

Qualitative properties of modified equations

O. GONZALEZ†

Department of Mathematics, Swiss Federal Institute of Technology, Lausanne, Switzerland

D. J. HIGHAM‡

Department of Mathematics, University of Strathclyde, Glasgow G1 1XH, UK

AND

A. M. STUART§

Program in Scientific Computing and Computational Mathematics & Division of Mechanics and Computation, Durand Building 257, Department of Mechanical Engineering, Stanford University, Stanford, CA 94305-4040, USA

[Received 7 July 1997 and in revised form 20 April 1998]

Suppose that a consistent one-step numerical method of order r is applied to a smooth system of ordinary differential equations. Given any integer $m \geq 1$, the method may be shown to be of order $r + m$ as an approximation to a certain modified equation. If the method and the system have a particular qualitative property then it is important to determine whether the modified equations inherit this property. In this article, a technique is introduced for proving that the modified equations inherit qualitative properties from the method and the underlying system. The technique uses a straightforward contradiction argument applicable to arbitrary one-step methods and does not rely on the detailed structure of associated power series expansions. Hence the conclusions apply, but are not restricted, to the case of Runge–Kutta methods. The new approach unifies and extends results of this type that have been derived by other means: results are presented for integral preservation, reversibility, inheritance of fixed points, Hamiltonian problems and volume preservation. The technique also applies when the system has an integral that the method preserves not exactly, but to order greater than r . Finally, a negative result is obtained by considering a gradient system and gradient numerical method possessing a global property that is not shared by the associated modified equations.

1. Introduction

In this article we consider the relationship between solutions to a given system of ordinary differential equations, numerical approximations to them, and solutions to associated modified equations. Our goal is to show that if an underlying system and an approximation scheme possess solutions sharing certain qualitative properties, then there is a family of associated modified equations possessing solutions that also share these properties.

†Work supported by a National Science Foundation Graduate Fellowship.

‡Work supported by the Engineering and Physical Sciences Research Council of the UK under grant GR/K80228.

§Work supported by the National Science Foundation under grant DMS-9201727 and by the Office of Naval Research under grant N00014-92-J-1876.

Although many results of this type may already be found in the literature, we provide a general framework and a unified method of proof. Our method of proof, which is analytical, is not as succinct as other more algebraic approaches, but it is readily accessible to readers with a basic knowledge of numerical methods for initial value problems.

The presentation is structured as follows. In Section 2 we establish the notation for the semigroups generated by the underlying differential equation and the numerical method. We work within a general class of one-step methods that satisfy a certain local approximation property: the expansions of the true and approximate semigroups in powers of the time-step Δt agree up to order r . (Since we are dealing with ordinary differential equations the local semigroup property may be extended to a local group property, but in many interesting applications global existence is only assured in forward time so that we retain the concept of *semigroup*.) Standard error estimates for such methods show that the error over a finite time interval is of $\mathcal{O}(\Delta t^r)$. Runge–Kutta methods are included in our framework, together with a variety of non-standard one-step methods used in practice, such as those that ensure conservation of invariants for Hamiltonian systems or volume for systems defined by divergence-free vector fields.

In Section 3 we discuss modified equations. Here, given an integer $m \geq 1$, the idea is to find an $\mathcal{O}(\Delta t^r)$ modification of the original ordinary differential equation with the property that the numerical method is $\mathcal{O}(\Delta t^{r+m})$ accurate as an approximation of this *modified equation*. We prove a general result concerning the existence and approximation properties of modified equations for the general class of one-step methods introduced in Section 2. This straightforward result is well known in the numerical analysis literature, and also has related counterparts in the mathematical physics literature where interpolation of near-identity maps is important. Our proof simply sets the notation and methodology used in the remainder of the article. Modified equations were first studied in detail in the numerical analysis community by Warming & Hyett (1974) within the context of partial differential equations. In this area the modified equation approach is often useful in interpreting the qualitative properties of errors introduced by numerical approximation, such as numerical dissipation or dispersion for wave propagation problems. For further results on the usefulness and applicability of the modified equation approach, especially in the context of ordinary differential equations, see Griffiths & Sanz-Serna (1986), Beyn (1991), Reich (1993, 1996), Calvo, Murua & Sanz-Serna (1994), Hairer (1994), Reddien (1995), Fiedler & Scheurle (1996) and Sanz-Serna & Murua (1997). A major application has been the derivation of exponentially small error estimates, starting with the work of Neishtadt (1984) and continued in, for example, Benettin & Giorgilli (1994), Hairer & Lubich (1997), and Reich (1996).

The main contribution of this article is contained in Section 4 where the qualitative properties of modified equations are studied. By use of a straightforward contradiction argument we show that if the numerical method shares a certain structural property with the underlying system, then the family of associated modified equations inherits this property. Such results motivate the use of numerical methods that respect qualitative features of the ordinary differential equation. The specific structural properties that we consider are preservation of a scalar function, reversibility, inheritance of fixed points, conservation of the canonical symplectic two-form for Hamiltonian systems and conservation of volume. A negative example is also given: we show that the global limit set behaviour

of gradient systems is not necessarily shared by modified equations for gradient numerical methods.

The fact that structure-preserving numerical methods may possess modified equations with analogous structure has been known for some time. For early references in the numerical analysis literature see Mackay (1992), Sanz-Serna (1992) and Sanz-Serna & Calvo (1994) and for some early applications see Auerbach & Friedman (1991) and Yoshida (1993). More recently, modified equations for problems with special structure have been explored in greater detail using several different approaches, as in Hairer (1994), Calvo, Murua & Sanz-Serna (1994), Benettin & Giorgilli (1994), Reich (1993, 1996), Hairer & Stoffer (1997) and Sanz-Serna & Murua (1997).

The work of Reich, in particular, describes a very general algebraic approach to the study of structure in modified equations which we now outline for the purposes of comparison with our analytic approach. For differential equations whose vector fields lie in a certain linear subspace g of the infinite-dimensional Lie algebra of smooth vector fields, the semigroup generated lies in a corresponding subset G of the Fréchet manifold of smooth diffeomorphisms. Assuming G is a submanifold and that the tangent space to G at the identity is g , Reich shows that general one-step methods with semigroups in G possess modified equations with vector fields in g . In this sense, the modified equations possess the same structure as the underlying differential equation. The primary goal of this article is to demonstrate that a straightforward contradiction argument may be applied to a general one-step method to obtain similar results, without directly relying upon the geometrical relationship between G and g . Our approach is less succinct than that of Reich, but fits more naturally into a traditional numerical analysis framework. The approach also allows us to prove some new results: Theorems 4.2 and 4.4. Theorem 4.2 concerns methods that ‘almost’ preserve an integral of the system; that is, they preserve the integral to a higher order than the classic order suggests. Examples of such methods have been proposed by Calvo, Iserles & Zanna (1996). Theorem 4.4 concerns the preservation of fixed points.

Most of the work mentioned in the references above, and also the analysis presented here, applies to fixed-stepsize implementations. Reich (1996), Hairer & Stoffer (1997) and Hairer (1997) have recently explored the idea of developing customized variable-stepsize strategies for which an appropriate modified equation theory exists.

2. Background

Consider a system of ordinary differential equations in \mathbb{R}^p of the form

$$\frac{du}{dt} = f(u) \quad (2.1)$$

where the vector field $f : \mathbb{R}^p \rightarrow \mathbb{R}^p$ is assumed to be of class C^∞ . For any $u_0 \in \mathbb{R}^p$ we denote by $S : B \times [0, T] \rightarrow \mathbb{R}^p$ the local evolution semigroup generated by (2.1) where B is a closed ball at u_0 and $T > 0$. In particular, for any $U \in B$ the curve

$$u(t) = S(U, t) = S_t(U) \quad (2.2)$$

is a solution to (2.1) with initial condition $u(0) = U$, defined for all $t \in [0, T]$. Furthermore, for each $t \in [0, T]$ the mapping $S_t : B \rightarrow \mathbb{R}^p$ is a C^∞ diffeomorphism

onto its image, and we denote its derivative at a point $U \in B$ by $dS_t(U) \in \mathbb{R}^{p \times p}$. We will use the fact that the mapping $B \times [0, T] \ni (U, t) \mapsto dS_t(U) \in \mathbb{R}^{p \times p}$ is continuous in U and continuously differentiable in t , and we note that $dS_t(U)$ is invertible for each $U \in B$ and $t \in [0, T]$. Hence, by compactness, there exist real numbers $C_i > 0$ ($i = 1, \dots, 4$) such that

$$C_1 \leq \|dS_t(U)\| \leq C_2 \quad \text{and} \quad C_3 \leq \|dS_t(U)^{-1}\| \leq C_4, \quad (2.3)$$

for all $U \in B$ and $t \in [0, T]$, where $\|\cdot\|$ denotes the Frobenius norm on $\mathbb{R}^{p \times p}$.

We will consider one-step numerical methods for (2.1) of the form

$$\mathcal{G}_{\Delta t}(U_n, U_{n+1}) = 0, \quad (2.4)$$

where $\mathcal{G}_{\Delta t} : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}^p$ is a given C^∞ map that depends smoothly on the parameter Δt . For any $u_0 \in \mathbb{R}^p$ we assume the numerical scheme generates a local evolution semigroup in the sense that there is a closed ball \mathcal{B} at u_0 , real numbers $h, T > 0$, and a mapping $\bar{S}_{\Delta t} : \mathcal{B} \rightarrow \mathbb{R}^p$ such that for any $U \in \mathcal{B}$ and $\Delta t \in [0, h]$ the sequence $\{U_n\}$ generated by

$$U_n = \bar{S}_{\Delta t}^n(U) \quad (2.5)$$

satisfies (2.4) for all $n\Delta t \in [0, T]$. Here $\bar{S}_{\Delta t}^n$ denotes the n -fold composition of the map $\bar{S}_{\Delta t}$.

Given any $u_0 \in \mathbb{R}^p$ we assume without loss of generality that $\mathcal{B} = B$ and $T = T$. Furthermore, we assume the numerical scheme is consistent of order r as an approximation to (2.1); that is, for any $U \in B$ we have

$$\left. \frac{\partial^i}{\partial t^i} \right|_{t=0} \bar{S}_t(U) = \left. \frac{\partial^i}{\partial t^i} \right|_{t=0} S_t(U), \quad i = 1, \dots, r, \quad (2.6)$$

where $r \geq 1$ by consistency.

For any $n\Delta t \in [0, T]$ with $\Delta t \in [0, h]$ let $d\bar{S}_{\Delta t}^n(U) \in \mathbb{R}^{p \times p}$ denote the derivative of $\bar{S}_{\Delta t}^n : B \rightarrow \mathbb{R}^p$ at a point $U \in B$, and let $\|\cdot\|$ denote the Euclidean norm on \mathbb{R}^p . Then, by standard results from the numerical analysis of ordinary differential equations (see, for example, Stuart & Humphries (1996, Theorem 6.2.1)) there exist real numbers $C_5 > 0$ and $C_6 > 0$ depending on $U \in B$ and T such that

$$\|S_t(U) - \bar{S}_{\Delta t}^n(U)\| \leq C_5 \Delta t^r \quad (2.7)$$

and

$$\|dS_t(U) - d\bar{S}_{\Delta t}^n(U)\| \leq C_6 \Delta t^r \quad (2.8)$$

for any $t = n\Delta t \in [0, T]$ with $\Delta t \in [0, h]$. Additionally, in view of (2.3) and (2.8), there is a real number $C_7 > 0$ depending on $U \in B$ and T such that the derivative of the mapping $\bar{S}_{\Delta t}^n : B \rightarrow \mathbb{R}^p$ satisfies

$$\|d\bar{S}_{\Delta t}^n(U)\| \leq C_7. \quad (2.9)$$

3. Associated modified equations

To any ordinary differential equation of the form (2.1), and numerical approximation scheme (2.4) of order r , we can associate a *modified equation* of index N of the form

$$\frac{dv}{dt} = \tilde{f}_{\Delta t}^{(N)}(v), \quad (3.1)$$

where $N \geq 1$ is an integer and the *modified vector field* $\tilde{f}_{\Delta t}^{(N)} : \mathbb{R}^p \rightarrow \mathbb{R}^p$ is defined as

$$\tilde{f}_{\Delta t}^{(N)}(v) = f(v) + \sum_{i=1}^N \Delta t^{r+i-1} q_i(v) \quad (3.2)$$

for some functions $q_i : \mathbb{R}^p \rightarrow \mathbb{R}^p$ ($i = 1, \dots, N$). It is convenient to define the modified equation of index zero to be the original equation (2.1). Thus we have

$$\tilde{f}_{\Delta t}^{(0)}(u) = f(u), \quad \forall u \in \mathbb{R}^p. \quad (3.3)$$

For any $v_0 \in \mathbb{R}^p$ we denote by $\tilde{S}^{(N)} : \tilde{B} \times [0, \tilde{T}] \times [0, \tilde{h}] \rightarrow \mathbb{R}^p$ the local evolution semigroup generated by (3.1) where \tilde{B} is a closed ball at v_0 and $\tilde{h}, \tilde{T} > 0$. In particular, for any $V \in \tilde{B}$ and $\Delta t \in [0, \tilde{h}]$ the curve defined by

$$v_{\Delta t}(t) = \tilde{S}_{\Delta t}^{(N)}(V, t) = \tilde{S}_{t, \Delta t}^{(N)}(V) \quad (3.4)$$

is a solution to (3.1) with initial condition $v_{\Delta t}(0) = V$, defined for all $t \in [0, \tilde{T}]$. For any $t \in [0, \tilde{T}]$ and $\Delta t \in [0, \tilde{h}]$ we denote by $d\tilde{S}_{t, \Delta t}^{(N)}(V) \in \mathbb{R}^{p \times p}$ the derivative of the mapping $\tilde{S}_{t, \Delta t}^{(N)} : \tilde{B} \rightarrow \mathbb{R}^p$ at a point $V \in \tilde{B}$. As for the underlying system, we will use the fact that, for any $\Delta t \in [0, \tilde{h}]$, the mapping $\tilde{B} \times [0, \tilde{T}] \ni (V, t) \mapsto d\tilde{S}_{t, \Delta t}^{(N)}(V) \in \mathbb{R}^{p \times p}$ is continuous in V and continuously differentiable in t , and we note that $d\tilde{S}_{t, \Delta t}^{(N)}(V)$ is invertible for each $V \in \tilde{B}$ and $t \in [0, \tilde{T}]$. In what follows we will consider the semigroups generated by (2.1), (2.4) and (3.1) in the neighbourhood of a common point in \mathbb{R}^p . In this case, without loss of generality, we take $\tilde{B} = B$ and $\tilde{T} = T$.

With the above notation in hand we next show that, given an integer $N \geq 1$, it is possible to construct C^∞ functions q_i ($i = 1, \dots, N$) such that the numerical scheme (2.4) is an order $r + N$ approximation to (3.1). The local construction of the functions q_i is outlined in the lemma below and is based on the observation that the N th modified vector field (3.2) differs from the $(N + 1)$ st by a single term, that is

$$\tilde{f}_{\Delta t}^{(N+1)}(v) = \tilde{f}_{\Delta t}^{(N)}(v) + \Delta t^{r+N} q_{N+1}(v).$$

With this observation the basic strategy for constructing the functions q_i becomes clear: choose the q_i such that each new term in (3.2) improves the order of approximation by one power of Δt .

LEMMA 3.1 Given $u_0 \in \mathbb{R}^p$ suppose that for some integer $s > r$ the expansions

$$\begin{aligned} S_{\Delta t}(u_0) &= u_0 + \sum_{i=1}^s \Delta t^i F_i(u_0) + \mathcal{O}(\Delta t^{s+1}) \\ \bar{S}_{\Delta t}(u_0) - S_{\Delta t}(u_0) &= \sum_{i=r+1}^s \Delta t^i F_i^{(0)}(u_0) + \mathcal{O}(\Delta t^{s+1}) \end{aligned}$$

are valid, where $\{F_i\}_{i=1}^s$ and $\{F_i^{(0)}\}_{i=r+1}^s$ are C^∞ . Then, for every integer N such that $0 \leq N \leq s - r - 1$, there exists a modified equation (3.1) with the property

$$\bar{S}_{\Delta t}(u_0) - \tilde{S}_{\Delta t, \Delta t}^{(N)}(u_0) = \sum_{i=r+N+1}^s \Delta t^i F_i^{(N)}(u_0) + \mathcal{O}(\Delta t^{s+1}) \quad (3.5)$$

where $\{F_i^{(N)}\}_{i=r+N+1}^s$ are C^∞ .

Proof. By assumption, the result is true for $N = 0$. Assuming $s > r + 1$ we proceed by induction. Suppose the result is true for some N with $0 \leq N < s - r - 1$ and let $q_{N+1}(v) := F_{r+N+1}^{(N)}(v)$, so that

$$\tilde{f}_{\Delta t}^{(N+1)}(v) = \tilde{f}_{\Delta t}^{(N)}(v) + \Delta t^{r+N} F_{r+N+1}^{(N)}(v). \quad (3.6)$$

Expanding $\tilde{S}_{t, \Delta t}^{(N+1)}(u_0)$ for small t and setting $t = \Delta t$ leads to an expression of the form

$$\tilde{S}_{\Delta t, \Delta t}^{(N+1)}(u_0) = u_0 + \sum_{i=1}^s \Delta t^i \hat{F}_i(u_0) + \mathcal{O}(\Delta t^{s+1}) \quad (3.7)$$

for some $\{\hat{F}_i\}_{i=1}^s$ in C^∞ . By comparing Taylor expansions of $\tilde{S}_{t, \Delta t}^{(N)}(u_0)$ and $\tilde{S}_{t, \Delta t}^{(N+1)}(u_0)$ with $t = \Delta t$, we find from (3.6) that

$$\tilde{S}_{\Delta t, \Delta t}^{(N)}(u_0) - \tilde{S}_{\Delta t, \Delta t}^{(N+1)}(u_0) = -\Delta t^{r+N+1} F_{r+N+1}^{(N)}(u_0) + \mathcal{O}(\Delta t^{r+N+2}).$$

Hence, using the induction hypothesis (3.5),

$$\begin{aligned} \bar{S}_{\Delta t}(u_0) - \tilde{S}_{\Delta t, \Delta t}^{(N+1)}(u_0) &= [\bar{S}_{\Delta t}(u_0) - \tilde{S}_{\Delta t, \Delta t}^{(N)}(u_0)] + [\tilde{S}_{\Delta t, \Delta t}^{(N)}(u_0) - \tilde{S}_{\Delta t, \Delta t}^{(N+1)}(u_0)] \\ &= \Delta t^{r+N+1} F_{r+N+1}^{(N)}(u_0) - \Delta t^{r+N+1} F_{r+N+1}^{(N)}(u_0) \\ &\quad + \mathcal{O}(\Delta t^{r+N+2}) \\ &= \mathcal{O}(\Delta t^{r+N+2}). \end{aligned}$$

This, along with (3.7), gives the required result. \square

Henceforth, we will always assume that the expansions in the statement of Lemma 3.1 are valid for all integers $s > r$ and that the functions q_i , and hence the modified equations (3.1), are chosen so that (3.5) holds. Since the original vector field f is C^∞ this assumption automatically holds for $S_{\Delta t}$ and will also hold for most methods $\bar{S}_{\Delta t}$ used in practice. The following result then follows from a standard Gronwall-based convergence analysis.

THEOREM 3.2 Given any $u_0 \in \mathbb{R}^p$ and any integer $N \geq 1$ there exists a ball B at u_0 , real numbers $h, T > 0$, and smooth functions q_i ($i = 1, \dots, N$) such that the local evolution semigroups $S_t, \bar{S}_{\Delta t}, \tilde{S}_{t, \Delta t}^{(N)} : B \rightarrow \mathbb{R}^p$ for (2.1), (2.4) and (3.1), respectively, are defined for all $t = n\Delta t \in [0, T]$ with $\Delta t \in [0, h]$. Furthermore, there is a constant $C_8 = C_8(T, N, B) > 0$ such that for each $U \in B$

$$\| \|d\tilde{S}_{t, \Delta t}^{(N)}(U) - d\bar{S}_{\Delta t}^n(U)\| \| + \|\tilde{S}_{t, \Delta t}^{(N)}(U) - \bar{S}_{\Delta t}^n(U)\| \leq C_8 \Delta t^{r+N}, \quad (3.8)$$

for all $t = n\Delta t \in [0, T]$ with $\Delta t \in [0, h]$.

REMARKS

- (1) Note that combining (2.7), (2.8) and (3.8) (and without loss of generality taking the same constant) gives the bound

$$\| \|d\tilde{S}_{t, \Delta t}^{(N)}(U) - dS_t(U)\| \| + \|\tilde{S}_{t, \Delta t}^{(N)}(U) - S_t(U)\| \leq C_8 \Delta t^r, \quad (3.9)$$

which compares solutions of the original and modified equations.

- (2) The key result (3.8) shows that the numerical method applied to the original system (2.1) behaves like an order $r + N$ method with respect to the modified equation of index N . In other words, the method computes a very accurate approximation to a perturbed problem. This idea has links with the concept of backward error in numerical linear algebra, see Golub & Van Loan (1996). It is natural to ask whether the perturbed problem has the same structure as the original problem; in the linear algebra context this is known as structured backward error analysis. In the next section we address this question for a variety of different families of ordinary differential equations.
- (3) Note that the constant C_8 in (3.8) and (3.9) depends upon T, N and B . Hence, in particular, the bounds are valid only for finite time intervals and they are not uniform in the index N . However, by optimizing over the index, Neishtadt (1984), Benettin & Giorgilli (1994), Hairer & Lubich (1997) and Reich (1996) have shown that the difference between the numerical approximation and a modified equation remains exponentially small over arbitrarily long time intervals as $\Delta t \rightarrow 0$ for a variety of problems of interest.

4. Qualitative properties of the modified equations

In this section we employ an induction on N to establish various properties of the modified equation (3.1). In view of Theorem 3.2 we see that, given any $u_0 \in \mathbb{R}^p$, the ball B at u_0 and the values $h, T > 0$ will in general depend upon N . In the following induction arguments we will choose a ball B and numbers $h, T > 0$ such that all the local evolution semigroups $S_t, \bar{S}_{\Delta t}, \tilde{S}_{t, \Delta t}^{(m)} : B \rightarrow \mathbb{R}^p$ ($m = 1, \dots, N + 1$) are defined for any $t = n\Delta t \in [0, T]$ with $\Delta t \in [0, h]$. Note that h may shrink to zero as $N \rightarrow \infty$, but will be non-zero for every fixed integer $N \geq 1$.

By virtue of Theorem 3.2 we may assume, without loss of generality, that the same constants C_1, C_2, C_3 and C_4 which appear in (2.3) may be used to bound the derivatives of the semigroups for the modified equations up to index $N + 1$. Thus, for $1 \leq m \leq N + 1$,

$$C_1 \leq \|d\tilde{S}_{t,\Delta t}^{(m)}(U)\| \leq C_2 \quad \text{and} \quad C_3 \leq \|d\tilde{S}_{t,\Delta t}^{(m)}(U)^{-1}\| \leq C_4. \quad (4.1)$$

4.1 Integrals for the modified semigroup

Suppose that the underlying system (2.1) and the approximation scheme (2.4) share an integral $\mathcal{F} \in C^1(\mathbb{R}^p, \mathbb{R})$. That is, for any $u_0 \in \mathbb{R}^p$ the function \mathcal{F} is invariant under the local semigroups S_t and $\tilde{S}_{\Delta t}$ in the sense that, for any $U \in B$ and $\Delta t \in [0, h]$, we have $\mathcal{F}(S_t(U)) = \mathcal{F}(U)$ and $\mathcal{F}(\tilde{S}_{\Delta t}^n(U)) = \mathcal{F}(U)$ for all $t \in [0, T]$ and $n\Delta t \in [0, T]$. We then have the following result, which is also proved in Reich (1993, 1996) by Lie algebraic methods.

THEOREM 4.1 Suppose the underlying system (2.1) and the approximation scheme (2.4) share an integral $\mathcal{F} \in C^1(\mathbb{R}^p, \mathbb{R})$. Then \mathcal{F} is an integral for the associated modified equation (3.1) of index N for any integer $N \geq 1$. Hence, the modified equation (3.1) has the form

$$\frac{dv}{dt} = f(v) + \Delta t^r \sum_{i=1}^N \Delta t^{i-1} q_i(v),$$

where

$$\nabla \mathcal{F}(v) \cdot q_i(v) = 0, \quad \forall v \in \mathbb{R}^p, \quad i = 1, \dots, N.$$

Proof. For induction assume the modified equation of index N has $\mathcal{F} : \mathbb{R}^p \rightarrow \mathbb{R}$ as an integral. Note that this is true for $N = 0$ since the modified equation of order zero is the original equation (2.1).

Consider any $u_0 \in \mathbb{R}^p$. Then, for any $U \in B$ and $\Delta t \in [0, h]$ we have

$$\mathcal{F}(\tilde{S}_{t,\Delta t}^{(N)}(U)) = \mathcal{F}(U), \quad (4.2)$$

for all $t \in [0, T]$. Equivalently, for any $\Delta t \in [0, h]$, we have

$$\nabla \mathcal{F}(u) \cdot \tilde{f}_{\Delta t}^{(N)}(u) = 0, \quad \forall u \in \mathcal{I}m(\tilde{S}_{\Delta t}^{(N)}), \quad (4.3)$$

where

$$\mathcal{I}m(\tilde{S}_{\Delta t}^{(N)}) = \{u \in \mathbb{R}^p \mid u = \tilde{S}_{t,\Delta t}^{(N)}(U), \quad U \in B, \quad t \in [0, T]\}. \quad (4.4)$$

Now assume, for contradiction, that \mathcal{F} is not an integral for the modified equation of index $N + 1$, which is of the form

$$\frac{dv}{dt} = \tilde{f}_{\Delta t}^{(N+1)}(v) = \tilde{f}_{\Delta t}^{(N)}(v) + \Delta t^{r+N} q_{N+1}(v). \quad (4.5)$$

Then there exists $u_0 \in \mathbb{R}^p$ such that

$$\nabla \mathcal{F}(u_0) \cdot q_{N+1}(u_0) \neq 0. \quad (4.6)$$

Otherwise, $\nabla \mathcal{F}(u) \cdot \tilde{f}_{\Delta t}^{(N+1)}(u) = 0$ for all $u \in \mathbb{R}^p$ and \mathcal{F} would be an integral.

Let $C_9(u_0) = \nabla \mathcal{F}(u_0) \cdot q_{N+1}(u_0)/2 \neq 0$ and assume, without loss of generality, that $C_9 > 0$; otherwise, if $C_9 < 0$, then one can redefine \mathcal{F} by changing sign. By continuity there is a closed ball D at u_0 such that

$$\nabla \mathcal{F}(U) \cdot q_{N+1}(U) \geq C_9 > 0, \quad \forall U \in D. \quad (4.7)$$

Consider a point $U \in D \cap B$ and let $h, T > 0$ be such that, for any $\Delta t \in [0, h]$, the evolution semigroups satisfy $\tilde{S}_{t, \Delta t}^{(N+1)}(U) \in D$ for all $t \in [0, T]$ and $U_n = \tilde{S}_{\Delta t}^n(U) \in D$ for all $n \Delta t \in [0, T]$. Then, for any $\Delta t \in [0, h]$ and $t \in [0, T]$, we have by (4.3) and (4.5)

$$\begin{aligned} \frac{\partial}{\partial \tau} \Big|_{\tau=t} \mathcal{F}(\tilde{S}_{\tau, \Delta t}^{(N+1)}(U)) &= \nabla \mathcal{F}(\tilde{S}_{t, \Delta t}^{(N+1)}(U)) \cdot \tilde{f}_{\Delta t}^{(N+1)}(\tilde{S}_{t, \Delta t}^{(N+1)}(U)) \\ &= \Delta t^{r+N} \nabla \mathcal{F}(\tilde{S}_{t, \Delta t}^{(N+1)}(U)) \cdot q_{N+1}(\tilde{S}_{t, \Delta t}^{(N+1)}(U)) \\ &\geq C_9 \Delta t^{r+N}, \end{aligned} \quad (4.8)$$

which implies

$$|\mathcal{F}(\tilde{S}_{T, \Delta t}^{(N+1)}(U)) - \mathcal{F}(U)| \geq C_9 T \Delta t^{r+N}, \quad (4.9)$$

for all $\Delta t \in [0, h]$.

By compactness of the closed ball D , since $\mathcal{F} \in C^1(\mathbb{R}^p, \mathbb{R})$, there is a real number $C_{10} > 0$ such that

$$|\mathcal{F}(U) - \mathcal{F}(V)| \leq C_{10} \|U - V\|, \quad \forall U, V \in D. \quad (4.10)$$

Furthermore, in view of (3.8), the modified equation of index $N + 1$ and the numerical scheme (2.4) have solutions satisfying

$$\|\tilde{S}_{T, \Delta t}^{(N+1)}(U) - \tilde{S}_{\Delta t}^n(U)\| \leq C_8 \Delta t^{r+N+1}, \quad (4.11)$$

for all $\Delta t = T/n$ and $n \geq n^*$, where n^* is any positive integer such that $T/n^* \in [0, h]$. Since by hypothesis \mathcal{F} is an integral for the local numerical semigroup we use (4.10) and (4.11) to write

$$\begin{aligned} |\mathcal{F}(\tilde{S}_{T, \Delta t}^{(N+1)}(U)) - \mathcal{F}(U)| &= |\mathcal{F}(\tilde{S}_{T, \Delta t}^{(N+1)}(U)) - \mathcal{F}(\tilde{S}_{\Delta t}^n(U))| \\ &\leq C_{10} \|\tilde{S}_{T, \Delta t}^{(N+1)}(U) - \tilde{S}_{\Delta t}^n(U)\| \\ &\leq C_8 C_{10} \Delta t^{r+N+1}, \end{aligned} \quad (4.12)$$

for all $\Delta t = T/n$ and $n \geq n^*$. This yields a contradiction, since for $\Delta t < TC_9/C_8C_{10}$ both (4.9) and (4.12) cannot hold. Hence \mathcal{F} must be an integral for the modified equation of index $N + 1$. The result follows by induction. \square

The same argument can be used in the case where a method fails to preserve an integral exactly, but preserves it to higher accuracy than the classic order would suggest. See, for example, Calvo, Iserles & Zanna (1996) for instances of such methods.

THEOREM 4.2 Suppose that the system (2.1) preserves an integral $\mathcal{F} \in C^1(\mathbb{R}^p, \mathbb{R})$, and that given any $U \in B$, with B compact, there exist $h, T > 0$ such that

$$|\mathcal{F}(\bar{S}_{\Delta t}^n(U)) - \mathcal{F}(U)| \leq C_{11} \Delta t^l, \quad \forall \Delta t \in [0, h], \quad n \Delta t \in [0, T],$$

for some constant $C_{11} = C_{11}(B, T)$ and integer $l > r$. Then all modified equations of index up to and including $l - r$ also preserve \mathcal{F} ; that is,

$$\nabla \mathcal{F}(v) \cdot q_i(v) = 0, \quad \forall v \in \mathbb{R}^p, \quad i = 1, \dots, l - r.$$

Proof. Applying the induction argument used in the previous theorem, we find that the bound (4.12) is degraded by a term of $\mathcal{O}(\Delta t^l)$, and hence the contradiction remains for indices up to $l - r$. \square

We note that Theorem 3.2 yields two corollaries to the above result:

- (i) Over any finite time interval the numerical solution approximates to $\mathcal{O}(\Delta t^l)$ a modified equation that preserves \mathcal{F} (namely, the modified equation of index $l - r$).
- (ii) For any $N \geq 1$ the modified equation of index N approximates the numerical solution to $\mathcal{O}(\Delta t^{r+N})$ and preserves \mathcal{F} to within $\mathcal{O}(\Delta t^l)$ over any finite time interval.

4.2 Reversibility and the modified semigroup

The system (2.1) is said to be *reversible* if there exists an invertible linear transformation $\rho : \mathbb{R}^p \rightarrow \mathbb{R}^p$ such that

$$f \cdot \rho(y) = -\rho \cdot f(y), \quad \forall y \in \mathbb{R}^p, \quad (4.13)$$

where \cdot denotes composition. (This definition, which is also used by Stoffer (1995) and Hairer & Stoffer (1997), is more general than the definition of reversibility that is found in some texts, such as Strogatz (1994).) The implication of property (4.13) on local semigroups can be summarized as follows. Consider two compact sets B_1 and B_2 such that $B_2 = \rho(B_1)$, and let $S_{1,t}, t \in [0, T_1]$ and $S_{2,t}, t \in [0, T_2]$ be the local semigroups for (2.1) defined on B_1 and B_2 . (Note that $S_{1,t}$ and $S_{2,t}$ can be identified by extending their domain to a compact set including both B_1 and B_2 and possibly reducing the time-interval on which they are defined; similar considerations apply to the numerical method and to the modified equations considered below.) The reversibility property (4.13) implies that the time domains of $S_{1,t}$ and $S_{2,t}$ can be extended to $[-T_2, T_1]$ and $[-T_1, T_2]$, respectively, and implies that the semigroups enjoy the property

$$\rho \cdot S_{1,t}(y) = S_{2,-t} \cdot \rho(y), \quad \forall y \in B_1, \quad t \in [-T_2, T_1]. \quad (4.14)$$

Moreover, for any $y \in B_1$ and $T > 0$ such that $S_{1,t}(y) \in B_1$ for all $t \in [-T, T]$, we have

$$S_{2,t} \cdot \rho \cdot S_{1,t}(y) = \rho(y), \quad \forall t \in [-T, T]. \quad (4.15)$$

Another way to characterize the above relations that will prove useful in our development is the following. Define a function $\psi : [-T, T] \rightarrow \mathbb{R}^p$ by

$$\psi(t) = \rho \cdot S_{1,t}(y) - S_{2,-t} \cdot \rho(y). \quad (4.16)$$

Differentiating with respect to time and using (4.16) we see that $\psi(t)$ satisfies an equation of the form

$$\frac{d\psi}{dt} = g(t, \psi) \quad \text{where} \quad g(t, \psi) = \rho \cdot f(S_{1,t}(y)) + f(\rho \cdot S_{1,t}(y) - \psi). \quad (4.17)$$

Since $\psi(0) = 0$, and (4.13) implies $g(t, 0) = 0$ for all $t \in [-T, T]$, we deduce that $\psi(t) \equiv 0$ is the unique solution to (4.17). Hence (4.14) holds, and (4.15) follows from the properties of the semigroup.

Extending the notation established above, a one-step method such as (2.4) is said to be reversible if, whenever the structure (4.13) is present, we have

$$\bar{S}_{2,\Delta t} \cdot \rho \cdot \bar{S}_{1,\Delta t}(y) = \rho(y)$$

for any $y \in B_1$ and $\Delta t \in [0, h]$ such that $\bar{S}_{1,\Delta t}(y) \in B_1$. Here $\bar{S}_{1,\Delta t}$ and $\bar{S}_{2,\Delta t}$ are the local semigroups on B_1 and B_2 generated by the method. Stoffer (1995) showed that all symmetric Runge–Kutta methods are reversible in this sense, and Hairer & Stoffer (1997) showed that when a symmetric Runge–Kutta method is applied to a reversible system, all modified equations are reversible. Below, we use a different technique to show that the same result holds for all one-step methods of the form given in Section 2. A similar result can also be proved by modifying the techniques presented in Reich (1996).

THEOREM 4.3 If equation (2.1) and the numerical method (2.4) are reversible, then so are all modified equations. More precisely, assume (4.13) holds and assume for any compact sets B_1 and B_2 such that $B_2 = \rho(B_1)$ the method defined by (2.4) has the property

$$\bar{S}_{2,\Delta t} \cdot \rho \cdot \bar{S}_{1,\Delta t}(y) = \rho(y) \quad (4.18)$$

for any $y \in B_1$ and $\Delta t \in [0, h]$ such that $\bar{S}_{1,\Delta t}(y) \in B_1$. Then, for any $N \geq 1$,

$$q_i \cdot \rho(y) = -\rho \cdot q_i(y), \quad \forall y \in \mathbb{R}^p, \quad i = 1, \dots, N$$

so that

$$\tilde{f}_{\Delta t}^{(N)} \cdot \rho(y) = -\rho \cdot \tilde{f}_{\Delta t}^{(N)}(y) \quad \forall y \in \mathbb{R}^p.$$

Proof. For induction assume the modified equations up to index N are reversible, that is

$$\tilde{f}_{\Delta t}^{(i)} \cdot \rho(y) = -\rho \cdot \tilde{f}_{\Delta t}^{(i)}(y), \quad 0 \leq i \leq N, \quad y \in \mathbb{R}^p. \quad (4.19)$$

Note that this is true for $N = 0$ since the modified equation of index zero is the original equation (2.1).

Now suppose the modified equation of index $N + 1$ is not reversible. This implies that

$$q_{N+1} \cdot \rho(u_0) \neq -\rho \cdot q_{N+1}(u_0) \quad (4.20)$$

at some point $u_0 \in \mathbb{R}^p$. By continuity, there exists a ball D around u_0 and some constant $C_{12} > 0$ such that

$$\|q_{N+1} \cdot \rho(y) + \rho \cdot q_{N+1}(y)\|_\infty \geq C_{12}, \quad \forall y \in D. \quad (4.21)$$

Let $B_1 = D$ and $B_2 = \rho(B_1)$, and for $i = 1, 2$ let $S_{i,t}, \bar{S}_{i,\Delta t}^n, \tilde{S}_{i,t,\Delta t}^{(m)} : B_i \rightarrow \mathbb{R}^p$ ($m = 1, \dots, N+1$) be the semigroups generated by (2.1), (2.4) and (3.1), respectively. Given $y = u_0$ let $h, T > 0$ be sufficiently small such that, for any $\Delta t \in [0, h]$, we have $S_{1,t}(y), \tilde{S}_{1,t,\Delta t}^{(N+1)}(y) \in B_1$ for all $t \in [-T, T]$ and $\tilde{S}_{1,\Delta t}^n(y) \in B_1$ for any $n\Delta t \in [0, T]$.

Define a function $\psi : [-T, T] \rightarrow \mathbb{R}^p$ by

$$\psi(t) = \rho \cdot \tilde{S}_{1,t,\Delta t}^{(N+1)}(y) - \tilde{S}_{2,-t,\Delta t}^{(N+1)} \cdot \rho(y).$$

Then differentiating with respect to time gives

$$\frac{d\psi}{dt}(t) = \rho \cdot \tilde{f}_{\Delta t}^{(N+1)}(\tilde{S}_{1,t,\Delta t}^{(N+1)}(y)) + \tilde{f}_{\Delta t}^{(N+1)}(\tilde{S}_{2,-t,\Delta t}^{(N+1)} \cdot \rho(y)),$$

and by the definition of $\psi(t)$ we obtain

$$\frac{d\psi}{dt}(t) = \rho \cdot \tilde{f}_{\Delta t}^{(N+1)}(\tilde{S}_{1,t,\Delta t}^{(N+1)}(y)) + \tilde{f}_{\Delta t}^{(N+1)}(\rho \cdot \tilde{S}_{1,t,\Delta t}^{(N+1)}(y) - \psi(t)). \quad (4.22)$$

Using (3.2), (4.13) and (4.19) we have

$$\rho \cdot \tilde{f}_{\Delta t}^{(N+1)}(v) = \Delta t^{r+N} [\rho \cdot q_{N+1}(v) + q_{N+1} \cdot \rho(v)] - \tilde{f}_{\Delta t}^{(N+1)} \cdot \rho(v)$$

and substituting into (4.22) yields

$$\begin{aligned} \frac{d\psi}{dt}(t) &= \Delta t^{r+N} [\rho \cdot q_{N+1}(\tilde{S}_{1,t,\Delta t}^{(N+1)}(y)) + q_{N+1} \cdot \rho(\tilde{S}_{1,t,\Delta t}^{(N+1)}(y))] \\ &\quad + [\tilde{f}_{\Delta t}^{(N+1)}(\rho \cdot \tilde{S}_{1,t,\Delta t}^{(N+1)}(y) - \psi(t)) - \tilde{f}_{\Delta t}^{(N+1)} \cdot \rho(\tilde{S}_{1,t,\Delta t}^{(N+1)}(y))]. \end{aligned}$$

This expression may be written as

$$\frac{d\psi}{dt}(t) = \Delta t^{r+N} [\rho \cdot q_{N+1}(\tilde{S}_{1,t,\Delta t}^{(N+1)}(y)) + q_{N+1} \cdot \rho(\tilde{S}_{1,t,\Delta t}^{(N+1)}(y))] + A(t)\psi(t),$$

where $A(t)$ is obtained from the integral form of the mean value theorem for vector-valued functions. If $B(t)$ is the fundamental matrix solving

$$\frac{dB}{dt}(t) = A(t)B(t), \quad B(0) = I,$$

then

$$\psi(t) = \Delta t^{r+N} B(t) \int_0^t B(s)^{-1} [\rho \cdot q_{N+1}(\tilde{S}_{1,s,\Delta t}^{(N+1)}(y)) + q_{N+1} \cdot \rho(\tilde{S}_{1,s,\Delta t}^{(N+1)}(y))] ds.$$

Now assume for contradiction that

$$B(t) \int_0^t B(s)^{-1} [\rho \cdot q_{N+1}(S_{1,s}(y)) + q_{N+1} \cdot \rho(S_{1,s}(y))] ds = 0$$

for all $t \in [0, T]$, with T sufficiently small. Since $B(t)$ is close to the identity for small t , a straightforward continuity argument shows that

$$\rho \cdot q_{N+1}(S_{1,s}(y)) + q_{N+1} \cdot \rho(S_{1,s}(y)) = 0, \quad \forall s \in [0, T].$$

This contradicts (4.21). Thus there are $\tau \in [0, T]$ and $\delta > 0$ such that

$$\left\| B(\tau) \int_0^\tau B(s)^{-1} [\rho \cdot q_{N+1}(S_{1,s}(y)) + q_{N+1} \cdot \rho(S_{1,s}(y))] ds \right\| \geq 2\delta.$$

In view of (3.9) we have

$$\left\| B(\tau) \int_0^\tau B(s)^{-1} [\rho \cdot q_{N+1}(\tilde{S}_{1,s,\Delta t}^{(N+1)}(y)) + q_{N+1} \cdot \rho(\tilde{S}_{1,s,\Delta t}^{(N+1)}(y))] ds \right\| \geq \delta$$

for all $\Delta t \in [0, h]$ with h sufficiently small, and thus

$$\|\psi(\tau)\| \geq \Delta t^{r+N} \delta, \quad \forall \Delta t \in [0, h]. \quad (4.23)$$

Now define

$$\Phi(\tau) = \tilde{S}_{2,\tau,\Delta t}^{(N+1)} \cdot \rho \cdot \tilde{S}_{1,\tau,\Delta t}^{(N+1)}(y) - \rho(y)$$

so that

$$\psi(\tau) = \tilde{S}_{2,-\tau,\Delta t}^{(N+1)}(\Phi(\tau) + \rho(y)) - \tilde{S}_{2,-\tau,\Delta t}^{(N+1)} \cdot \rho(y).$$

Since $\tilde{S}_{2,-t,\Delta t}^{(N+1)}$ is close to the identity mapping for t small, the mean value theorem and (4.23) yield

$$\|\Phi(\tau)\| \geq \frac{\Delta t^{r+N} \delta}{2}, \quad \forall \Delta t \in [0, h], \quad (4.24)$$

after reducing T if necessary. However, since the numerical method is reversible, we have

$$\Phi(\tau) = \tilde{S}_{2,\tau,\Delta t}^{(N+1)} \cdot \rho \cdot \tilde{S}_{1,\tau,\Delta t}^{(N+1)}(y) - \tilde{S}_{2,\Delta t}^n \cdot \rho \cdot \tilde{S}_{1,\Delta t}^n(y)$$

for $n\Delta t = \tau$ with $\Delta t \in [0, h]$. By Theorem 3.2 we deduce that

$$\|\Phi(\tau)\| \leq C_{13} \Delta t^{r+N+1}, \quad \forall \Delta t \in [0, h]. \quad (4.25)$$

Since (4.24) and (4.25) give a contradiction, we deduce that (4.20) cannot hold. This completes the inductive step. \square

4.3 Fixed points of the modified semigroup

The next theorem shows that any numerical method which inherits the equilibria of (2.1) as fixed points has modified equations which inherit these as equilibria.

THEOREM 4.4 Suppose that there is a point $y^* \in \mathbb{R}^p$ such that $f(y^*) = 0$ and $\tilde{S}_{\Delta t}(y^*) = y^*$ for all $\Delta t > 0$. Then, for any $N \geq 1$, we have $q_i(y^*) = 0$ ($i = 1, \dots, N$) and hence

$$\tilde{f}_{\Delta t}^{(N)}(y^*) = 0, \quad \forall \Delta t > 0.$$

Proof. For induction assume that the modified vector fields up to index N have the property $\tilde{f}_{\Delta t}^{(N)}(y^*) = 0$ for all $\Delta t > 0$. Note that this is true for $N = 0$ since $\tilde{f}_{\Delta t}^{(0)} = f$. Now suppose $\tilde{f}_{\Delta t}^{(N+1)}(y^*) \neq 0$ for some $\Delta t > 0$. Since

$$\tilde{f}_{\Delta t}^{(N+1)}(v) = \tilde{f}_{\Delta t}^{(N)}(v) + \Delta t^{r+N} q_{N+1}(v)$$

this implies that $q_{N+1}(y^*) \neq 0$, and thus $\tilde{f}_{\Delta t}^{(N+1)}(y^*) \neq 0$ for all $\Delta t > 0$.

To obtain a contradiction, let D be a closed ball around y^* and let $h, T > 0$ be sufficiently small such that $\tilde{S}_{t, \Delta t}^{(N+1)}(y^*) \in D$ for all $t \in [0, T]$ and $\Delta t \in [0, h]$. Introduce a maximum Lipschitz constant in D by

$$L := \max\{\text{Lip}[f], \text{Lip}[q_1], \dots, \text{Lip}[q_{N+1}]\},$$

and let $y(t) = \tilde{S}_{t, \Delta t}^{(N+1)}(y^*)$ and $z(t) = y(t) - y^*$.

Now, since $\tilde{S}_{\Delta t}^n(y^*) = y^*$, Theorem 3.2 shows that

$$\|z(t)\| \leq C_8 \Delta t^{r+N+1} \quad (4.26)$$

for any $t = n\Delta t \in [0, T]$ with $\Delta t \in [0, h]$. By the definition of L we thus have

$$\|f(y(t)) - f(y^*)\| \leq L\|z(t)\| \leq LC_8 \Delta t^{r+N+1} \quad (4.27)$$

and

$$\|q_i(y(t)) - q_i(y^*)\| \leq LC_8 \Delta t^{r+N+1}, \quad 1 \leq i \leq N+1. \quad (4.28)$$

Using the fact that

$$\tilde{f}_{\Delta t}^{(N+1)}(y(t)) - \tilde{f}_{\Delta t}^{(N+1)}(y^*) = f(y(t)) - f(y^*) + \sum_{i=1}^N \Delta t^{r+i-1} [q_i(y(t)) - q_i(y^*)],$$

together with (4.27) and (4.28), it follows that there is a constant C_{14} such that

$$\|\tilde{f}_{\Delta t}^{(N+1)}(y(t)) - \tilde{f}_{\Delta t}^{(N+1)}(y^*)\| \leq C_{14} \Delta t^{r+N+1}. \quad (4.29)$$

Finally, since $dz(t)/dt = \tilde{f}_{\Delta t}^{(N+1)}(y(t))$ and $z(0) = 0$, we have from (4.29) that

$$\begin{aligned} \|z(T)\| &= \left\| \int_0^T \tilde{f}_{\Delta t}^{(N+1)}(y(t)) dt \right\| \\ &= \left\| \int_0^T \tilde{f}_{\Delta t}^{(N+1)}(y^*) + [\tilde{f}_{\Delta t}^{(N+1)}(y(t)) - \tilde{f}_{\Delta t}^{(N+1)}(y^*)] dt \right\| \\ &\geq \left\| \int_0^T \tilde{f}_{\Delta t}^{(N+1)}(y^*) dt \right\| - \int_0^T \|\tilde{f}_{\Delta t}^{(N+1)}(y(t)) - \tilde{f}_{\Delta t}^{(N+1)}(y^*)\| dt \\ &\geq T \Delta t^{r+N} \|q_{N+1}(y^*)\| - TC_{14} \Delta t^{r+N+1}. \end{aligned}$$

Hence, after reducing h if necessary, there exists a constant $C_{15} > 0$ such that

$$\|z(T)\| \geq C_{15} \Delta t^{r+N}.$$

This contradicts (4.26), and thus we must have $\tilde{f}_{\Delta t}^{(N+1)}(y^*) = 0$ for all $\Delta t > 0$. \square

4.4 Symplecticity of the modified semigroup

Suppose now that the dimension p is even, say $p = 2m$, and the vector field $f : \mathbb{R}^p \rightarrow \mathbb{R}^p$ is the Hamiltonian with respect to the canonical symplectic structure; that is,

$$f(u) = J \nabla H(u) \quad (4.30)$$

for some smooth function $H : \mathbb{R}^p \rightarrow \mathbb{R}$, where $J \in \mathbb{R}^{p \times p}$ is of the form

$$J = \begin{pmatrix} O_m & I_m \\ -I_m & O_m \end{pmatrix}. \quad (4.31)$$

Here O_m and I_m denote the zero and identity matrices in $\mathbb{R}^{m \times m}$, respectively. We now show that when the numerical approximation scheme generates a symplectic semigroup, the associated modified equations share the same property. This result is already well known (see Mackay (1992), Sanz-Serna (1992), Reich (1993, 1996), Benettin & Giorgilli (1994), Calvo, Murua & Sanz-Serna (1994) and Hairer (1994) for example), and we simply provide a new proof using the unifying technique of this article.

THEOREM 4.5 Suppose the underlying system (2.1) and the approximation scheme (2.4) both generate symplectic semigroups. Then the associated modified equation (3.1) of index N generates a symplectic semigroup for any integer $N \geq 1$. Thus the modified equation (3.1) has the form

$$\frac{dv}{dt} = J \nabla [H(v) + \Delta t^r Q^{(N)}(v; \Delta t)]. \quad (4.32)$$

Proof. For any $u_0 \in \mathbb{R}^p$ let $S : B \times [0, T] \rightarrow \mathbb{R}^p$ denote the local semigroup generated by (4.30) with image defined as

$$\mathcal{I}m(S) = \{u \in \mathbb{R}^p \mid u = S(U, t), \quad U \in B, \quad t \in [0, T]\}, \quad (4.33)$$

and let $df(u) \in \mathbb{R}^{p \times p}$ denote the derivative at any point $u \in \mathbb{R}^p$ of the Hamiltonian vector field $f : \mathbb{R}^p \rightarrow \mathbb{R}^p$ given by (4.30).

Since the vector field we are considering is Hamiltonian, the mapping $S_t : B \rightarrow \mathbb{R}^p$ is symplectic for each $t \in [0, T]$ in the sense that

$$dS_t(U)^T J dS_t(U) = J \quad (4.34)$$

for all $U \in B$. An equivalent statement to (4.34) is that $df(u) \in \mathbb{R}^{p \times p}$ is *infinitesimally symplectic* for each $u \in \mathcal{I}m(S)$; that is,

$$df(u)^T J + J df(u) = 0 \quad (4.35)$$

for all $u \in \mathcal{I}m(S)$. That (4.35) holds follows from (4.30) since $R^T J + J R = 0$ for any matrix $R = J A$ where $A^T = A$.

For induction assume that, given any $u_0 \in \mathbb{R}^p$, the modified equation of index N generates a local evolution semigroup $\tilde{S}^{(N)}$ which is symplectic, and note that this is true for $N = 0$, the unperturbed equation (2.1). Then, for any $\Delta t \in [0, h]$ and $t \in [0, T]$, the mapping $\tilde{S}_{t, \Delta t}^{(N)} : B \rightarrow \mathbb{R}^p$ satisfies

$$d\tilde{S}_{t, \Delta t}^{(N)}(U)^T J d\tilde{S}_{t, \Delta t}^{(N)}(U) = J \quad (4.36)$$

for all $U \in B$. Equivalently, for any $\Delta t \in [0, h]$, the modified vector field $\tilde{f}_{\Delta t}^{(N)}$ satisfies

$$d\tilde{f}_{\Delta t}^{(N)}(u)^T J + J d\tilde{f}_{\Delta t}^{(N)}(u) = 0 \quad (4.37)$$

for all $u \in \text{Im}(\tilde{S}_{\Delta t}^{(N)})$.

Now assume, for contradiction, that the modified equation of index $N + 1$, which is of the form

$$\frac{dv}{dt} = \tilde{f}_{\Delta t}^{(N+1)}(v) = \tilde{f}_{\Delta t}^{(N)}(v) + \Delta t^{r+N} q_{N+1}(v), \quad (4.38)$$

generates a semigroup which is not symplectic. Then there exists $u_0 \in \mathbb{R}^p$ such that

$$\Phi(u_0) \neq 0 \quad (4.39)$$

where

$$\Phi(u) := dq_{N+1}(u)^T J + J dq_{N+1}(u). \quad (4.40)$$

Otherwise, $d\tilde{f}_{\Delta t}^{(N+1)}(u)$ would be infinitesimally symplectic for all $u \in \mathbb{R}^p$ and, for any $u_0 \in \mathbb{R}^p$, the local semigroup $\tilde{S}_{t, \Delta t}^{(N+1)}$ would be symplectic.

Let $C_{16} = C_{16}(u_0) > 0$ be such that

$$\|\Phi(u_0)\| = 2C_{16} > 0. \quad (4.41)$$

Since q_{N+1} is C^∞ smooth, it follows by continuity that there is a closed ball D at u_0 such that

$$\|\Phi(U)\| \geq C_{16} > 0, \quad \forall U \in D. \quad (4.42)$$

Given $U \in D \cap B$ let $h, T > 0$ be such that, for any $\Delta t \in [0, h]$, the semigroups satisfy $S(U, t), \tilde{S}_{\Delta t}^{(N+1)}(U, t) \in D$ for all $t \in [0, T]$ and $U_n = \tilde{S}_{\Delta t}^n(U) \in D$ for all $n\Delta t \in [0, T]$. Furthermore, let $V(t) = d\tilde{S}_{\Delta t}^{(N+1)}(U, t)$ and note that $V(t)$ satisfies the matrix equation

$$\frac{dV}{dt} = d\tilde{f}_{\Delta t}^{(N+1)}(u_{\Delta t}(t))V, \quad V(0) = I_p \quad (4.43)$$

where we have used the notation $u_{\Delta t}(t) = \tilde{S}_{\Delta t}^{(N+1)}(U, t)$. Using this relation together with (4.37) and (4.38) it follows that, for any $\Delta t \in [0, h]$, we have

$$\begin{aligned} & \left. \frac{\partial}{\partial \tau} \right|_{\tau=t} (d\tilde{S}_{\Delta t}^{(N+1)}(U, \tau)^T J d\tilde{S}_{\Delta t}^{(N+1)}(U, \tau)) \\ &= d\tilde{S}_{\Delta t}^{(N+1)}(U, t)^T (d\tilde{f}_{\Delta t}^{(N+1)}(u_{\Delta t}(t))^T J + J d\tilde{f}_{\Delta t}^{(N+1)}(u_{\Delta t}(t))) d\tilde{S}_{\Delta t}^{(N+1)}(U, t) \\ &= \Delta t^{r+N} d\tilde{S}_{\Delta t}^{(N+1)}(U, t)^T (dq_{N+1}(u_{\Delta t}(t))^T J + J dq_{N+1}(u_{\Delta t}(t))) d\tilde{S}_{\Delta t}^{(N+1)}(U, t). \end{aligned} \quad (4.44)$$

For convenience we introduce the notation

$$\mathcal{A}_{\Delta t}(U, t) = d\tilde{S}_{\Delta t}^{(N+1)}(U, t)^T J d\tilde{S}_{\Delta t}^{(N+1)}(U, t) \quad (4.45)$$

and for future reference we note that $\mathcal{A}_{\Delta t}(U, 0) = J$. By (4.44) and (3.9) we deduce that

$$\left. \frac{\partial}{\partial \tau} \right|_{\tau=t} \mathcal{A}_{\Delta t}(U, \tau) = \Delta t^{r+N} dS(U, t)^T \Phi(S(U, t)) dS(U, t) + \hat{r}(U, t) \quad (4.46)$$

where

$$\|\hat{r}(U, t)\| \leq K(U, t) \Delta t^{2r+N} \quad (4.47)$$

for some $K = K(U, t)$ independent of Δt . Hence, for any $t \in [0, T]$ we have

$$\mathcal{A}_{\Delta t}(U, t) = J + \Delta t^{r+N} \int_0^t dS(U, s)^T \Phi(S(U, s)) dS(U, s) ds + \int_0^t \hat{r}(U, s) ds. \quad (4.48)$$

If we assume that

$$\int_0^t dS(U, s)^T \Phi(S(U, s)) dS(U, s) ds = 0 \quad (4.49)$$

for all $t \in [0, T]$, then standard continuity arguments lead to the conclusion

$$dS(U, s)^T \Phi(S(U, s)) dS(U, s) = 0 \quad (4.50)$$

for all $s \in [0, T]$. By the invertibility of $dS(U, s)$, we then deduce that

$$\Phi(S(U, s)) = 0 \quad (4.51)$$

for all $s \in [0, T]$. However, since $S(U, s) \in D$ for all $s \in [0, T]$, we have a contradiction with (4.42), and so there must exist $\tau \in [0, T]$ (independent of Δt) and $\delta > 0$ such that

$$\left\| \int_0^\tau dS(U, s)^T \Phi(S(U, s)) dS(U, s) ds \right\| \geq \delta. \quad (4.52)$$

Using the above result in (4.48) gives

$$\|\mathcal{A}_{\Delta t}(U, \tau) - J\| \geq \Delta t^{r+N} \delta - \left\| \int_0^\tau \hat{r}(U, s) ds \right\|, \quad (4.53)$$

and by (4.47) it follows that

$$\|\mathcal{A}_{\Delta t}(U, \tau) - J\| \geq \Delta t^{r+N} \delta / 2, \quad (4.54)$$

for all $\Delta t \in [0, h]$, possibly by further reduction of h .

Now, using the identity

$$\begin{aligned} & d\tilde{S}_{\Delta t}^{(N+1)}(U, t)^T J d\tilde{S}_{\Delta t}^{(N+1)}(U, t) - d\bar{S}_{\Delta t}^n(U)^T J d\bar{S}_{\Delta t}^n(U) \\ &= \frac{1}{2} (d\tilde{S}_{\Delta t}^{(N+1)}(U, t) + d\bar{S}_{\Delta t}^n(U))^T J (d\tilde{S}_{\Delta t}^{(N+1)}(U, t) - d\bar{S}_{\Delta t}^n(U)) \\ & \quad + \frac{1}{2} (d\tilde{S}_{\Delta t}^{(N+1)}(U, t) - d\bar{S}_{\Delta t}^n(U))^T J (d\tilde{S}_{\Delta t}^{(N+1)}(U, t) + d\bar{S}_{\Delta t}^n(U)) \end{aligned} \quad (4.55)$$

we obtain the bound

$$\begin{aligned} & \left\| d\tilde{S}_{\Delta t}^{(N+1)}(U, \tau)^T J d\tilde{S}_{\Delta t}^{(N+1)}(U, \tau) - d\bar{S}_{\Delta t}^n(U)^T J d\bar{S}_{\Delta t}^n(U) \right\| \\ & \leq (C_2 + C_7) C_{17} \left\| d\tilde{S}_{\Delta t}^{(N+1)}(U, \tau) - d\bar{S}_{\Delta t}^n(U) \right\| \\ & \leq (C_2 + C_7) C_{17} C_8 \Delta t^{r+N+1} \end{aligned} \quad (4.56)$$

for all $\Delta t = \tau/n$ with $n \geq n^*$, where n^* is any positive integer such that $\tau/n^* \in [0, h]$. The above expression follows from the bounds in (2.9), Theorem 3.2 and (4.1), and the

notation $C_{17} = \|J\| > 0$. By the definition of $\mathcal{A}_{\Delta t}$, and the fact that the numerical scheme generates a symplectic semigroup, we have

$$\|\mathcal{A}_{\Delta t}(U, \tau) - J\| \leq (C_2 + C_7)C_{17}C_8\Delta t^{r+N+1}$$

for all $\Delta t = \tau/n$ with $n \geq n^*$. This contradicts (4.54). Hence $dq_{N+1}(u)$ must be infinitesimally symplectic for all $u \in \mathbb{R}^p$, which implies that the modified equation of index $N + 1$ generates a symplectic semigroup. The case $N = 0$ holds since it gives the equation (2.1), and thus by induction the modified equation of index N generates a symplectic semigroup for any integer $N \geq 1$. By the results of Dragt & Finn (1976) we deduce that, for any integer $N \geq 1$, the modified equation is Hamiltonian as claimed. \square

4.5 Volume preservation and the modified semigroup

All equations (2.1) with f divergence-free have semigroups which preserve phase volume. Since Hamiltonian vector fields are divergence-free this result holds for all Hamiltonian problems; indeed, for dimension $p = 2$, it is equivalent to symplecticity of the semigroup. It is possible to construct numerical methods that automatically inherit the property of volume preservation; see, for example, Feng & Wang (1994), Feng & Shang (1995), Shang (1994), Quispel (1995) and Suris (1996) (and earlier references cited therein to the Soviet literature). We now give a result that applies to the modified equations of one-step volume-preserving methods; this result may also be proved by the techniques in Reich (1996).

THEOREM 4.6 If the system and method preserve volume, then so do all modified equations. More precisely, assume that for any compact set B we have

$$\int_{S_t(B)} dv = \int_{\tilde{S}_{\Delta t}^n(B)} dv = \int_B dv$$

for all $t = n\Delta t \in [0, T]$ and $\Delta t \in [0, h]$ where $h, T > 0$ are sufficiently small. Then, for any $N \geq 1$, we have $\nabla \cdot \tilde{f}_{\Delta t}^{(N)}(y) = 0$ for all $y \in \mathbb{R}^p$ and $\Delta t > 0$.

Proof. For induction assume the modified vector fields up to index N are divergence-free, that is

$$\nabla \cdot \tilde{f}_{\Delta t}^{(i)}(y) = 0, \quad 0 \leq i \leq N, \quad y \in \mathbb{R}^p, \quad \Delta t > 0. \tag{4.57}$$

Since (2.1) preserves phase volume if and only if $\nabla \cdot f \equiv 0$, we note that the inductive hypothesis is true for $N = 0$. Now suppose the modified vector field of index $N + 1$ is not divergence-free for some $\Delta t > 0$. This implies that

$$\nabla \cdot q_{N+1}(u_0) \neq 0 \tag{4.58}$$

for some $u_0 \in \mathbb{R}^p$.

Since $\tilde{S}_{t,\Delta t}^{(N+1)}$ is close to the identity for t small and q_{N+1} is smooth, we deduce from (4.58) that there are numbers $\gamma, h, T > 0$ sufficiently small and a number $\delta > 0$ such that

$$|\nabla \cdot q_{N+1}(y)| \geq \delta, \quad \forall y \in \bigcup_{\substack{t \in [0, T] \\ \Delta t \in [0, h]}} \tilde{S}_{t,\Delta t}^{(N+1)}(B, t),$$

where B is the closed ball at u_0 of radius γ . Also, if we let

$$V^{(N+1)}(t) = \int_{\tilde{S}_{t,\Delta t}^{(N+1)}(B)} dv,$$

then

$$\begin{aligned} \frac{dV^{(N+1)}}{dt} &= \int_{\tilde{S}_{t,\Delta t}^{(N+1)}(B)} \nabla \cdot \tilde{f}_{\Delta t}^{(N+1)}(v) dv \\ &= \int_{\tilde{S}_{t,\Delta t}^{(N+1)}(B)} \Delta t^{r+N} \nabla \cdot q_{N+1}(v) dv. \end{aligned}$$

Assuming first that $\nabla \cdot q_{N+1}(u_0) > 0$, we obtain

$$\frac{dV^{(N+1)}}{dt} \geq \delta \Delta t^{r+N} V^{(N+1)}(t), \quad t \in [0, T], \quad \Delta t \in [0, h].$$

This inequality leads to the bounds

$$V^{(N+1)}(t) \geq \exp(\delta \Delta t^{r+N} t) V^{(N+1)}(0) \geq (1 + \delta \Delta t^{r+N} t) V^{(N+1)}(0),$$

and hence

$$|V^{(N+1)}(t) - V^{(N+1)}(0)| \geq \delta \Delta t^{r+N} t V^{(N+1)}(0), \quad t \in [0, T], \quad \Delta t \in [0, h]. \quad (4.59)$$

Alternatively, in the case where $\nabla \cdot q_{N+1}(u_0) < 0$, we have

$$\frac{dV^{(N+1)}}{dt} \leq -\delta \Delta t^{r+N} V^{(N+1)}(t), \quad t \in [0, T], \quad \Delta t \in [0, h],$$

and

$$V^{(N+1)}(t) \leq \exp(-\delta \Delta t^{r+N} t) V^{(N+1)}(0) \leq \left(1 - \frac{\delta \Delta t^{r+N} t}{2}\right) V^{(N+1)}(0), \quad (4.60)$$

for all $t \in [0, T]$ and $\Delta t \in [0, h]$, after reducing T if necessary. Thus, regardless of the sign of $\nabla \cdot q_{N+1}(u_0)$, we deduce from (4.59) and (4.60) that

$$|V^{(N+1)}(t) - V^{(N+1)}(0)| \geq \frac{\delta \Delta t^{r+N} t V^{(N+1)}(0)}{2}, \quad t \in [0, T], \quad \Delta t \in [0, h]. \quad (4.61)$$

To obtain a contradiction, let

$$V_{\Delta t}(n) = \int_{\tilde{S}_{\Delta t}^n(B)} dv.$$

Then, since the numerical method conserves volume, we have $V_{\Delta t}(n) = V^{(N+1)}(0)$, and hence

$$|V^{(N+1)}(t) - V^{(N+1)}(0)| = |V^{(N+1)}(t) - V_{\Delta t}(n)| \quad (4.62)$$

for any $t = n\Delta t \in [0, T]$ with $\Delta t \in [0, h]$. Using the relations

$$V^{(N+1)}(t) = \int_{\tilde{S}_{t,\Delta t}^{(N+1)}(B)} dv = \int_B \det[d\tilde{S}_{t,\Delta t}^{(N+1)}(v)] dv$$

and

$$V_{\Delta t}(n) = \int_{\tilde{S}_{\Delta t}^n(B)} dv = \int_B \det[d\tilde{S}_{\Delta t}^n(v)] dv,$$

together with (4.62) and the derivative bounds in Theorem 3.2, we find that there exists a constant C_{18} such that

$$|V^{(N+1)}(t) - V^{(N+1)}(0)| \leq C_{18}\Delta t^{r+N+1}$$

for any $t = n\Delta t \in [0, T]$ with $\Delta t \in [0, h]$. This contradicts (4.61) and thus we must have $\nabla \cdot \tilde{f}_{\Delta t}^{(N+1)}(y) = 0$ for all $y \in \mathbb{R}^p$ and $\Delta t > 0$. \square

4.6 A negative example

The results presented in this article fit into a general framework that may be summarized as follows: *if the differential equations have a certain structural property and the numerical method has an analogous structural property for $\Delta t > 0$, then all modified equations share this property for $\Delta t > 0$* . However, the results proved are strongly tied to structural properties associated with certain subspaces of $\mathcal{V}(\mathbb{R}^p)$, the infinite-dimensional Lie algebra of smooth vector fields on \mathbb{R}^p , as outlined in the introduction. Vector fields sharing a particular structural property do not necessarily form a subspace of $\mathcal{V}(\mathbb{R}^p)$, and it is natural to ask about the properties enjoyed by modified equations in this case. As the next example shows, modified equations do not generically inherit structural properties shared by the numerical method and underlying system.

Consider the scalar problem with $f(u) = -u^3$. This is an example of a gradient system with Lyapunov function $F(u) = u^4/4$. For this problem we have

$$\frac{d}{dt}F(S_t(u)) = -f(S_t(u))^2. \quad (4.63)$$

Hence, the Lyapunov function F decreases along every non-constant solution trajectory. It follows that

$$S_t(u) \rightarrow 0, \quad \text{as } t \rightarrow \infty, \quad \forall u \in \mathbb{R}^p. \quad (4.64)$$

For the implicit Euler method $U_{n+1} = U_n + \Delta t f(U_{n+1})$ it is readily shown that, for $f(u) = -u^3$,

$$\frac{F(U_{n+1}) - F(U_n)}{\Delta t} \leq -\left(\frac{U_{n+1} - U_n}{\Delta t}\right)^2.$$

It follows that, independently of Δt ,

$$\tilde{S}_{\Delta t}^n(u) \rightarrow 0, \quad \text{as } n \rightarrow \infty, \quad \forall u \in \mathbb{R}^p, \quad (4.65)$$

which is the discrete analogue of (4.64).

The first modified equation associated with the implicit Euler method on this problem has a vector field of the form

$$\tilde{f}_{\Delta t}^{(1)}(u) = -u^3 + \frac{3\Delta t}{2}u^5.$$

In this case,

$$\tilde{S}_{t,\Delta t}^{(1)}(u) \rightarrow 0, \quad \text{as } t \rightarrow \infty, \quad \text{for } |u| < \sqrt{\frac{2}{3\Delta t}},$$

and

$$|\tilde{S}_{t,\Delta t}^{(1)}(u)| \rightarrow \infty, \quad \text{as } t \rightarrow \infty, \quad \text{for } |u| > \sqrt{\frac{2}{3\Delta t}}.$$

Hence, the first modified equation does not have the property of global convergence to the origin, even though this property is shared by the original equation and the discretization.

REFERENCES

- AUERBACH, S. P. & FRIEDMAN, A. 1991 Long-time behaviour of numerically computed orbits: small and intermediate time-step analysis of one-dimensional systems. *J. Comput. Phys.* **93**, 189–223.
- BENETTIN, G. & GIORGILLI, A. 1994 On the Hamiltonian interpolation of near-to-the-identity symplectic mappings with application to symplectic integrators. *J. Stat. Phys.* **74**, 1117–1143.
- BEYN, W.-J. 1991 Numerical methods for dynamical systems. *Advances in Numerical Analysis; Volume 1*. Oxford: Clarendon.
- CALVO, M. P., ISERLES, A., & ZANNA, A. 1996 Runge–Kutta methods for orthogonal and isospectral flows. *Numerical Analysis Report NA08*, DAMTP, University of Cambridge.
- CALVO, M. P., MURUA, A., & SANZ-SERNA, J. M. 1994 Modified equations for ODEs. *Contemp. Math.* **172**, 63–74.
- DRAGT, A. J. & FINN, J. M. 1976 Lie series and invariant functions for analytic symplectic maps. *J. Math. Phys.* **17**, 2215–2227.
- FENG KANG & SHANG ZAI-JIU 1995 Volume-preserving difference schemes for source free dynamical systems. *Numer. Math.* **71**, 451–463.
- FENG KANG & WANG DAO-LIU 1994 Dynamical systems and geometric construction of algorithms. *Contemporary Mathematics 163* (Z. Shi and C. Yang, eds).
- FIEDLER, B. & SCHEURLE, J. 1996 Discretization of homoclinic orbits, rapid forcing and ‘invisible’ chaos. *Memoirs of the AMS*, Providence, RI.
- GOLUB, G. H. & Van Loan, C. F. 1996 *Matrix Computations* 3rd edn. Johns Hopkins University Press.
- GRIFFITHS, D. F. & SANZ-SERNA, J. M. 1986 On the scope of the method of modified equations. *SIAM J. Sci. Stat. Comput.* **7**, 994–1008.
- HAIRER, E. 1994 Backward analysis of numerical integrators and symplectic methods. *Ann. Numer. Math.* **1**, 107–132.
- HAIRER, E. 1997 Variable time step integration with symplectic methods. *Appl. Numer. Math.* **25**, 219–227.
- HAIRER, E. & LUBICH, CH. 1997 The life-span of backward error analysis for numerical integrators. *Numer. Math.* **76**, 441–462.
- HAIRER, E. & STOFFER, D. 1997 Reversible long-term integration with variable step-sizes. *SIAM J. Sci. Comput.* **18**, 257–269.

- MACKAY, R. S. 1992 Some aspects of the dynamics and numerics of Hamiltonian systems. *Proc. IMA Conf. on The Dynamics of Numerics and the Numerics of Dynamics, 1990* (D. Broomhead and A. Iserles, eds). Cambridge: Cambridge University Press.
- NEISHTADT, A. I. 1984 The separation of motions in systems with rapidly rotating phase. *J. Appl. Math. Mech.* **48**, 133–139.
- QUISPÉL, G. R. W. 1995 Volume preserving integrators. *Phys. Lett.* **206A**, 26–30.
- REDDIEN, G. W. 1995 On the stability of numerical methods of Hopf points using backward error analysis. *Computing* **55**, 163–180.
- REICH, S. 1993 Numerical integration of the generalized Euler equations. *Technical Report 93-20*, Department of Computer Science, University of British Columbia.
- REICH, S. 1996 Backward error analysis for numerical integrators. *Preprint SC 96-21*, Konrad-Zuse-Zentrum für Informationstechnik, Berlin.
- SANZ-SERNA, J. M. 1992 Symplectic integrators for Hamiltonian problems: an overview. *Acta Numerica 1992*, vol. 1, 243–286.
- SANZ-SERNA, J. M. & CALVO, M. P. 1994 *Numerical Hamiltonian Problems*. London: Chapman and Hall.
- SANZ-SERNA, J. M. & MURUA, A. 1997 *NSF-CBMS Lecture, Golden, Colorado, 1997*. Philadelphia: SIAM. To appear.
- SHANG ZAI-JIU 1994 Construction of volume-preserving difference schemes for source free systems via generating functions. *J. Comput. Math.* **12**, 265–272.
- STOFFER, D. 1995 Variable steps for reversible integration methods. *Computing* **55**, 1–22.
- STROGATZ, S. H. 1994 *Nonlinear Dynamics and Chaos*. Reading, MA: Addison-Wesley.
- STUART, A. M. & HUMPHRIES, A. R. 1996 *Dynamical Systems and Numerical Analysis*. Cambridge: Cambridge University Press.
- SURIS, Y. B. 1996 Partitioned Runge–Kutta methods as phase volume preserving integrators. *Phys. Lett.* **220A**, 63–69.
- WARMING, R. F. & HYETT, B. J. 1974 The modified equation approach to the stability and accuracy analysis of finite difference methods. *J. Comput. Phys.* **14**, 159–179.
- YOSHIDA, H. 1993 Recent progress in the theory and application of symplectic integrators. *Celestial Mech. Dyn. Astron.* **56**, 27–43.