nearby invariant torus of the numerical method is shown under weak assumptions on the step size in Stoffer (1998) and Hairer and Lubich (1999).

### 5.3. Statistical behaviour

The equation of motion of a system of 864 particles interacting through a Lennard-Jones potential has been integrated for various values of the temperature and density, relative, generally, to a fluid state. The equilibrium properties have been calculated and are shown to agree very well with the corresponding properties of argon. L. Verlet (1967)

In molecular dynamics, it is the computation of statistical or thermodynamic quantities, such as temperature, which is of interest, rather than single trajectories. The success of the Störmer–Verlet method in this field lies in the observation that the method is apparently able to reproduce the correct statistical behaviour over long times. Since Verlet (1967), this has been confirmed in countless computational experiments. Backward error analysis gives indications as to why this might be so, but to our knowledge there are as yet no rigorous mathematical results in the literature explaining the favourable statistical behaviour.

In the following we derive a result which is a discrete analogue of the virial theorem of statistical mechanics; *cf.* Abraham and Marsden (1978, p. 243) and Gallavotti (1999, p. 129). It comes as a consequence of the long-time near-conservation of energy. Consider the Poisson bracket  $\{F, H\} = \nabla_q F^T \nabla_p H - \nabla_p F^T \nabla_q H$  of an arbitrary differentiable function  $F(p, q)$  with the Hamiltonian. Along every solution  $(p(t), q(t))$  of the Hamiltonian system we have

$$\{F, H\}(p(t), q(t)) = \frac{d}{dt} F(p(t), q(t)),$$

and hence the time average of the Poisson bracket along a solution is

$$\frac{1}{T} \int_0^T \{F, H\}(p(t), q(t)) dt = \frac{1}{T} (F(p(T), q(T)) - F(p(0), q(0))).$$

If  $F$  is bounded along the solution, this shows that the average is of size  $\mathcal{O}(1/T)$  as  $T \rightarrow \infty$ . In particular, this condition is satisfied if the energy level set  $\{(p, q) : H(p, q) = H(p(0), q(0))\}$  is compact.

**Example 5.6.** For a separable Hamiltonian (2.8) the choice  $F(p, q) = p^T q$  yields the virial theorem of Clausius (Gallavotti 1999, p. 129),

$$\lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T p(t)^T M^{-1} p(t) dt = \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T q(t)^T \nabla U(q(t)) dt,$$

*i.e.*, the time average of twice the kinetic energy equals that of the *virial function*  $q^T \nabla U(q)$ .

For the numerical discretization there is the following result.

**Theorem 5.7.** Let  $H(p, q)$  be a real-analytic Hamiltonian for which

$$K_\delta = \{(p, q) : |H(p, q) - H_0| \leq \delta\} \text{ is compact}$$

for some  $\delta > 0$ . Let  $F(p, q)$  be any smooth real-valued function, bounded by  $\mu$  on  $K_\delta$ . Then the numerical solution  $(p_n, q_n)$  obtained by the Störmer–Verlet method satisfies

$$\left| \frac{1}{N} \sum'_{n=0} \{F, H\}(p_n, q_n) \right| \leq \frac{2\mu}{Nh} + Ch^2 \quad \text{for } Nh \leq e^{c/h} \text{ and } h \leq h_0, \tag{5.11}$$

where the prime on the sum indicates that the first and last term are taken with weight  $\frac{1}{2}$ . The constants  $C, c, h_0 > 0$  depend on bounds of  $H$  on a complex neighbourhood of  $K_\delta$  and on bounds of the first three derivatives of  $F$  on  $K_\delta$ , but they are independent of  $h$  and  $(p_0, q_0) \in K_{\delta/2}$ .

In particular, the left-hand side of (5.11) is  $\mathcal{O}(h^2)$  for  $h^{-2} \leq Nh \leq e^{c/h}$ .

*Proof.* By Theorem 5.1 and the remark thereafter, we know that

$$y_n := (p_n, q_n) \in K_\delta \quad \text{for } nh \leq e^{c/h}.$$

Since  $y_{n+1} = \varphi_h(y_n) + \mathcal{O}(h^3)$  and  $\{F, H\}(\varphi_t(y)) = \frac{d}{dt} F(\varphi_t(y))$ , we have

$$\begin{aligned} \frac{h}{2} \{F, H\}(y_n) + \frac{h}{2} \{F, H\}(y_{n+1}) &= \int_0^h \{F, H\}(\varphi_t(y_n)) dt + \mathcal{O}(h^3) \\ &= F(y_{n+1}) - F(y_n) + \mathcal{O}(h^3). \end{aligned}$$

Hence

$$\frac{1}{N} \sum'_{n=0} \{F, H\}(y_n) = \frac{1}{Nh} (F(y_N) - F(y_0)) + \mathcal{O}(h^2),$$

which yields the stated estimate. □

**Example 5.8.** We give a numerical experiment with a small-scale version of Verlet’s argon model. It considers  $N_A$  atoms interacting by the Lennard–Jones potential

$$V(r) = 4\varepsilon \left( \left( \frac{\sigma}{r} \right)^{12} - \left( \frac{\sigma}{r} \right)^6 \right).$$

The Hamiltonian of the system is (3.2) with  $V_{ij} = V$ . We choose  $N_A = 7$

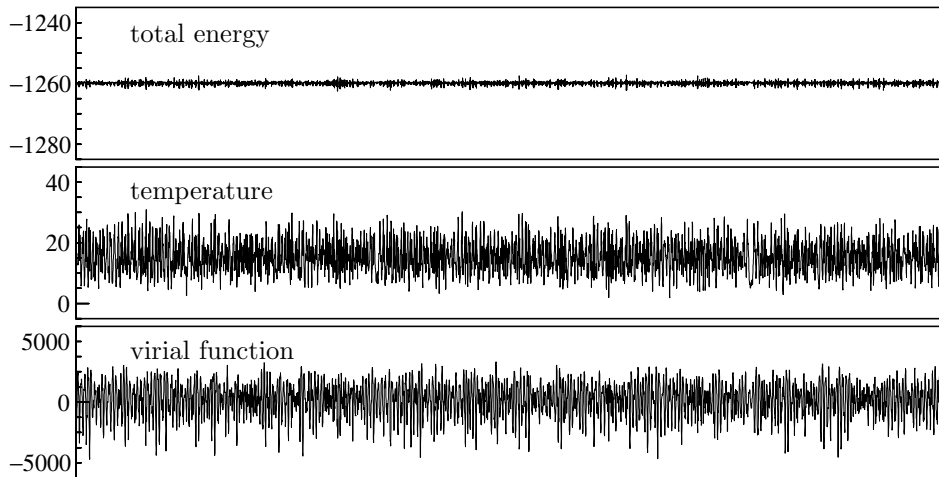


Figure 5.5. Computed total energy, temperature and virial function of the argon crystal, 10 000 steps of size  $h = 40$  [fsec].

and the data of Biesiadecki and Skeel (1993); see also Hairer *et al.* (2002, p. 15). Figure 5.5 shows the Hamiltonian, the temperature

$$T(p) = \frac{1}{N_A k_B} \frac{1}{2m} \sum_{i=1}^{N_A} \|p_i\|^2$$

( $k_B$  is Boltzmann's constant), and the virial function

$$C(q) = \sum_{i=2}^{N_A} \sum_{j=1}^{i-1} V'(r_{ij}) r_{ij}$$

(with  $r_{ij} = \|q_i - q_j\|$ ) over an interval of length  $4 \cdot 10^5$  [fsec], obtained by the Störmer–Verlet method with step size  $h = 40$  [fsec]. The units in the figure are such that  $k_B = 1$ . The size of the oscillations in the Hamiltonian is proportional to  $h^2$ , whereas that in the temperature and in the virial function is independent of  $h$ . At the end of the integration (after 10 000 steps) the averages of twice the kinetic energy and of the virial function are 217.1 and 217.8, respectively.

#### 5.4. Oscillatory differential equations

Nonlinear mass-spring models have traditionally been very useful in explaining various phenomena of more complicated 'real' physical systems. We have already mentioned the Toda lattice. An equally famous problem is the Fermi–Pasta–Ulam model (Fermi, Pasta and Ulam 1955, Ford 1992), where

a nonlinear perturbation to a primarily linear problem is studied over long times. Here we use a variant of this problem for gaining insight into the long-time energy behaviour of the Störmer–Verlet method applied to oscillatory systems with multiple time scales. We are interested in using step sizes  $h$  for which the product with the highest frequency  $\omega$  in the system is bounded away from zero. (Values of  $h\omega \approx 1/2$  are routinely used in molecular dynamics.) In this situation, backward error analysis is no longer applicable, since the ‘exponentially small’ error terms are then of size  $\mathcal{O}(e^{-c/h\omega}) = \mathcal{O}(1)$ .

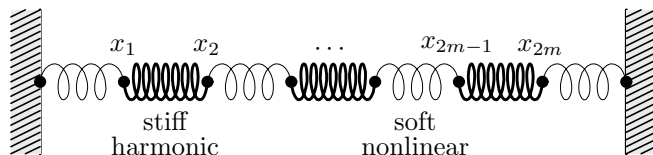


Figure 5.6. Chain with alternating soft nonlinear and stiff linear springs.

**Example 5.9.** Consider a chain of  $2m$  mass points, connected with alternating soft nonlinear and stiff linear springs, and fixed at the end points; see Galgani, Giorgilli, Martinoli and Vanzini (1992) and Figure 5.6. The variables  $x_1, \dots, x_{2m}$  stand for the displacements of the mass points. In terms of the new variables

$$q_i = (x_{2i} + x_{2i-1})/\sqrt{2}, \quad q_{m+i} = (x_{2i} - x_{2i-1})/\sqrt{2}$$

(which represent a scaled displacement and a scaled expansion/compression of the  $i$ th stiff spring) and the momenta  $p_i = \dot{q}_i$ , the motion is described by a Hamiltonian system with

$$H(p, q) = \frac{1}{2} \sum_{i=1}^{2m} p_i^2 + \frac{\omega^2}{2} \sum_{i=1}^m q_{m+i}^2 + \frac{1}{4} \left( (q_1 - q_{m+1})^4 + \sum_{i=1}^{m-1} (q_{i+1} - q_{m+i+1} - q_i - q_{m+i})^4 + (q_m + q_{2m})^4 \right),$$

where  $\omega \gg 1$  is a large parameter. Here we assume cubic nonlinear springs, but the special form of the nonlinearity is not important.

For an illustration we consider  $m = 3$  and choose  $\omega = 30$ . In Figure 5.7 we have plotted the following quantities as functions of time: the Hamiltonian  $H$  (actually we plot  $H - 0.8$  for graphical reasons), the oscillatory energy  $I$  defined as

$$I = I_1 + I_2 + I_3 \quad \text{with} \quad I_j = \frac{1}{2} p_{m+j}^2 + \frac{1}{2} \omega^2 q_{m+j}^2,$$

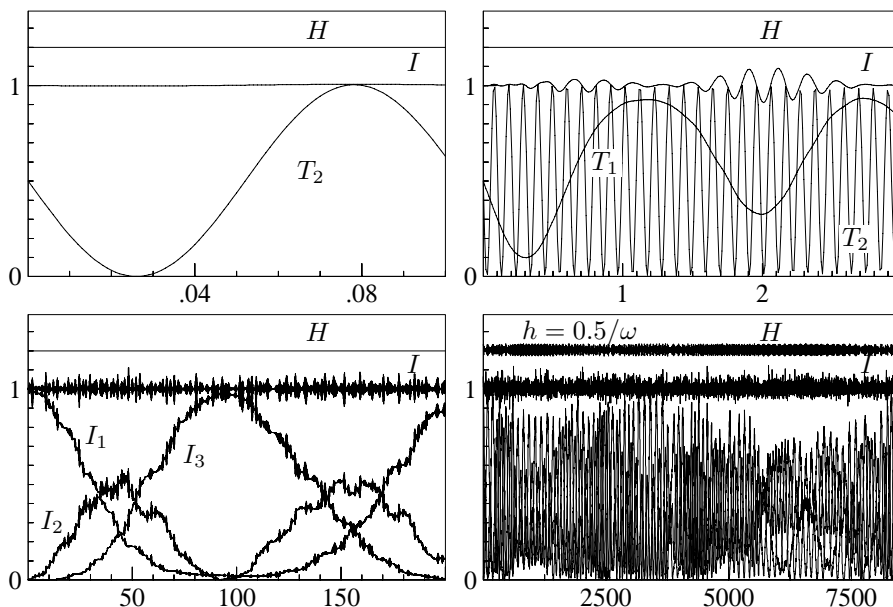


Figure 5.7. Different time scales in a Fermi–Pasta–Ulam problem, and energy conservation of the Störmer–Verlet method (last picture).

and the kinetic energies of the mass centre motion and of the relative motion of masses joined by a stiff spring,

$$T_1 = \frac{1}{2}(p_1^2 + p_2^2 + p_3^2), \quad T_2 = \frac{1}{2}(p_4^2 + p_5^2 + p_6^2).$$

The system has different dynamics on several time scales: on the fast scale  $\omega^{-1}$  the motion is nearly harmonic in the stiff springs, on scale  $\omega^0$  there is the motion of the soft springs driven by the nonlinearity, and on the slow scale  $\omega$  there is an energy exchange between the stiff linear springs. For the first three pictures the solutions were computed with high accuracy.

In the last picture we show the results obtained by the Störmer–Verlet method with step size  $h = 0.5/\omega$ . We note that both  $H$  and  $I$  are approximately conserved over long times. For fixed  $\omega$  the size of the oscillations in  $H$  is proportional to  $h^2$ . However, the oscillations remain of the same size if  $h$  decreases and  $\omega$  increases such that  $h\omega$  remains constant. The oscillations in  $I$  are of size  $\mathcal{O}(\omega^{-1})$  uniformly for  $h \rightarrow 0$ .

The equations of motion for the above example are of the form

$$\ddot{q} = -\Omega^2 q - \nabla U(q) \quad \text{with} \quad \Omega = \begin{pmatrix} 0 & 0 \\ 0 & \omega I \end{pmatrix}, \quad (5.12)$$

with a single high frequency  $\omega \gg 1$  and with a smooth potential  $U(q)$  whose derivatives are bounded independently of  $\omega$ . In addition to the total energy



as a conserved quantity,

$$H(p, q) = \frac{1}{2} p^T p + \frac{1}{2} q^T \Omega^2 q + U(q),$$

the system has an adiabatic invariant: over times exponentially long in  $\omega$ , the *oscillatory energy*

$$I(p, q) = \frac{1}{2} p^T \begin{pmatrix} 0 & 0 \\ 0 & I \end{pmatrix} p + \frac{\omega^2}{2} q^T \begin{pmatrix} 0 & 0 \\ 0 & I \end{pmatrix} q \tag{5.13}$$

is preserved up to  $\mathcal{O}(\omega^{-1})$ . This holds uniformly for all initial values for which the total energy is bounded by a constant independent of  $\omega$ , *i.e.*, for bounded  $(p, q)$  with  $(0 \ I)q = \mathcal{O}(\omega^{-1})$ .

Now consider applying the Störmer–Verlet method to such a system. The step size is then restricted to  $h\omega < 2$  for linear stability, as Example 3.4 shows. The Hamiltonian  $H(p_n, q_n)$  and the oscillatory energy  $I(p_n, q_n)$  of (5.13) oscillate rapidly, but stay within an  $\mathcal{O}((h\omega)^2)$  band over long times. The oscillations do not become smaller when  $h$  is decreased but  $\omega$  is increased such that their product  $h\omega$  is kept fixed. Nevertheless, the following result shows that the *time averages* of the total and oscillatory energies

$$\begin{aligned} \bar{H}_n &= \frac{h}{T} \sum_{|jh| \leq T/2} H(p_{n+j}, q_{n+j}), \\ \bar{I}_n &= \frac{h}{T} \sum_{|jh| \leq T/2} I(p_{n+j}, q_{n+j}), \end{aligned}$$

for an arbitrary fixed  $T > 0$ , remain constant up to  $\mathcal{O}(h)$  over long times even when  $h\omega$  is bounded away from zero, but within the range of linear stability.

**Theorem 5.10.** Let the Störmer–Verlet method be applied to the problem (5.12) with a step size  $h$  for which  $0 < c_0 \leq h\omega \leq c_1 < 2$ . Let  $\tilde{\omega}$  be defined by the relation  $\sin(\frac{1}{2}h\tilde{\omega}) = \frac{1}{2}h\omega$  and suppose  $|\sin(\frac{1}{2}kh\tilde{\omega})| \geq c\sqrt{h}$  for  $k = 1, \dots, N$  for some  $N \geq 2$  and  $c > 0$ . Suppose further that the total energy at the initial value  $(p_0, q_0)$  is bounded independently of  $\omega$ , and that the numerical solution values  $q_n$  stay in a region where all derivatives of the potential  $U$  are bounded. Then, the time averages of the total and the oscillatory energy along the numerical solution satisfy

$$\begin{aligned} \bar{H}_n &= \bar{H}_0 + \mathcal{O}(h), \\ \bar{I}_n &= \bar{I}_0 + \mathcal{O}(h), \end{aligned} \quad \text{for } 0 \leq nh \leq h^{-N+1}.$$

The constants symbolized by  $\mathcal{O}$  are independent of  $n, h, \omega$  with the above conditions.

It should, however, be noted that the time averages  $\overline{H}_n$  and  $\overline{I}_n$  do not, in general, remain  $\mathcal{O}(h)$  close to the initial values  $H(p_0, q_0)$  and  $I(p_0, q_0)$ .

The estimates of Theorem 5.10 can be improved to  $\mathcal{O}(h^2)$  if a weighted time average is taken, replacing the characteristic function of the interval  $[-T/2, T/2]$  by a smooth windowing function with bounded support, and if the oscillatory energy  $I$  is replaced by  $J(p, q) = I(p, q) + q^T \begin{pmatrix} 0 & 0 \\ 0 & I \end{pmatrix} \nabla U(q)$ , which is preserved up to  $\mathcal{O}(\omega^{-2})$  over exponentially long time intervals.

For  $h\omega \rightarrow 0$ , the long-time near-preservation of the adiabatic invariant  $I$  can be shown using backward error analysis (Reich 1999b), but this argument breaks down for  $h\omega$  bounded away from zero as in Theorem 5.10.

We comment only briefly on the proof of Theorem 5.10; see Hairer *et al.* (2002, Chapter XIII) for the full proof. It is based on representing the numerical solution locally (on bounded intervals) by a *modulated Fourier expansion*

$$q_n = \sum_{|k| < N} z_k(t) e^{ik\tilde{\omega}t} + \mathcal{O}(h^N) \quad \text{for } t = nh,$$

where the coefficients  $z_k(t)$  together with all their derivatives (up to some arbitrarily fixed order) are bounded by  $\mathcal{O}(\tilde{\omega}^{-|k|})$ . A similar representation holds for  $p_n$ . The expansion coefficients  $y_k(t) = z_k(t) e^{ik\tilde{\omega}t}$  satisfy a system of equations similar in structure to (5.2) (but of higher dimension). This permits us to use similar arguments to the second proof of Theorem 5.1 to infer the existence of certain modified energies  $H^*$  and  $I^*$ , which the numerical method preserves up to  $\mathcal{O}(h)$  over times  $h^{-N+1}$ . Finally, the time averages  $\overline{H}_n$  and  $\overline{I}_n$  can be expressed, up to  $\mathcal{O}(h)$ , in terms of these modified energies.

## 6. Constrained Hamiltonian systems

A minimal set of coordinates of a mechanical system is often difficult to find. The minimal coordinates may be defined only implicitly, or frequent changes of charts are necessary along a solution of the system. In this situation it is favourable to formulate the problem as a Hamiltonian system with constraints.

### 6.1. Formulation as differential-algebraic equations

We consider a mechanical system with coordinates  $q \in \mathbb{R}^d$  that are subject to constraints  $g(q) = 0$ . The equations of motion are then of the form

$$\begin{aligned} \dot{p} &= -\nabla_q H(p, q) - \nabla_q g(q) \lambda, \\ \dot{q} &= \nabla_p H(p, q), \quad 0 = g(q), \end{aligned} \tag{6.1}$$

where the Hamiltonian  $H(p, q)$  is usually given by (2.8). Here,  $p$  and  $q$  are vectors in  $\mathbb{R}^d$ ,  $g(q) = (g_1(q), \dots, g_m(q))^T$  is the vector of constraints, and  $\nabla_q g = (\nabla_q g_1, \dots, \nabla_q g_m)$  is the transposed Jacobian matrix of  $g(q)$ .

To compute the Lagrange multiplier  $\lambda$ , we differentiate the constraint  $0 = g(q(t))$  with respect to time. This yields the so-called hidden constraint

$$0 = \nabla_q g(q)^T \nabla_p H(p, q), \tag{6.2}$$

which is an invariant of the flow of (6.1). A further differentiation gives

$$0 = \frac{\partial}{\partial q} (\nabla_q g(q)^T \nabla_p H(p, q)) \nabla_p H(p, q) - \nabla_q g(q)^T \nabla_p^2 H(p, q) (\nabla_q H(p, q) + \nabla_q g(q) \lambda), \tag{6.3}$$

which allows us to express  $\lambda$  in terms of  $(p, q)$ , if the matrix

$$\nabla_q g(q)^T \nabla_p^2 H(p, q) \nabla_q g(q) \quad \text{is invertible} \tag{6.4}$$

( $\nabla_p^2 H$  denotes the Hessian matrix of  $H$ ). Inserting the so-obtained function  $\lambda(p, q)$  into (6.1) gives the ordinary differential equation

$$\begin{aligned} \dot{p} &= -\nabla_q H(p, q) - \nabla_q g(q) \lambda(p, q), \\ \dot{q} &= \nabla_p H(p, q) \end{aligned} \tag{6.5}$$

for  $(p, q)$ , which is well defined on the domain where  $H(p, q)$  and  $g(q)$  are defined, and not only for  $g(q) = 0$ . The standard theory for ordinary differential equations can be used to deduce existence and uniqueness of the solution. Important properties of the system (6.1) are the following.

- Whenever the initial values satisfy  $(p_0, q_0) \in \mathcal{M}$  with

$$\mathcal{M} = \{(p, q) : g(q) = 0, \nabla_q g(q)^T \nabla_p H(p, q) = 0\}, \tag{6.6}$$

the solution stays on the manifold  $\mathcal{M}$  for all  $t$ ; hence, the flow of (6.1) is a mapping  $\varphi_t : \mathcal{M} \rightarrow \mathcal{M}$ .

- The flow  $\varphi_t$  is a symplectic transformation on  $\mathcal{M}$ , which means that

$$(\varphi'_t(p, q)\xi)^T J \varphi'_t(p, q)\eta = \xi^T J \eta \quad \text{for } \xi, \eta \in T_{(p, q)}\mathcal{M}. \tag{6.7}$$

Here,  $T_{(p, q)}\mathcal{M}$  denotes the tangent space of  $\mathcal{M}$  at  $(p, q) \in \mathcal{M}$ , and the product  $\varphi'_t(p, q)\xi$  has to be interpreted as the directional derivative on the manifold.

- For Hamiltonians satisfying

$$H(-p, q) = H(p, q),$$

the flow  $\varphi_t$  is  $\rho$ -reversible for  $\rho(p, q) = (-p, q)$  in the sense that (2.5) holds for  $(p, q) \in \mathcal{M}$ .

The first of these properties follows from the definition of  $\lambda(p, q)$ . For  $(p_0, q_0) \in \mathcal{M}$ , a first integration of (6.3) gives (6.2) and a second integration yields  $g(q) = 0$  along the solution of (6.5).

To prove the symplecticity, we consider the (unconstrained) Hamiltonian system with  $K(p, q) = H(p, q) + g(q)^T \lambda(p, q)$ . Its flow is symplectic and coincides with that of (6.5) on the manifold  $\mathcal{M}$ .

The reversibility is a consequence of the fact that  $H(-p, q) = H(p, q)$  implies  $\lambda(-p, q) = \lambda(p, q)$ . The flow of (6.5) and hence also its restriction onto  $\mathcal{M}$  is thus  $\rho$ -reversible.

**Example 6.1. (Kepler and two-body problems on the sphere)**

Following Kozlov and Harin (1992), we consider a particle moving on the unit sphere attracted by a fixed point  $a$  on the sphere. The potential is given as a fundamental solution of the Laplace–Beltrami equation on the sphere:

$$U(q, a) = -\frac{\cos \vartheta}{\sin \vartheta}, \quad \cos \vartheta = \langle q, a \rangle. \tag{6.8}$$

The Kepler problem on the sphere is then of the form (6.1) with

$$H(p, q) = \frac{1}{2} p^T p + U(q, a), \quad g(q) = q^T q - 1.$$

The left picture of Figure 6.1 shows the solution corresponding to the point  $a = (0.3\sqrt{2}, 0.3\sqrt{2}, 0.8)^T$  and to initial values given in spherical coordinates by  $\varphi_0 = 1, \theta_0 = 1.1$  and  $\dot{\varphi}_0 = 1.2, \dot{\theta}_0 = -1.1$ . The point  $a$  and the initial value are indicated by a larger symbol in Figure 6.1.

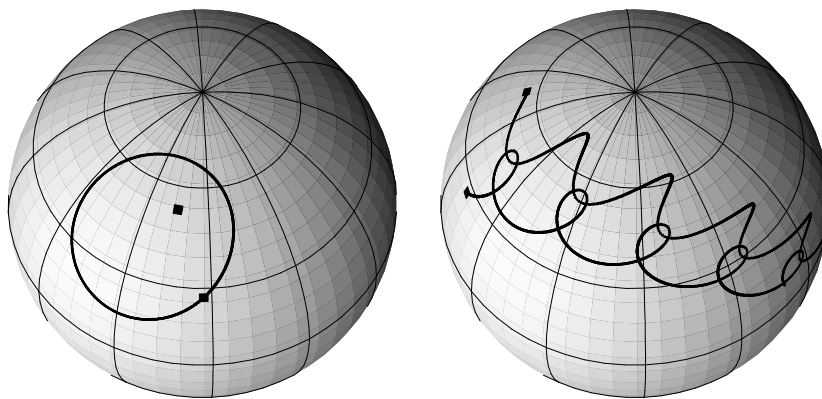


Figure 6.1. Solutions of the Kepler problem (left) and the two-body problem on the sphere (right).

Whereas in Euclidean space the two-body problem reduces to the Kepler problem, this is not the case on the sphere. For the two-body problem the Hamiltonian is

$$H(p_1, p_2, q_1, q_2) = \frac{1}{2} p_1^T p_1 + \frac{1}{2} p_2^T p_2 + U(q_1, q_2)$$

with  $U(q_1, q_2)$  given by (6.8). The constraints are  $g_i(q_1, q_2) = q_i^T q_i - 1$  for  $i = 1, 2$ . The solution with initial values  $\varphi_{10} = -0.3, \theta_{10} = 1.1, \varphi_{20} = -0.8, \theta_{20} = 0.6$  and  $\dot{\varphi}_{10} = 0.9, \dot{\theta}_{10} = -0.5, \dot{\varphi}_{20} = 0.3, \dot{\theta}_{20} = -0.1$  is plotted in Figure 6.1 (right).

**Example 6.2. (Rigid body)** The motion of a rigid body with a fixed point chosen at the origin can be described by an orthogonal matrix  $Q(t)$ . Letting  $I_1, I_2, I_3$  denote the moments of inertia of the body, its kinetic energy is

$$T = \frac{1}{2} (I_1 \Omega_1^2 + I_2 \Omega_2^2 + I_3 \Omega_3^2),$$

where the angular velocity  $\Omega = (\Omega_1, \Omega_2, \Omega_3)^T$  of the body is defined by

$$\widehat{\Omega} = \begin{pmatrix} 0 & -\Omega_3 & \Omega_2 \\ \Omega_3 & 0 & -\Omega_1 \\ -\Omega_2 & \Omega_1 & 0 \end{pmatrix} = Q^T \dot{Q},$$

(Arnold 1989, Chapter 6). In terms of  $Q$ , the kinetic energy on the manifold  $O(3) = \{Q \mid Q^T Q = I\}$  becomes

$$T = \frac{1}{2} \text{trace}(\widehat{\Omega} D \widehat{\Omega}^T) = \frac{1}{2} \text{trace}(Q^T \dot{Q} D \dot{Q}^T Q) = \frac{1}{2} \text{trace}(\dot{Q} D \dot{Q}^T),$$

where  $D = \text{diag}(d_1, d_2, d_3)$  is given by the relations  $I_1 = d_2 + d_3, I_2 = d_3 + d_1,$  and  $I_3 = d_1 + d_2$ . With  $P = \partial T / \partial \dot{Q} = \dot{Q} D$ , we are thus concerned with

$$H(P, Q) = \frac{1}{2} \text{trace}(P D^{-1} P^T) + U(Q),$$

and the constrained Hamiltonian system becomes

$$\begin{aligned} \dot{P} &= -\nabla_Q U(Q) - Q \Lambda, \\ \dot{Q} &= P D^{-1}, \quad 0 = Q^T Q - I, \end{aligned} \tag{6.9}$$

where  $\Lambda$  is a symmetric matrix consisting of Lagrange multipliers. This is of the form (6.1) and satisfies the regularity condition (6.4).

### 6.2. Development of the Rattle algorithm

The most important numerical algorithm for the solution of constrained Hamiltonian systems is an adaptation of the Störmer–Verlet method. Its historical development is in three main steps.

**First step.** For Hamiltonians  $H(p, q) = \frac{1}{2}p^T M^{-1}p + U(q)$  with constant mass matrix  $M$  (cf. Section 2.2), the problem is a second-order differential equation  $M\ddot{q} = -\nabla_q U(q) - \nabla_q g(q)\lambda$  with constraint  $g(q) = 0$ . The most natural extension of (1.2) is

$$\begin{aligned} q_{n+1} - 2q_n + q_{n-1} &= -h^2 M^{-1}(\nabla_q U(q_n) + \nabla_q g(q_n)\lambda_n), \\ 0 &= g(q_{n+1}). \end{aligned} \quad (6.10)$$

This algorithm (called *Shake*) was originally proposed by Ryckaert, Ciccotti and Berendsen (1977) for computations in molecular dynamics. The  $p$ -components, not used in the recursion, are approximated by  $p_n = M(q_{n+1} - q_{n-1})/2h$ .

**Second step.** A one-step formulation of this method, obtained by a formal analogy to formula (1.5), reads

$$\begin{aligned} p_{n+1/2} &= p_n - \frac{h}{2}(\nabla_q U(q_n) + \nabla_q g(q_n)\lambda_n), \\ q_{n+1} &= q_n + hM^{-1}p_{n+1/2}, \quad 0 = g(q_{n+1}), \\ p_{n+1} &= p_{n+1/2} - \frac{h}{2}(\nabla_q U(q_{n+1}) + \nabla_q g(q_{n+1})\lambda_{n+1}). \end{aligned} \quad (6.11)$$

This formula cannot be implemented, because  $\lambda_{n+1}$  is not yet available at this step (it is computed together with  $q_{n+2}$ ). As a remedy, Andersen (1983) suggests replacing the last line in (6.11) with the projection step

$$\begin{aligned} p_{n+1} &= p_{n+1/2} - \frac{h}{2}(\nabla_q U(q_{n+1}) + \nabla_q g(q_{n+1})\mu_n), \\ 0 &= \nabla_q g(q_{n+1})^T M^{-1}p_{n+1}. \end{aligned} \quad (6.12)$$

This modification, called *Rattle*, is motivated by the fact that the numerical approximation  $(p_{n+1}, q_{n+1})$  lies on the solution manifold  $\mathcal{M}$ .

**Third step.** Jay (1994) and Reich (1993) observed independently that the Rattle method can be interpreted as a partitioned Runge–Kutta method and thus allows the extension to general Hamiltonians

$$\begin{aligned} p_{n+1/2} &= p_n - \frac{h}{2}(\nabla_q H(p_{n+1/2}, q_n) + \nabla_q g(q_n)\lambda_n), \\ q_{n+1} &= q_n + \frac{h}{2}(\nabla_p H(p_{n+1/2}, q_n) + \nabla_p H(p_{n+1/2}, q_{n+1})), \\ 0 &= g(q_{n+1}), \\ p_{n+1} &= p_{n+1/2} - \frac{h}{2}(\nabla_q H(p_{n+1/2}, q_{n+1}) + \nabla_q g(q_{n+1})\mu_n), \\ 0 &= \nabla_q g(q_{n+1})^T \nabla_p H(p_{n+1}, q_{n+1}) \end{aligned} \quad (6.13)$$

whenever  $(p_n, q_n) \in \mathcal{M}$ . The first three equations of (6.13) determine  $(p_{n+1/2}, q_{n+1}, \lambda_n)$ , whereas the remaining two are equations for  $(p_{n+1}, \mu_n)$ .

For a sufficiently small step size, these equations have a locally unique solution (Hairer *et al.* 2002, p. 214).

**Example 6.3. (Kepler problem on the sphere)** We apply the Rattle method with a large step size  $h = 0.07$  to the problem of Example 6.1. The numerical solution, plotted in Figure 6.2, shows a precession as it appears in computations with symplectic integrators for the Kepler problem in Euclidean space; see Figure 1.5. We remark that the value of the Hamiltonian along the numerical solution oscillates around the correct value and the energy error remains bounded by 0.114 on very long time intervals.

Since the constraint  $g(q)$  is quadratic and the Hamiltonian is separable, the formulae (6.13) are explicit with exception of the computation of  $\lambda_n$ , for which a scalar quadratic equation needs to be solved.

**Example 6.4. (Rigid body)** The Rattle method (6.13) applied to (6.9) yields

$$\begin{aligned} P_{1/2} &= P_0 - \frac{h}{2} \nabla_Q V(Q_0) - \frac{h}{2} Q_0 \Lambda_1, \\ Q_1 &= Q_0 + h P_{1/2} D^{-1}, \quad Q_1^T Q_1 = I, \\ P_1 &= P_{1/2} - \frac{h}{2} \nabla_Q V(Q_1) - \frac{h}{2} Q_1 \Lambda_2, \quad D^{-1} P_1^T Q_1 + Q_1^T P_1 D^{-1} = 0, \end{aligned} \tag{6.14}$$

where both  $\Lambda_1$  and  $\Lambda_2$  are symmetric matrices. For consistent initial values,  $Q_0$  is orthogonal and  $Q_0^T P_0 D^{-1} = \widehat{\Omega}_0$  is skew-symmetric. Working with

$$\widehat{\Omega}_0 = Q_0^T \dot{Q}_0 = Q_0^T P_0 D^{-1}, \quad \widehat{\Omega}_{1/2} = Q_0^T P_{1/2} D^{-1}, \quad \widehat{\Omega}_1 = Q_1^T P_1 D^{-1},$$

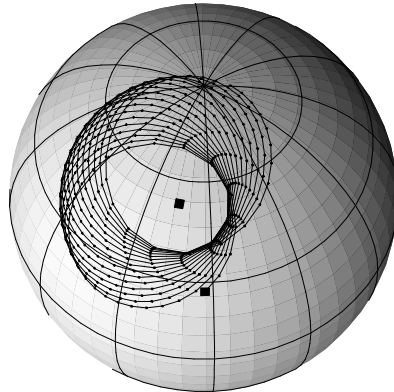


Figure 6.2. Numerical solution of the Kepler problem on the sphere, obtained with the Rattle method using step size  $h = 0.07$ .

instead of  $P_0, P_{1/2}, P_1$ , the equations (6.14) become the following integrator  $(Q_0, \widehat{\Omega}_0) \mapsto (Q_1, \widehat{\Omega}_1)$ :

- find an orthogonal matrix  $I + h\widehat{\Omega}_{1/2}$  such that

$$\widehat{\Omega}_{1/2} = \widehat{\Omega}_0 - \frac{h}{2} Q_0^T \nabla_Q V(Q_0) D^{-1} - \frac{h}{2} \Lambda_1 D^{-1}$$

holds with a symmetric matrix  $\Lambda_1$ ;

- compute  $Q_1 = Q_0(I + h\widehat{\Omega}_{1/2})$ ;
- compute a skew-symmetric matrix  $\widehat{\Omega}_1$  such that

$$\widehat{\Omega}_1 = \widehat{\Omega}_{1/2} - \frac{h}{2} Q_1^T \nabla_Q V(Q_1) D^{-1} - (\widehat{\Omega}_{1/2} + \widehat{\Omega}_{1/2}^T) - \frac{h}{2} \Lambda_2 D^{-1}$$

holds with a symmetric matrix  $\Lambda_2$ .

This algorithm for the simulation of the heavy top is proposed in McLachlan and Scovel (1995). An efficient implementation uses the representation of the appearing orthogonal matrices by quaternions (Hairer 2003).

### 6.3. Geometric properties of Rattle

For consistent initial values  $(p_n, q_n) \in \mathcal{M}$ , the Rattle method (6.13) yields an approximation  $(p_{n+1}, q_{n+1})$  which is again on  $\mathcal{M}$ . We thus have a numerical flow  $\Phi_h : \mathcal{M} \rightarrow \mathcal{M}$ . The geometric properties of Section 2 for the Störmer–Verlet method extend to this algorithm.

**Theorem 6.5.** The Rattle method is symmetric, that is,  $\Phi_h = \Phi_{-h}^{-1}$  on  $\mathcal{M}$ . For Hamiltonians satisfying  $H(-p, q) = H(p, q)$ , the method is reversible with respect to the reflection  $\rho(p, q) = (-p, q)$ , that is, it satisfies  $\rho \circ \Phi_h = \Phi_h^{-1} \circ \rho$  on  $\mathcal{M}$ .

The proof is by straightforward verification, as for the Störmer–Verlet method.

**Theorem 6.6.** The Rattle method is symplectic, that is,

$$(\Phi_h'(p, q)\xi)^T J \Phi_h'(p, q)\eta = \xi^T J \eta \quad \text{for } \xi, \eta \in T_{(p, q)}\mathcal{M}. \quad (6.15)$$

This result was first proved by Leimkuhler and Skeel (1994) for the method (6.11)–(6.12), and by Jay (1994) and Reich (1993) for the general case (6.13).

One proof of Theorem 6.6 is by computing  $\Phi_h'(p, q)\xi$  using implicit differentiation, and by verifying the identity (6.15). Further proofs are based on the interpretation as a variational integrator (Marsden and West 2001), and on explicit formulae of a generating function as in (2.18); see Hairer (2003).



## 7. Geometric integration beyond Störmer–Verlet

In this article we deliberately considered only the Störmer–Verlet method and a few selected geometric properties. Even within the class of ordinary differential equations, we have not mentioned important topics of geometric integration, such as

- higher-order methods, for instance, symmetric composition, partitioned Runge–Kutta, and linear multistep methods,
- the structure-preserving use of variable step sizes,
- differential equations with further geometric properties such as differential equations on Lie groups, problems with multiple time scales, *etc.*

The reader will find more on these topics in the monographs by Sanz-Serna and Calvo (1994) and Hairer *et al.* (2002), in the special journal issue Budd and Iserles (1999), and in the survey articles by Iserles, Munthe-Kaas, Nørsett and Zanna (2000), Marsden and West (2001) and McLachlan and Quispel (2002).

## REFERENCES

- R. Abraham and J. E. Marsden (1978), *Foundations of Mechanics*, 2nd edn, Benjamin/Cummings, Reading, MA.
- H. C. Andersen (1983), Rattle: A ‘velocity’ version of the Shake algorithm for molecular dynamics calculations, *J. Comput. Phys.* **52**, 24–34.
- V. I. Arnold (1963), Small denominators and problems of stability of motion in classical and celestial mechanics, *Russian Math. Surveys* **18**, 85–191.
- V. I. Arnold (1989), *Mathematical Methods of Classical Mechanics*, 2nd edn, Springer, New York.
- G. Benettin and A. Giorgilli (1994), On the Hamiltonian interpolation of near to the identity symplectic mappings with application to symplectic integration algorithms, *J. Statist. Phys.* **74**, 1117–1143.
- J. J. Biesiadecki and R. D. Skeel (1993), Dangers of multiple time step methods, *J. Comput. Phys.* **109**, 318–328.
- P. B. Bochev and C. Scovel (1994), On quadratic invariants and symplectic structure, *BIT* **34**, 337–345.
- C. J. Budd and A. Iserles, eds (1999), *Geometric Integration: Numerical Solution of Differential Equations on Manifolds*, special issue of *R. Soc. Lond. Philos. Trans. Ser. A* **357**(1754).
- R. De Vogelaere (1956), Methods of integration which preserve the contact transformation property of the Hamiltonian equations, Report No. 4, Department of Mathematics, University of Notre Dame, IN.
- K. Feng (1985), On difference schemes and symplectic geometry, in *Proc. Fifth Intern. Symposium on Differential Geometry and Differential Equations, August 1984, Beijing*, pp. 42–58.

- K. Feng (1991), Formal power series and numerical algorithms for dynamical systems, in *Proc. International Conference on Scientific Computation, Hangzhou, China* (T. Chan and Z.-C. Shi, eds), Vol. 1 of *Series on Appl. Math.*, pp. 28–35.
- K. Feng and Z. Shang (1995), Volume-preserving algorithms for source-free dynamical systems, *Numer. Math.* **71**, 451–463.
- E. Fermi, J. Pasta and S. Ulam (1955), Studies of nonlinear problems, Los Alamos Report No. LA-1940. Later published in *E. Fermi: Collected Papers*, Vol. II, Chicago University Press (1965), pp. 978–988.
- R. Feynman (1965), *The Character of Physical Law*, first published by the BBC (1965). MIT Press (1967).
- H. Flaschka (1974), The Toda lattice, II: Existence of Integrals, *Phys. Rev. B* **9**, 1924–1925.
- J. Ford (1992), The Fermi–Pasta–Ulam problem: Paradox turns discovery, *Physics Reports* **213**, 271–310.
- L. Galgani, A. Giorgilli, A. Martinoli and S. Vanzini (1992), On the problem of energy equipartition for large systems of the Fermi–Pasta–Ulam type: Analytical and numerical estimates, *Physica D* **59**, 334–348.
- G. Gallavotti (1999), *Statistical Mechanics: A Short Treatise*, Springer, Berlin.
- E. Hairer (2003), Global modified Hamiltonian for constrained symplectic integrators, *Numer. Math.* To appear.
- E. Hairer and Ch. Lubich (1997), The life-span of backward error analysis for numerical integrators, *Numer. Math.* **76**, 441–462.
- E. Hairer and Ch. Lubich (1999), Invariant tori of dissipatively perturbed Hamiltonian systems under symplectic discretization, *Appl. Numer. Math.* **29**, 57–71.
- E. Hairer and Ch. Lubich (2000a), Long-time energy conservation of numerical methods for oscillatory differential equations, *SIAM J. Numer. Anal.* **38**, 414–441.
- E. Hairer and Ch. Lubich (2000b), Energy conservation by Störmer-type numerical integrators, in *Numerical Analysis 1999* (G. F. Griffiths and G. A. Watson, eds), CRC Press LLC, pp. 169–190.
- E. Hairer, S. P. Nørsett and G. Wanner (1993), *Solving Ordinary Differential Equations I: Nonstiff Problems*, Springer, Heidelberg.
- E. Hairer, Ch. Lubich and G. Wanner (2002), *Geometric Numerical Integration: Structure-Preserving Algorithms for Ordinary Differential Equations*, Springer, Berlin.
- A. Iserles, H. Z. Munthe-Kaas, S. P. Nørsett and A. Zanna (2000), Lie-group methods, in *Acta Numerica*, Vol. 9, Cambridge University Press, pp. 215–365.
- L. Jay (1994), Runge–Kutta type methods for index three differential-algebraic equations with applications to Hamiltonian systems, Thesis No. 2658, University of Genève.
- V. V. Kozlov and A. O. Harin (1992), Kepler’s problem in constant curvature spaces, *Celestial Mech. Dynam. Astronom.* **54**, 393–399.
- F. M. Lasagni (1988), Canonical Runge–Kutta methods, *ZAMP* **39**, 952–953.
- B. J. Leimkuhler and R. D. Skeel (1994), Symplectic numerical integrators in constrained Hamiltonian systems, *J. Comput. Phys.* **112**, 117–125.

- R. MacKay (1992), Some aspects of the dynamics of Hamiltonian systems, in *The Dynamics of Numerics and the Numerics of Dynamics* (D. S. Broomhead and A. Iserles, eds), Clarendon Press, Oxford, pp. 137–193.
- S. Maeda (1980), Canonical structure and symmetries for discrete systems, *Math. Japonica* **25**, 405–420.
- J. E. Marsden and M. West (2001), Discrete mechanics and variational integrators, in *Acta Numerica*, Vol. 10, Cambridge University Press, pp. 357–514.
- R. I. McLachlan and P. Atela (1992), The accuracy of symplectic integrators, *Nonlinearity* **5**, 541–562.
- R. I. McLachlan and G. R. W. Quispel (2002), Splitting methods, in *Acta Numerica*, Vol. 11, Cambridge University Press, pp. 341–434.
- R. I. McLachlan and C. Scovel (1995), Equivariant constrained symplectic integration, *J. Nonlinear Sci.* **5**, 233–256.
- P. C. Moan (2002), On backward error analysis and Nekhoroshev stability in the numerical analysis of conservative systems of ODEs, PhD thesis, University of Cambridge.
- J. Moser (1968), Lectures on Hamiltonian systems, *Mem. Amer. Math. Soc.* **81**, 1–60.
- H. Poincaré (1892/1893/1899), *Les Méthodes Nouvelles de la Mécanique Céleste, Tome I–III*, Gauthier-Villars, Paris.
- S. Reich (1993), Symplectic integration of constrained Hamiltonian systems by Runge–Kutta methods, Technical Report 93-13, Department of Computer Science, University of British Columbia.
- S. Reich (1999a), Backward error analysis for numerical integrators, *SIAM J. Numer. Anal.* **36**, 1549–1570.
- S. Reich (1999b), Preservation of adiabatic invariants under symplectic discretization, *Appl. Numer. Math.* **29**, 45–56.
- R. D. Ruth (1983), A canonical integration technique, *IEEE Trans. Nuclear Science* **NS-30**, 2669–2671.
- J.-P. Ryckaert, G. Ciccotti and H. J. C. Berendsen (1977), Numerical integration of the cartesian equations of motion of a system with constraints: Molecular dynamics of  $n$ -alkanes, *J. Comput. Phys.* **23**, 327–341.
- J. M. Sanz-Serna (1988), Runge–Kutta schemes for Hamiltonian systems, *BIT* **28**, 877–883.
- J. M. Sanz-Serna (1992), Symplectic integrators for Hamiltonian problems: An overview, in *Acta Numerica*, Vol. 1, Cambridge University Press, pp. 243–286.
- J. M. Sanz-Serna and M. P. Calvo (1994), *Numerical Hamiltonian Problems*, Chapman and Hall, London.
- Z. Shang (1999), KAM theorem of symplectic algorithms for Hamiltonian systems, *Numer. Math.* **83**, 477–496.
- Z. Shang (2000), Resonant and diophantine step sizes in computing invariant tori of Hamiltonian systems, *Nonlinearity* **13**, 299–308.
- C. L. Siegel and J. K. Moser (1971), Lectures on Celestial Mechanics, Vol. 187 of *Grundlehren der Mathematischen Wissenschaft*, Springer, Heidelberg.
- D. Stoffer (1988), On reversible and canonical integration methods, SAM-Report No. 88-05, ETH Zürich.

- D. Stoffer (1998), On the qualitative behaviour of symplectic integrators, III: Perturbed integrable systems, *J. Math. Anal. Appl.* **217**, 521–545.
- G. Strang (1968), On the construction and comparison of difference schemes, *SIAM J. Numer. Anal.* **5**, 506–517.
- Y. B. Suris (1988), On the conservation of the symplectic structure in the numerical solution of Hamiltonian systems, in *Numerical Solution of Ordinary Differential Equations* (S. S. Filippov, ed.), Keldysh Institute of Applied Mathematics, USSR Academy of Sciences, Moscow, pp. 148–160. (In Russian.)
- Y.-F. Tang (1994), Formal energy of a symplectic scheme for Hamiltonian systems and its applications, I, *Comput. Math. Appl.* **27**, 31–39.
- M. Toda (1970), Waves in nonlinear lattice, *Progr. Theor. Phys. Suppl.* **45**, 174–200.
- L. Verlet (1967), Computer ‘experiments’ on classical fluids, I: Thermodynamical properties of Lennard-Jones molecules, *Phys. Rev.* **159**, 98–103.
- A. P. Veselov (1991), Integrable maps, *Russ. Math. Surv.* **46**, 1–51.
- H. Yoshida (1993), Recent progress in the theory and application of symplectic integrators, *Celestial Mech. Dynam. Astronom.* **56**, 27–43.