

## BACKWARD ERROR ANALYSIS FOR NUMERICAL INTEGRATORS\*

SEBASTIAN REICH†

**Abstract.** Backward error analysis has become an important tool for understanding the long time behavior of numerical integration methods. This is true in particular for the integration of Hamiltonian systems where backward error analysis can be used to show that a symplectic method will conserve energy over exponentially long periods of time. Such results are typically based on two aspects of backward error analysis: (i) It can be shown that the modified vector fields have some qualitative properties which they share with the given problem and (ii) an estimate is given for the difference between the best interpolating vector field and the numerical method. These aspects have been investigated recently, for example, by Benettin and Giorgilli in [*J. Statist. Phys.*, 74 (1994), pp. 1117–1143], by Hairer in [*Ann. Numer. Math.*, 1 (1994), pp. 107–132], and by Hairer and Lubich in [*Numer. Math.*, 76 (1997), pp. 441–462]. In this paper we aim at providing a unifying framework and a simplification of the existing results and corresponding proofs. Our approach to backward error analysis is based on a simple recursive definition of the modified vector fields that does not require explicit Taylor series expansion of the numerical method and the corresponding flow maps as in the above-cited works. As an application we discuss the long time integration of chaotic Hamiltonian systems and the approximation of time averages along numerically computed trajectories.

**Key words.** numerical integrators, differential equations, error analysis, Hamiltonian systems, long time dynamics

**AMS subject classifications.** 65L05, 34A50, 65L70, 70H05

**PII.** S0036142997329797

**1. Introduction.** In this paper, we consider the relationship between solutions to a given system of ordinary differential equations (vector fields)

$$\frac{d}{dt} \mathbf{x} = \mathbf{Z}(\mathbf{x}),$$

numerical approximations

$$\mathbf{x}_{n+1} = \Psi_{\delta t}(\mathbf{x}_n)$$

to them, and solutions to associated modified equations

$$\frac{d}{dt} \mathbf{x} = \tilde{\mathbf{X}}_i(\mathbf{x}; \delta t), \quad (i \geq 1).$$

The vector fields  $\tilde{\mathbf{X}}_i(\delta t)$  are formulated in terms of an asymptotic expansion in the step size  $\delta t$ ; i.e., they are chosen such that the numerical solution can formally be interpreted, with increasing index  $i$ , as the increasingly accurate solution of the modified equation. Previous papers on backward error analysis for differential equations include those by Warming and Hyett [38], Griffiths and Sanz-Serna [12], Beyn [6], Feng [10], Eirola [9], Fiedler and Scheurle [11], and Sanz-Serna [31]. Another early reference to

---

\*Received by the editors November 7, 1997; accepted for publication (in revised form) September 30, 1998; published electronically September 8, 1999.

<http://www.siam.org/journals/sinum/36-5/32979.html>

†Konrad-Zuse-Zentrum, Takustr. 7, D-14195 Berlin, Germany. Current address: Department of Mathematics and Statistics, University of Surrey, Guildford, Surrey GU2 5XH, UK (S.Reich@surrey.ac.uk).

related ideas is by Moser [24] who discusses the approximation of a symplectic map near an equilibrium by the flow map of a Hamiltonian vector field.

More recently, general formulas for the computation of the modified vector fields  $\tilde{\mathbf{X}}_i(\delta t)$  have been derived by Hairer [16]; Calvo, Murua, and Sanz-Serna [7]; Benettin and Giorgilli [5]; and Reich [26]. In papers by Neishtadt [25], Benettin and Giorgilli [5], and Hairer and Lubich [18], the question of closeness of the numerical approximations and the solutions of the modified equations has been addressed. In particular, it has been shown in these papers that the difference can be made exponentially small in the step size  $\delta t$ ; i.e.,

$$(1.1) \quad \|\Phi_{\delta t, \tilde{\mathbf{X}}_{i_*}}(\mathbf{x}) - \Psi_{\delta t}(\mathbf{x})\| \leq c_1 \delta t e^{-c_2/\delta t},$$

provided the vector field  $\mathbf{Z}$  and the numerical one-step method  $\Psi_{\delta t}$  are real analytic [22]. Here  $c_1, c_2 > 0$  are appropriate constants,  $\Phi_{\delta t, \tilde{\mathbf{X}}_{i_*}}$  denotes the time- $\delta t$ -flow map of the vector field  $\tilde{\mathbf{X}}_{i_*}$ , and the index  $i_*(\delta t)$  has been chosen such that the difference is minimized.

Backward error analysis is of utmost importance for an understanding of the qualitative behavior of symplectic methods [32] for Hamiltonian problems. It has been shown by Hairer [16]; Calvo, Murua, and Sanz-Serna [7]; Reich [26]; and Benettin and Giorgilli [5] that for symplectic discretizations, the modified vector fields  $\tilde{\mathbf{X}}_i(\delta t)$  are Hamiltonian. For special cases see also the papers by Auerbach and Friedman [4] and Yoshida [40]. The Hamiltonian structure of the modified equations implies that a symplectic integrator almost preserves the total energy over an exponentially long period of time [25, 23, 5, 18]. Similarly, the adiabatic invariant of a Hamiltonian system with a rapidly rotating phase is also preserved over an exponentially long period of time provided a symplectic method is used [29].

The fact that symplectic methods lead to modified equations that are Hamiltonian is a special instance of the so-called geometric properties of backward error analysis. By this we mean the following: If the vector field  $\mathbf{Z}$  belongs to a certain class of vector fields, like integral preserving or divergence-free vector fields, and the numerical approximation  $\Phi_{\delta t}$  also preserves the corresponding quantities, then the modified vector fields  $\tilde{\mathbf{X}}_i(\delta t)$  will be in the same class as  $\mathbf{Z}$ . Besides symplectic methods, special instances of these geometric aspects have been discussed before. See, for example, the papers by Reich [26]; Hairer and Stoffer [20]; and Gonzalez, Higham, and Stuart [13].

In this paper, we revisit backward error analysis by using a simple recursive scheme for the definition of the modified vector fields  $\tilde{\mathbf{X}}_i(\delta t)$  as first proposed by the author in the technical report [26]. The main advantage of this formulation is that it does not require Taylor series expansions of the numerical one-step method  $\Psi_{\delta t}$  and the flow maps  $\Phi_{\delta t, \tilde{\mathbf{X}}_i}$  in terms of the step size  $\delta t$  as it is used in the papers by Benettin and Giorgilli [5], Hairer [16], and Hairer and Lubich [18]. This in turn allows for a simple characterization<sup>1</sup> of the geometric properties of the modified vector fields  $\tilde{\mathbf{X}}_i(\delta t)$  and a rather simple proof for the exponentially small truncation error (1.1). Our approach is close to the one discussed by Benettin and Giorgilli [5] in the sense that we consider general one-step methods<sup>2</sup> and that we use a “direct” approach.<sup>3</sup>

<sup>1</sup>The general idea can already be found in the report [26].

<sup>2</sup>Hairer and Lubich [18] consider methods that can be represented by P-series [19]. Note that Runge–Kutta and partitioned Runge–Kutta methods fall under this category.

<sup>3</sup>This is in contrast to the “indirect” approach used by Neishtadt [25] where the one-step method is first interpolated by the flow of a time-dependent vector field, and averaging in time is then used to obtain an optimal approximating time-independent vector field.

However, different techniques are used and we will discuss this in more detail in section 4.

In section 5, we consider the numerical integration of a “chaotic” Hamiltonian system by a symplectic method and discuss the approximation of time averages along numerically computed trajectories. We assume that a Poincaré section [14] can be defined and that the corresponding Poincaré section is uniformly hyperbolic. Backward error analysis and the shadowing lemma [33] will be used to show that a numerically computed trajectory stays close to an exact solution over exponentially long periods of time. This and a large deviation theorem [36] allow us to discuss the convergence of long time averages along numerically computed trajectories. The anisotropic Kepler problem [15] will serve us as a numerical illustration. This problem requires the application of a symplectic variable step size method as first discussed by the author in the technical report [27] and independently by Hairer in [17].

**2. The modified vector field recursion.** Let us consider a smooth vector field

$$(2.2) \quad \frac{d}{dt} \mathbf{x} = \mathbf{Z}(\mathbf{x}),$$

$\mathbf{Z} : \mathcal{U} \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$  and its discretization by a one-step method [19]

$$(2.3) \quad \mathbf{x}_{n+1} = \Psi_{\delta t}(\mathbf{x}_n) = \mathbf{x}_n + \delta t \psi(\mathbf{x}_n, \delta t).$$

We assume that  $\Psi_{\delta t} : \mathcal{U} \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$  is a smooth map and a method of order  $p \geq 1$ ; i.e.,

$$\|\Phi_{\delta t, \mathbf{Z}}(\mathbf{x}) - \Psi_{\delta t}(\mathbf{x})\| = \mathcal{O}(\delta t^{p+1})$$

for all  $\mathbf{x} \in \mathcal{U}$  where  $\Phi_{\delta t, \mathbf{Z}}$  is the time- $\delta t$ -flow map of the differential equation (2.2).

As described in the introduction, we look for a family of vector fields  $\tilde{\mathbf{X}}(\delta t)$  such that

$$\Phi_{\delta t, \tilde{\mathbf{X}}(\delta t)} \approx \Psi_{\delta t}$$

or, equivalently,

$$\Phi_{1, \mathbf{X}(\delta t)} \approx \Psi_{\delta t}, \quad \mathbf{X}(\delta t) := \delta t \tilde{\mathbf{X}}(\delta t)$$

for all  $\delta t$  sufficiently small. Here  $\Phi_{1, \mathbf{X}}$  denotes the time-one-flow map of the vector field  $\mathbf{X}(\delta t)$ . The family of *modified vector fields*  $\mathbf{X}(\delta t)$ ,  $\delta t \geq 0$ , is formally defined in terms of an asymptotic expansion in the step size  $\delta t$ , i.e.,

$$\mathbf{X}(\delta t) = \delta t \Delta \mathbf{X}_1 + \delta t^2 \Delta \mathbf{X}_2 + \delta t^3 \Delta \mathbf{X}_3 + \dots$$

The formally infinite sequence of vector fields  $\{\Delta \mathbf{X}_i\}_{i=1, \dots, \infty}$  can be obtained by Taylor series expansion of the one-step method  $\Psi_{\delta t}$ , i.e.,

$$\Psi_{\delta t} = \mathbf{id} + \delta t \Psi_1 + \delta t^2 \Psi_2 + \dots,$$

$\mathbf{id}(\mathbf{x}) = \mathbf{x}$  the identity map, and comparison of this series with the expansion of the time-one-flow map  $\Phi_{1, \mathbf{X}(\delta t)}$  in terms of  $\delta t$ . The vector fields  $\Delta \mathbf{X}_i$  are chosen such that these two series coincide term by term. This is the general approach followed by Benettin and Giorgilli [5] and Hairer [16]. The two papers differ in the way the

Taylor series expansions are written down. But they lead to exactly the same sequence of vector fields  $\{\Delta \mathbf{X}_i\}_{i=1,\dots,\infty}$ . We obviously have  $\Delta \mathbf{X}_1 = \mathbf{Z}$  and  $\Delta \mathbf{X}_i = \mathbf{0}$ ,  $i = 2, \dots, p$ , for a method of order  $p$ .

We now give a recursive definition of the modified vector field  $\mathbf{X}(\delta t)$  that does not require an explicit Taylor series expansion. This recursion was introduced by the author in the technical report [26]. First we formally introduce the “truncated” expansions  $\mathbf{X}_i(\delta t)$ ,  $i = 1, 2, \dots, \infty$ , by means of

$$\mathbf{X}_i(\delta t) = \sum_{j=1}^i \delta t^j \Delta \mathbf{X}_j.$$

We obviously have

$$\mathbf{X}_{i+1}(\delta t) := \mathbf{X}_i(\delta t) + \delta t^{i+1} \Delta \mathbf{X}_{i+1}.$$

Let us assume that  $\mathbf{X}_i(\delta t)$  has been chosen such that the difference between the time-one-flow map of  $\mathbf{X}_i(\delta t)$  and the numerical one-step method  $\Psi_{\delta t}$  is  $\mathcal{O}(\delta t^{i+1})$ . This suggests that we consider the following recursion:

$$(2.4) \quad \mathbf{X}_{i+1}(\delta t) := \mathbf{X}_i(\delta t) + \delta t^{i+1} \Delta \mathbf{X}_{i+1},$$

$$(2.5) \quad \Delta \mathbf{X}_{i+1} := \lim_{\delta t \rightarrow 0} \frac{\Psi_{\delta t} - \Phi_{1, \mathbf{X}_i(\delta t)}}{\delta t^{i+1}}.$$

Indeed, this definition of  $\Delta \mathbf{X}_{i+1}$  implies that  $\mathbf{X}_{i+1}(\delta t)$ , defined by (2.4), generates a time-one-flow map that is  $\mathcal{O}(\delta t^{i+2})$  away from the numerical method  $\Psi_{\delta t}$ . This can be seen from

$$\begin{aligned} \Phi_{1, \mathbf{X}_{i+1}(\delta t)} - \Psi_{\delta t} &= \Phi_{1, \mathbf{X}_i(\delta t)} + \delta t^{i+1} \Delta \mathbf{X}_{i+1} - \Psi_{\delta t} + \mathcal{O}(\delta t^{i+2}) \\ &= \delta t^{i+1} \Delta \mathbf{X}_{i+1} - \delta t^{i+1} \lim_{\delta t \rightarrow 0} \frac{\Psi_{\delta t} - \Phi_{1, \mathbf{X}_i(\delta t)}}{\delta t^{i+1}} + \mathcal{O}(\delta t^{i+2}) \\ &= \mathcal{O}(\delta t^{i+2}). \end{aligned}$$

Thus (2.4) and (2.5) recursively define the modified vector fields  $\mathbf{X}_i(\delta t)$ ,  $i = 1, \dots, \infty$ . The recursion is started with  $\mathbf{X}_1(\delta t) = \delta t \mathbf{Z}$ . The generated sequence  $\{\Delta \mathbf{X}_i\}_{i=1,\dots,\infty}$  is, of course, equivalent to the sequences obtained by using Taylor series expansions as described in [16, 5].

Throughout this paper, we will exclusively work with the recursion (2.4)–(2.5). In section 3, it will be shown that this leads to a simple characterization of the geometric properties of the modified vector fields and, in section 4, explicit estimates for the difference between the time-one-flow map of the modified vector field  $\mathbf{X}_i(\delta t)$  and the numerical method will be given. We like to point out that these results can also be (and have been [16, 5, 18, 20, 13]) derived using an explicit Taylor series expansion of the flow map and the numerical method. However, we feel that the application of the recursion (2.4)–(2.5) leads to a simplification in the presentation of these results.

**3. Geometric properties of backward error analysis.** In this section, we consider differential equations (2.2) whose corresponding vector field  $\mathbf{Z}$  belongs to a certain linear subspace  $\mathfrak{g}$  of the infinite-dimensional Lie algebra<sup>4</sup> of smooth vector fields on  $\mathbb{R}^n$  [21, 1].

<sup>4</sup>The algebraic operation is the Lie bracket  $[\mathbf{X}, \mathbf{Y}]$  of two vector fields  $\mathbf{X}$  and  $\mathbf{Y}$  [3].

ASSUMPTION. Given a linear subspace  $\mathfrak{g}$  of the infinite-dimensional Lie algebra of smooth vector fields on  $\mathbb{R}^n$ , let us assume that there is a corresponding subset  $\mathfrak{G}$  of the infinite-dimensional Frechet manifold [21] of diffeomorphisms on  $\mathbb{R}^n$  such that

$$\mathfrak{g} = T_{\text{id}} \mathfrak{G}.$$

Here  $T_{\text{id}} \mathfrak{G}$  is defined as the set of all vector fields  $\mathbf{X} := \partial_\tau [\Psi_\tau]_{\tau=0}$  for which the one-parametric family of diffeomorphisms  $\Psi_\tau \in \mathfrak{G}$  is smooth in  $\tau$  and  $\Psi_{\tau=0} = \text{id}$ .

For the linear space (Lie algebra) of Hamiltonian vector fields on  $\mathbb{R}^n$  this is, for example, the subset of canonical transformations [1]. An important aspect of those differential equations is that the corresponding flow map  $\Phi_{t,\mathbf{Z}}$  forms a one-parametric subgroup in  $\mathfrak{G}$  [21, 1]. Especially in the context of long-term integration, it is desirable to discretize differential equations of this type in such a way that the corresponding iteration map  $\Psi_{\delta t}$  belongs to the same subset  $\mathfrak{G}$  as  $\Phi_{t,\mathbf{Z}}$ . We will call those integrators *geometric integrators*.

The following result concerning the backward error analysis of geometric integrators has been first stated in the technical report [26] as follows in Theorem 1.

THEOREM 1. Let us assume that the vector field  $\mathbf{Z}$  in

$$\frac{d}{dt} \mathbf{x} = \mathbf{Z}(\mathbf{x})$$

belongs to a linear subspace  $\mathfrak{g}$  of the Lie algebra of all smooth vector fields on  $\mathbb{R}^n$ . Let us assume furthermore that

$$\mathbf{x}_{n+1} = \Psi_{\delta t}(\mathbf{x}_n) = \mathbf{x}_n + \delta t \psi(\mathbf{x}_n, \delta t)$$

is a geometric integrator for this subspace  $\mathfrak{g}$ ; i.e.,  $\Psi_{\delta t} \in \mathfrak{G}$  for all  $\delta t \geq 0$  sufficiently small. Then the perturbed vector fields  $\mathbf{X}_i(\delta t)$ ,  $i = 1, \dots, \infty$ , defined through the recursion (2.4)–(2.5) belong to  $\mathfrak{g}$ , i.e.,

$$\mathbf{X}_i(\delta t) \in \mathfrak{g}.$$

*Proof.* The statement is certainly true for  $\mathbf{X}_1(\delta t) = \delta t \mathbf{Z}$ . Let us assume that it also holds for  $\mathbf{X}_i(\delta t)$ ; i.e.,  $\mathbf{X}_i(\delta t) \in \mathfrak{g}$  for all  $\delta t \geq 0$  sufficiently small. Since

$$\Psi_{\delta t}(\mathbf{x}) = \mathbf{x} + \delta t \psi(\mathbf{x}, \delta t) \in \mathfrak{G}$$

and

$$\Phi_{1,\mathbf{X}_i(\delta t)} \in \mathfrak{G}$$

for all  $\delta t \geq 0$  sufficiently small as well as

$$\Psi_{\delta t=0} = \Phi_{1,\mathbf{X}_i(\delta t=0)} = \text{id},$$

we have

$$\Delta \mathbf{X}_{i+1} = \lim_{\delta t \rightarrow 0} \frac{\Psi_{\delta t} - \Phi_{1,\mathbf{X}_i(\delta t)}}{\delta t^{i+1}} \in T_{\text{id}} \mathfrak{G}$$

and  $\Delta \mathbf{X}_{i+1} \in \mathfrak{g}$ . This implies  $\mathbf{X}_{i+1}(\delta t) \in \mathfrak{g}$  as required.  $\square$

*Remark.* Often the linear subspace  $\mathfrak{g}$  is, in fact, a subalgebra under the Lie bracket [3]

$$(3.6) \quad [\mathbf{X}, \mathbf{Y}] := \frac{\partial}{\partial \mathbf{x}} \mathbf{X} \cdot \mathbf{Y} - \frac{\partial}{\partial \mathbf{x}} \mathbf{Y} \cdot \mathbf{X};$$

i.e.,  $\mathbf{X}, \mathbf{Y} \in \mathfrak{g}$  implies  $[\mathbf{X}, \mathbf{Y}] \in \mathfrak{g}$ . But this property is not needed in Theorem 1.

Let us discuss five examples.

*Example 1.* Consider the subspace  $\mathfrak{g}$  of all vector fields that preserve a particular first integral  $F : \mathbb{R}^n \rightarrow \mathbb{R}$ . In fact, this space is a subalgebra under the Lie bracket (3.6). In other words,

$$(3.7) \quad \partial_{\mathbf{x}} F \cdot \mathbf{X} = 0$$

and

$$(3.8) \quad \partial_{\mathbf{x}} F \cdot \mathbf{Y} = 0$$

imply that

$$(3.9) \quad \partial_{\mathbf{x}} F \cdot [\mathbf{X}, \mathbf{Y}] = 0.$$

To show this we differentiate (3.7) with respect to  $\mathbf{x}$ , which gives

$$\mathbf{X}^T \cdot \partial_{\mathbf{x}\mathbf{x}} F + \partial_{\mathbf{x}} F \cdot \partial_{\mathbf{x}} \mathbf{X} = 0.$$

The same procedure is applied to (3.8). Using these identities and the definition (3.6) in (3.9) yield the desired result. The corresponding subset  $\mathfrak{G}$  is given by the  $F$ -preserving diffeomorphisms  $\Psi$ , i.e.,

$$F \circ \Psi = F.$$

In fact, let  $\Psi_{\tau}$  be a smooth family of  $F$ -preserving diffeomorphisms with  $\Psi_{\tau=0} = \mathbf{id}$ ; then  $\mathbf{X} := \partial_{\tau} [\Psi_{\tau}]_{\tau=0} \in \mathfrak{g}$  since

$$\partial_{\tau} [F \circ \Psi_{\tau}]_{\tau=0} = \partial_{\mathbf{x}} F \cdot \mathbf{X} = 0.$$

Thus,  $T_{\mathbf{id}} \mathfrak{G} = \mathfrak{g}$  and we can apply Theorem 1. In particular, if a numerical method  $\Psi_{\delta t}$  satisfies

$$F \circ \Psi_{\delta t} = F,$$

then the modified vector fields  $\mathbf{X}_i(\delta t)$ ,  $i = 1, \dots, \infty$ , possess  $F$  as a first integral. The same result was recently derived by Gonzalez, Higham, and Stuart [13] using a contradiction argument.

*Example 2.* Consider the Lie subalgebra of all divergence-free vector fields  $\mathbf{Z}$ , i.e.,  $\text{div } \mathbf{Z} = 0$ . The corresponding subsets  $\mathfrak{G}$  are the volume preserving diffeomorphisms, i.e.,

$$\det \left[ \frac{\partial}{\partial \mathbf{x}} \Psi(\mathbf{x}) \right] = 1.$$

Again we have  $T_{\mathbf{id}} \mathfrak{G} = \mathfrak{g}$ . Namely,

$$\begin{aligned} 0 &= \partial_{\tau} \det \left[ \frac{\partial}{\partial \mathbf{x}} \Psi_{\tau}(\mathbf{x}) \right]_{\tau=0} \\ &= \text{trace} [\partial_{\mathbf{x}} \partial_{\tau} \Psi_{\tau}(\mathbf{x})]_{\tau=0} \\ &= \text{div } \mathbf{X}(\mathbf{x}), \end{aligned}$$

$\mathbf{X} := \partial_\tau [\Psi_\tau]_{\tau=0}$ . Thus, if the numerical method  $\Psi_{\delta t}$  is volume conserving, then the modified vector fields  $\mathbf{X}_i(\delta t)$ ,  $i = 1, \dots, \infty$ , are divergence-free. Again, the same result has been formulated by Gonzalez, Higham, and Stuart [13] using a contradiction argument.

*Example 3.* Let an involution<sup>5</sup>  $\mathbf{S} \in \mathbb{R}^{n \times n}$  be given and consider the subspace  $\mathfrak{g}$  of vector fields  $\mathbf{Z}$  on  $\mathbb{R}^n$  that satisfy the time-reversal symmetry

$$-\mathbf{Z}(\mathbf{x}) = \mathbf{S}\mathbf{Z}(\mathbf{S}\mathbf{x}).$$

This subspace is *not* a subalgebra under the Lie bracket (3.6). The corresponding subset  $\mathfrak{G}$  is given by the time-reversible diffeomorphisms  $\Psi$ , i.e.,  $\Psi^{-1}(\mathbf{x}) = \mathbf{S}\Psi(\mathbf{S}\mathbf{x})$ . Let  $\Psi_\tau \in \mathfrak{G}$  be smooth in  $\tau$  with  $\Psi_{\tau=0} = \mathbf{id}$ ; then

$$\begin{aligned} \mathbf{0} &= \partial_\tau [\mathbf{S}\Psi_\tau \circ \mathbf{S} - [\Psi_\tau]^{-1}]_{\tau=0} \\ &= \mathbf{S}\mathbf{X} \circ \mathbf{S} + \mathbf{X}, \end{aligned}$$

which implies that  $\mathbf{X} := \partial_\tau [\Psi_\tau]_{\tau=0} \in \mathfrak{g}$ . It follows that  $T_{\mathbf{id}} \mathfrak{G} = \mathfrak{g}$  and we can apply Theorem 1. Thus, if a numerical method  $\Psi_{\delta t}$  satisfies the time-reversal symmetry, then the modified vector fields  $\mathbf{X}_i(\delta t)$ ,  $i = 1, \dots, \infty$ , are time reversible. This result has been first proven by Hairer and Stoffer in [20] using ideas from [26].

*Example 4.* Let  $\{.,.\}$  denote the Poisson bracket of a (linear) Poisson manifold  $\mathcal{P} = \mathbb{R}^n$ . Then the Lie algebra of Hamiltonian vector fields on  $\mathcal{P}$  is given by

$$\frac{d}{dt} \mathbf{x} = \{ \mathbf{id}, H \}(\mathbf{x}),$$

where  $H : \mathcal{P} \rightarrow \mathbb{R}$  is a smooth function. The corresponding subset  $\mathfrak{G}$  is given by the set of smooth diffeomorphisms on  $\mathcal{P}$  that preserve the Poisson bracket  $\{.,.\}$  [1]. Let  $\Psi_\tau$  be a family of maps in  $\mathfrak{G}$  with  $\Psi_{\tau=0} = \mathbf{id}$ . Then

$$\begin{aligned} 0 &= \partial_\tau [\{F \circ \Psi_\tau, G \circ \Psi_\tau\} - \{F, G\}]_{\tau=0} \\ &= \{F, \partial_{\mathbf{x}}G \cdot \mathbf{X}\} + \{\partial_{\mathbf{x}}F \cdot \mathbf{X}, G\} \end{aligned}$$

for all smooth functions  $F, G : \mathcal{P} \rightarrow \mathbb{R}$ ,  $\mathbf{X} := \partial_\tau [\Psi_\tau]_{\tau=0}$ . This is the condition for a vector field  $\mathbf{X}$  to be locally Hamiltonian. Since  $\mathcal{P}$  is simply connected, the vector field is also globally Hamiltonian [3].

If the discrete evolution (2.3) satisfies  $\Psi_{\delta t} \in \mathfrak{G}$  for all  $\delta t > 0$ , then  $\Psi_{\delta t}$  is called a *symplectic method* and it follows from Theorem 1 that the modified vector fields  $\mathbf{X}_i(\delta t)$ ,  $i = 1, \dots, \infty$ , are Hamiltonian vector fields on  $\mathcal{P}$ . This result can also be found in [5, 16, 26].

If a symplectic method can be expanded as a P-series, then the vector fields  $\mathbf{X}_i(\delta t)$  are globally Hamiltonian even if the phase space  $\mathcal{P} \subset \mathbb{R}^n$  is not simply connected [16]. This result applies to all symplectic Runge–Kutta and partitioned Runge–Kutta methods. Furthermore, symplectic methods defined by a generating function of the third kind [32] are also always globally Hamiltonian [5]. The same statement is true for symplectic methods based on the composition of exact flow maps [40].

*Example 5.* Let us now consider differential equations on a matrix Lie group  $\mathcal{G} \subset \mathbb{R}^{n \times n}$  [35]. In general, time-independent differential equations on  $\mathcal{G}$  can be written in the form

$$\frac{d}{dt} \mathbf{Y} = \mathbf{A}(\mathbf{Y}) \mathbf{Y},$$

---

<sup>5</sup>An involution is a nonsingular matrix that satisfies  $\mathbf{S}^{-1} = \mathbf{S}$ .

where  $\mathbf{A} : \mathcal{G} \rightarrow \mathfrak{g}$ ,  $\mathfrak{g} \subset \mathbb{R}^{n \times n}$  the Lie algebra of  $\mathcal{G}$ . Many recent papers (see [8] and references therein) have been devoted to methods that preserve the Lie group structure; i.e.,

$$\mathbf{Y}_{n+1} = \Psi_{\delta t}(\mathbf{Y}_n),$$

and  $\mathbf{Y}_n \in \mathcal{G}$  implies  $\mathbf{Y}_{n+1} \in \mathcal{G}$ . Thus,  $\Psi_{\delta t}$  is a diffeomorphism defined on the submanifold  $\mathcal{G} \subset \mathbb{R}^{n \times n}$ . In fact, this submanifold can be characterized, at least locally, by a set of nonlinear equations which we denote by  $\mathbf{F}(\mathbf{Y}) = \mathbf{0}$ . Thus  $\Psi_{\delta t}$  is an  $\mathbf{F}$  integral preserving map; i.e.,  $\mathbf{F} \circ \Psi_{\delta t} = \mathbf{F}$  on  $\mathcal{G}$ . Following Example 1, we know then that the modified vector fields (as well as the given vector field) satisfy  $\partial_{\mathbf{Y}} \mathbf{F} \cdot \mathbf{X}_i(\delta t) = \mathbf{0}$ . Hence the modified vector fields  $\mathbf{X}_i(\delta t)$  are vector fields on  $\mathcal{G}$  and give rise to modified differential equations of type

$$\frac{d}{dt} \mathbf{Y} = \tilde{\mathbf{A}}_i(\mathbf{Y}; \delta t) \mathbf{Y},$$

with

$$\mathbf{X}_i(\mathbf{Y}; \delta t) = \delta t \tilde{\mathbf{A}}_i(\mathbf{Y}; \delta t) \mathbf{Y}$$

and  $\tilde{\mathbf{A}}_i(\delta t) : \mathcal{G} \rightarrow \mathfrak{g}$ . See [28] for further results on backward error analysis for numerical methods on manifolds.

**4. Truncation error of backward error analysis.** We would like to derive an explicit estimate for the norm of the vector fields  $\Delta \mathbf{X}_{i+1}$  and the difference between the time-one-flow map  $\Phi_{1, \mathbf{X}_i(\delta t)}$  and the numerical approximation  $\Psi_{\delta t}$ ,  $i = 1, \dots, \infty$ . To do so we assume from now on that the vector field  $\mathbf{Z}$  in (2.2) is real analytic. We also introduce the following notation: Let  $\mathcal{B}_r(\mathbf{x}_0) \subset \mathbb{C}^n$  denote the complex ball of radius  $r > 0$  around  $\mathbf{x}_0 \in \mathbb{R}^n$  and define

$$\|\mathbf{z}\| := \max_{i=1, \dots, n} |z_i| \quad (\mathbf{z} \in \mathbb{C}^n).$$

Let us consider a compact subset  $\mathcal{K} \subset \mathbb{R}^n$  of phase space and a constant  $r > 0$  such that a given real analytic vector field  $\mathbf{Y}$  is bounded on  $\mathcal{B}_r(\mathbf{x}_0)$  for all  $\mathbf{x}_0 \in \mathcal{K}$ . Then we define

$$\|\mathbf{Y}\|_r = \sup_{\mathbf{x} \in \mathcal{B}_r \mathcal{K}} \|\mathbf{Y}(\mathbf{x})\|$$

with

$$\mathcal{B}_r \mathcal{K} := \bigcup_{\mathbf{x}_0 \in \mathcal{K}} \mathcal{B}_r(\mathbf{x}_0).$$

We also define  $\mathcal{B}_0 \mathcal{K} = \mathcal{K}$  and

$$\|\mathbf{Y}\|_0 = \sup_{\mathbf{x} \in \mathcal{K}} \|\mathbf{Y}(\mathbf{x})\|.$$

To find an estimate for  $\Delta \mathbf{X}_{i+1}$ , as defined in (2.5), we need estimates for the mappings appearing on the right-hand side of (2.5). We start with an estimate for the map  $\Psi_{\delta t}$ .

LEMMA 1. *Let us assume that the vector field  $\mathbf{Z}$  in (2.2) is real analytic and that there is a compact subset  $\mathcal{K}$  of phase space and constants  $K, R > 0$  such that*

$$\|\mathbf{Z}\|_R \leq K.$$



We also assume that the numerical method  $\Psi_{\delta t}$  is real analytic. Then there exists a constant  $M \geq K$  such that

$$(4.10) \quad \|\Psi_\tau - \mathbf{id}\|_{\alpha R} \leq |\tau| M \leq (1 - \alpha) R \quad \text{for } |\tau| \leq \frac{(1 - \alpha) R}{M},$$

$\alpha \in [0, 1)$ .

*Proof.* Under the given assumptions, the flow map

$$\Phi_{\tau, \mathbf{Z}}(\mathbf{x}) = \mathbf{x} + \int_0^\tau \mathbf{Z}(\Phi_{t, \mathbf{Z}}(\mathbf{x})) dt$$

is defined for complex-valued  $\tau \in \mathbb{C}$ , where the integral on the right-hand side is independent of the path from zero to  $\tau$ . The complexified flow map satisfies

$$(4.11) \quad \begin{aligned} \|\Phi_{\tau, \mathbf{Z}} - \mathbf{id}\|_{\alpha R} &\leq \sup_{\mathbf{x} \in B_{\alpha R} \mathcal{K}} \int_0^\tau \|\mathbf{Z}(\Phi_{t, \mathbf{Z}}(\mathbf{x}))\| |dt| \\ &\leq |\tau| K \leq (1 - \alpha) R \quad \text{for } |\tau| \leq \frac{(1 - \alpha) R}{K}, \end{aligned}$$

$\alpha \in [0, 1)$ . Consistency of the numerical method implies that there exists a constant  $\Delta K > 0$  such that

$$\|\Psi_\tau - \mathbf{id}\|_{\alpha R} \leq |\tau| (K + \Delta K) \leq (1 - \alpha) R \quad |\tau| \leq \frac{(1 - \alpha) R}{K + \Delta K}$$

for the (complexified) map  $\Psi_\tau$ . Take  $M := K + \Delta K$ . □

*Remark.* Now let us consider an  $s$ -stage Runge–Kutta method with coefficients  $\{a_{ij}\}_{i,j=1,\dots,s}$  and  $\{b_i\}_{i=1,\dots,s}$  [19] satisfying

$$\sum_{j=1}^s |a_{ij}| \leq d \quad \text{and} \quad \sum_{i=1}^s |b_i| \leq d,$$

$d \geq 1$ , and assume that the Runge–Kutta method uniquely<sup>6</sup> defines a real analytic map  $\Psi_\tau$  for all step sizes  $\tau \in \mathbb{C}$  with  $|\tau| \leq R/K$ . Then we have  $M = dK$  in Lemma 1. This follows from the fact that, under the stated assumptions, all stage variables will be in  $B_R \mathcal{K}$ , where the vector field  $\mathbf{Z}$  is bounded by the constant  $K$ . A similar statement holds for partitioned Runge–Kutta methods.

**THEOREM 2.** *Let the assumptions of Lemma 1 be satisfied. Then there exists a family of real analytic vector fields  $\tilde{\mathbf{X}}(\delta t) : \mathcal{V} \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$ ,  $\mathcal{K} \subset \mathcal{V} \subset \mathcal{U}$ , such that*

$$\|\Psi_{\delta t} - \Phi_{\delta t, \tilde{\mathbf{X}}(\delta t)}\|_0 \leq 8 \delta t b M e^{-p} e^{-\gamma/\delta t}$$

with  $\gamma = R/(cMe)$ ,  $b = 20$ ,  $c = 300$ , and  $p \geq 1$  the order of the method. The family of modified vector fields  $\tilde{\mathbf{X}}(\delta t)$  satisfies the estimate

$$\|\tilde{\mathbf{X}}(\delta t) - \mathbf{Z}\|_0 \leq 2 d_p b M \left( \frac{c \delta t M}{R} \right)^p$$

with  $d_p \geq 1$  a constant depending on the order  $p$  of the method. For example,  $d_1 = 0.8$ ,  $d_2 = 1.5$ ,  $d_3 = 3.7$ , and  $d_4 = 12.8$ .

---

<sup>6</sup>For an implicit method, the solution can be obtained by fixed point iteration if  $|\tau|$  is sufficiently small.

*Proof.* We know that  $\mathbf{X}_1(\delta t) = \delta t \mathbf{Z}$  and that  $\Psi_{\delta t}$  is a method of order  $p \geq 1$ . Thus  $\|\mathbf{X}_1(\tau)\|_R \leq |\tau|M$  and  $\Delta \mathbf{X}_i = \mathbf{0}$  for  $i = 2, \dots, p$ . Next we find an estimate for the difference between the time-one-flow map  $\Phi_{1, \mathbf{X}_1(\tau)}$  and the map  $\Psi_\tau$ . Using (4.10) and (4.11), we obtain

$$\begin{aligned} \|\Phi_{1, \mathbf{X}_1(\tau)} - \Psi_\tau\|_{\alpha R} &= \|\Phi_{1, \mathbf{X}_1(\tau)} - \mathbf{id} + \mathbf{id} - \Psi_\tau\|_{\alpha R} \\ &\leq \|\Phi_{\tau, \mathbf{Z}} - \mathbf{id}\|_{\alpha R} + \|\Psi_\tau - \mathbf{id}\|_{\alpha R} \\ &\leq 2(1 - \alpha)R \end{aligned}$$

for  $|\tau| \leq \tau_1 := (1 - \alpha)R/M$ . Since the mappings are real analytic and their difference is  $\mathcal{O}(\delta t^{p+1})$ , we obtain the estimate [5]

$$\begin{aligned} \|\Phi_{1, \mathbf{X}_1(\delta t)} - \Psi_{\delta t}\|_{\alpha R} &\leq 2(1 - \alpha)R \left(\frac{\delta t}{\tau_1}\right)^{p+1} \\ &\leq 2\delta t M \left(\frac{\delta t M}{(1 - \alpha)R}\right)^p, \end{aligned}$$

$\alpha \in [0, 1)$ . Using this in (2.5) with  $i = p$ , we obtain

$$(4.12) \quad \|\Delta \mathbf{X}_{p+1}\|_{\alpha R} \leq 2M \left(\frac{M}{(1 - \alpha)R}\right)^p.$$

Next we show that

$$(4.13) \quad \|\Delta \mathbf{X}_i\|_{\alpha R} \leq bM \left(\frac{c(i - p)M}{(1 - \alpha)R}\right)^{i-1}$$

for  $i \geq p + 1$  with  $b = 20$ ,  $c = 300$ , and  $\alpha \in [0, 1)$ . The estimate is true for  $i = p + 1$  (compare (4.12)). We proceed by induction. First note that, for  $j \geq p + 2$ ,

$$\begin{aligned} \|\mathbf{X}_j(\tau)\|_{\alpha R} &\leq |\tau| \|\mathbf{Z}\|_{\alpha R} + \sum_{i=p+1}^j |\tau|^i \|\Delta \mathbf{X}_i\|_{\alpha R} \\ (4.14) \quad &\leq |\tau|M \left[ 1 + 2 \left(\frac{|\tau|M}{(1 - \alpha)R}\right)^p + \sum_{i=p+2}^j b \left(\frac{c(i - p)|\tau|M}{(1 - \alpha)R}\right)^{i-1} \right], \end{aligned}$$

$\alpha \in [0, 1)$ . We replace the parameter  $\alpha \in [0, 1)$  in this formula by  $\alpha + \delta_j(1 - \alpha) \in [\delta_j, 1)$ , where

$$\delta_j := \frac{b - 1}{(j - p + 1)c} = \frac{0.067}{j - p + 1}.$$

This yields

$$\|\mathbf{X}_j(\tau)\|_{(\alpha + \delta_j(1 - \alpha))R} \leq (b - 1)\tau_j M = \delta_j(1 - \alpha)R$$

for all  $\alpha \in [0, 1)$  and all  $\tau \in \mathbb{C}$  with

$$|\tau| \leq \frac{(1 - \alpha)R}{(j - p + 1)cM} =: \tau_j.$$

Here we have used that

$$(4.15) \quad \sum_{i=p+2}^j \left( \frac{i-p}{(1-\delta_j)(j-p+1)} \right)^{i-1} \leq 0.891$$

for  $p \geq 1$  and all  $j \geq p+2$ . In particular, substitute  $\tau_j$  for  $|\tau|$  and  $\alpha + \delta_j(1-\alpha)$  for  $\alpha$  in (4.14). Then use the identity  $1-\alpha-\delta_j(1-\alpha) = (1-\delta_j)(1-\alpha)$  and the inequality (4.15) to derive

$$\tau_j M \left[ 1 + 2 \left( \frac{1}{c\beta_j} \right)^p + b \sum_{i=p+2}^j \left( \frac{i-p}{\beta_j} \right)^{i-1} \right] \leq (b-1) \tau_j M,$$

$\beta_j := (1-\delta_j)(j-p+1)$ . Next we introduce the vector-valued, real analytic function

$$\begin{aligned} \mathbf{f}_\tau(\mathbf{x}) &:= \Phi_{1, \mathbf{X}_j(\tau)}(\mathbf{x}) - \mathbf{x} \\ &= \int_0^1 \mathbf{X}_j(\Phi_{t, \mathbf{X}_j(\tau)}(\mathbf{x}); \delta t) dt \end{aligned}$$

and observe that

$$(4.16) \quad \|\mathbf{f}_\tau\|_{\alpha R} \leq (b-1) \tau_j M = \delta(1-\alpha) R$$

for  $|\tau| \leq \tau_j$  and  $\alpha \in [0, 1)$ . Here we have used that  $\mathbf{x} \in \mathcal{B}_{\alpha R} \mathcal{K}$  implies  $\Phi_{t, \mathbf{X}_j(\tau)}(\mathbf{x}) \in \mathcal{B}_{(\alpha+\delta_j(1-\alpha))R} \mathcal{K}$  for all  $0 \leq |t| \leq 1$  and any  $|\tau| \leq \tau_j$ . Now we can find an estimate for the difference between the time-one-flow map  $\Phi_{1, \mathbf{X}_j(\tau)}$  and the map  $\Psi_\tau$ . Using (4.10) and (4.16), we obtain

$$\begin{aligned} \|\Phi_{1, \mathbf{X}_j(\tau)} - \Psi_\tau\|_{\alpha R} &\leq \|\mathbf{f}_\tau\|_{\alpha R} + \|\Psi_\tau - \mathbf{id}\|_{\alpha R} \\ &\leq (b-1) \tau_j M + \tau_j M \\ &\leq b \tau_j M \end{aligned}$$

for  $|\tau| \leq \tau_j$ . Since the mappings are real analytic and their difference is  $\mathcal{O}(\delta t^{j+1})$ , we obtain the estimate [5]

$$(4.17) \quad \begin{aligned} \|\Phi_{1, \mathbf{X}_j(\delta t)} - \Psi_{\delta t}\|_{\alpha R} &\leq b \tau_j M \left( \frac{\delta t}{\tau_j} \right)^{j+1} \\ &\leq b \delta t M \left( \frac{c(j-p+1) \delta t M}{(1-\alpha) R} \right)^j, \end{aligned}$$

$\alpha \in [0, 1)$ . Using this in (2.5) with  $i = j$ , we finally obtain

$$\|\Delta \mathbf{X}_{j+1}\|_{\alpha R} \leq b M \left( \frac{c(j-p+1) M}{(1-\alpha) R} \right)^j,$$

which verifies (4.13) for  $i = j + 1$ .

Next we need an estimate for the difference between the time-one-flow map  $\Phi_{1, \mathbf{X}_i(\delta t)}$  and the map  $\Psi_{\delta t}$  on the compact set  $\mathcal{K}$ . Using (4.17) with  $\alpha = 0$  and  $i = j$ , we immediately have

$$\|\Phi_{1, \mathbf{X}_i(\delta t)} - \Psi_{\delta t}\|_0 \leq \delta t b M \left( \frac{c \delta t (i-p+1) M}{R} \right)^i.$$

The family of vector fields  $\tilde{\mathbf{X}}(\delta t)$  is now defined by taking an optimal number  $i_*(\delta t)$  of iterations. We take  $i_*(\delta t)$  as the integer part of

$$i_o(\delta t) := \frac{R}{c \delta t M e} + p - 1.$$

Thus

$$\begin{aligned} \|\Phi_{\tau, \mathbf{X}_{i_*}(\delta t)} - \Psi_{\delta t}\|_0 &\leq \delta t b M e^{-i_*} \\ &\leq \delta t b M e^{-i_o+1} \\ &\leq 8 \delta t b M e^{-p} e^{-\gamma/\delta t}, \end{aligned}$$

$\gamma = R/(cMe)$ . We define  $\tilde{\mathbf{X}}(\delta t) := \delta t^{-1} \mathbf{X}_{i_*}(\delta t)$ . This completes the first part of the proof.

According to (4.14), the difference between the modified vector fields  $\tilde{\mathbf{X}}(\delta t)$  and  $\mathbf{Z}$  is given by

$$\|\tilde{\mathbf{X}}(\delta t) - \mathbf{Z}\|_0 \leq M \left( \frac{c \delta t M}{R} \right)^p \left[ \frac{2}{c^p} + \sum_{i=p+2}^{i_*} b (i-p)^p \left( \frac{c(i-p) \delta t M}{R} \right)^{i-p-1} \right].$$

Next we use

$$\delta t \leq \frac{R}{c(i_* - p + 1) M e}$$

to obtain

$$\begin{aligned} \|\tilde{\mathbf{X}}(\delta t) - \mathbf{Z}\|_0 &\leq M \left( \frac{c \delta t M}{R} \right)^p \left[ 0.0067 + b \sum_{i=p+2}^{i_*} \frac{(i-p)^p}{e^{i-p-1}} \left( \frac{i-p}{i_* - p + 1} \right)^{i-p-1} \right] \\ &\leq M \left( \frac{c \delta t M}{R} \right)^p [0.0067 + b d_p 1.38] \\ &\leq 2 d_p b M \left( \frac{c \delta t M}{R} \right)^p. \end{aligned}$$

Here  $d_p \geq 1$  is chosen such that

$$d_p \geq \frac{j^p}{e^{j-1}}$$

for all  $j \geq 2$ .  $\square$

*Remark.* The proof of Theorem 2 is similar in spirit to the one given by Benettin and Giorgilli [5] on the exponentially small difference between an optimal interpolating vector field and a near-to-the-identity map. However, there are a couple of important differences: (i) We explicitly take the order of a method into account. (ii) We directly derive estimates on the difference between the flow maps  $\Phi_{1, \mathbf{X}_i(\tau)}$  and  $\Psi_{\delta t}$  instead of using Taylor series expansions of  $\Phi_{1, \mathbf{X}_i(\tau)}$  and  $\Psi_{\delta t}$  and corresponding estimates for the elements in the series. We believe that this simplifies the proof of Theorem 2. (iii) By introducing the parameter  $\alpha \in [0, 1)$ , we do not have to shrink the domain of definition of the vector fields  $\mathbf{X}_i(\tau)$  as the iteration index  $i$  increases. Again we feel that this simplifies the proof. (iv) As in [18], we work directly with an estimate for the given vector field  $\mathbf{Z}$  instead of making assumptions on the map  $\Psi_{\delta t}$ . The rather pessimistic constants  $c = 300$  and  $b = 20$  entering the estimates seem to be the main disadvantage of our approach.

A more elaborate version of the proof of Theorem 2 can be found in [30].

**5. An application: Ergodic Hamiltonian systems.** Let us consider a (real analytic) Hamiltonian system

$$(5.18) \quad \frac{d}{dt} \mathbf{q} = \mathbf{M}^{-1} \mathbf{p},$$

$$(5.19) \quad \frac{d}{dt} \mathbf{p} = -\nabla_{\mathbf{q}} V(\mathbf{q}),$$

$\mathbf{q}, \mathbf{p} \in \mathbb{R}^n$ , together with a smooth function  $A : \mathbb{R}^{2n} \rightarrow \mathbb{R}$ . We are interested in evaluating the time average of  $A$  along a trajectory  $(\mathbf{q}(t), \mathbf{p}(t))$  of the Hamiltonian system (5.18)–(5.19); i.e.,

$$\langle A \rangle_T := \frac{1}{T} \int_0^T A(\mathbf{q}(t), \mathbf{p}(t)) dt, \quad T \gg 1.$$

We assume that

$$\langle A \rangle_\infty := \lim_{T \rightarrow \infty} \langle A \rangle_T$$

exists and is equal to the microcanonical ensemble average corresponding to the Hamiltonian

$$H(\mathbf{q}, \mathbf{p}) = \frac{\mathbf{p}^T \mathbf{M}^{-1} \mathbf{p}}{2} + V(\mathbf{q});$$

i.e., we assume that the system (5.18)–(5.19) is ergodic<sup>7</sup> (or even mixing) [37]. Thus

$$\langle A \rangle_\infty = \frac{\int A(\mathbf{q}, \mathbf{p}) \delta(E - H(\mathbf{q}, \mathbf{p})) d\mathbf{q}d\mathbf{p}}{\int \delta(E - H(\mathbf{q}, \mathbf{p})) d\mathbf{q}d\mathbf{p}} =: \frac{1}{C} \langle A, \delta(E - H) \rangle$$

with  $E = H(\mathbf{q}(0), \mathbf{p}(0))$ ,  $\delta(x)$  Dirac’s delta distribution,

$$C := \int \delta(E - H(\mathbf{q}, \mathbf{p})) d\mathbf{q}d\mathbf{p},$$

and

$$\langle A, \delta(E - H) \rangle := \int A(\mathbf{q}, \mathbf{p}) \delta(E - H(\mathbf{q}, \mathbf{p})) d\mathbf{q}d\mathbf{p}$$

the inner product of  $A$  and  $\delta(E - H)$ .

Let us write the equations (5.18)–(5.19) in more compact form as

$$\frac{d}{dt} \mathbf{x} = \mathbf{J} \nabla_{\mathbf{x}} H(\mathbf{x}) = \{\mathbf{id}, H\}(\mathbf{x}),$$

$\mathbf{x} := (\mathbf{q}^T, \mathbf{p}^T)^T \in \mathbb{R}^{2n}$ . The Hamiltonian  $H$  is preserved under the flow map  $\Phi_{t,H}$ . Let us assume that the hypersurface  $\mathcal{M}_0$  of constant energy  $H = 0$ ,

$$\mathcal{M}_0 := \{\mathbf{x} \in \mathbb{R}^{2n} : H(\mathbf{x}) = 0\},$$

---

<sup>7</sup>To be more precise, ergodicity of a system implies that the time average is equivalent to the ensemble average except for, at most, a set of initial conditions of measure zero.

is a compact subset of  $\mathbb{R}^{2n}$ . We also assume that there is a constant  $\gamma_1 > 0$  such that  $\|\nabla_{\mathbf{x}}H(\mathbf{x})\| > \gamma_1$  for all  $\mathbf{x} \in \mathcal{M}_0$ . This implies that  $\mathcal{M}_0$  is a smooth  $(2n - 1)$ -dimensional compact submanifold. Furthermore, the family of hypersurfaces

$$\mathcal{M}_E = \{\mathbf{x} \in \mathbb{R}^{2n} : H(\mathbf{x}) = E\}, \quad E \in (-\Delta E, +\Delta E),$$

$\Delta E > 0$  sufficiently small, are smooth and compact as well (in fact diffeomorphic to  $\mathcal{M}_0$ ). We define the open subset  $\mathcal{U}$  of phase space by

$$\mathcal{U} := \bigcup_{E \in (-\Delta E, +\Delta E)} \mathcal{M}_E.$$

So far we have made fairly generic assumptions. In the sequel, we become more specific to ensure that the Hamiltonian system (5.18)–(5.19) is ergodic/mixing.

In a first step we construct a Poincaré return map [14]. Let  $\psi : \mathcal{U} \rightarrow \mathbb{R}$  be a smooth function and  $\gamma_2 > 0$  a positive constant such that  $|\{\psi, H\}(\mathbf{x})| > \gamma_2$  on the level sets

$$\mathcal{S}_s := \{\mathbf{x} \in \mathcal{U} : \psi(\mathbf{x}) = s\}, \quad s \in (-\Delta s, +\Delta s),$$

$\Delta s > 0$  sufficiently small. Let us assume that  $\mathcal{S}_s$  defines a Poincaré section for each  $s \in (-\Delta s, +\Delta s)$  in the following way: For all  $\mathbf{x} \in \mathcal{S}_s$ , there is a positive number  $t_p(\mathbf{x}) > 0$  such that the solution  $\mathbf{x}(t)$ ,  $t \geq 0$ , with initial condition  $\mathbf{x}(0) = \mathbf{x}$  satisfies  $\mathbf{x}(t_p) \in \mathcal{S}_s$  and there is no  $0 < t'_p < t_p$  such that  $\mathbf{x}(t'_p) \in \mathcal{S}_s$ . The positive number  $t_p(\mathbf{x})$  is called the Poincaré return time of the point  $\mathbf{x} \in \mathcal{S}_s$ . Knowing the Poincaré return time for each  $\mathbf{x} \in \mathcal{S}_s$ , we define the “global” Poincaré map  $\mathbf{\Pi} : \mathcal{V} \rightarrow \mathcal{V}$  by

$$\mathbf{\Pi}(\mathbf{x}) := \Phi_{t_p(\mathbf{x}), H}(\mathbf{x})$$

and

$$\mathcal{V} := \bigcup_{s \in (-\Delta s, +\Delta s)} \mathcal{S}_s.$$

We assume that the Poincaré return times  $t_p(\mathbf{x})$ ,  $\mathbf{x} \in \mathcal{V}$ , are bounded by some constant  $K > 0$ .

We are interested in the solutions on a particular level set of constant energy. For simplicity, we take the level set  $\mathcal{M}_0$ . Then it is sufficient to consider the “restricted” Poincaré map  $\mathbf{\Pi}_0$ , which is defined as the restriction of  $\mathbf{\Pi}$  to

$$\mathcal{D} := \mathcal{S}_0 \cap \mathcal{M}_0.$$

Thus we have reduced the study of the dynamical properties of the Hamiltonian system (5.18)–(5.19) on the energy shell  $\mathcal{M}_0$  to the study of the properties of the Poincaré map  $\mathbf{\Pi}_0$ . If  $\mathbf{\Pi}_0$  is an ergodic (mixing) map, then the Hamiltonian system is ergodic (mixing) on  $\mathcal{M}_0$ . Note that  $\mathbf{\Pi}_0$  is volume preserving; i.e.,  $\det \partial_{\mathbf{x}}\mathbf{\Pi}_0(\mathbf{x}) = 1$ .

From now on we assume that  $\mathbf{\Pi}_0$  is a uniformly hyperbolic map; i.e., for each  $\mathbf{x} \in \mathcal{D}$ , the linearization  $\partial_{\mathbf{x}}\mathbf{\Pi}_0(\mathbf{x})$  at  $\mathbf{x}$  possesses strictly expanding and contracting directions only [14, 36]. The “stochastic” behavior of such a (deterministic) map has been investigated, for example, in [36]. Here we only point out the four main results.

- There is a unique invariant density  $\mu_0$  on  $\mathcal{D}$  that is invariant under  $\mathbf{\Pi}_0$ . Furthermore,  $\mu_0$  is given by the Lebesgue measure on  $\mathcal{D}$ .

- The autocorrelation function  $\langle A \circ [\mathbf{\Pi}_0]^n, A \rangle$  of a Hölder continuous function  $A : \mathcal{U} \rightarrow \mathbb{R}$  decays exponentially fast, i.e.,

$$|\langle A \circ [\mathbf{\Pi}_0]^n, A \rangle - (\langle A, \mu_0 \rangle)^2| \leq C \Lambda^n, \quad 0 < \Lambda < 1,$$

$C > 0$  an appropriate constant.

- The time averages

$$\langle A \rangle_N = \frac{1}{N} \sum_{i=1}^N A(\mathbf{x}_i)$$

of  $A$  along trajectories  $\{\mathbf{x}_i\}_{i=1, \dots, N}$  of  $\mathbf{\Pi}_0$  satisfy a central limit theorem.

- The time average  $\langle A \rangle_N$  of  $A$  along trajectories of  $\mathbf{\Pi}_0$  with initial value  $\mathbf{x}_0 \in \mathcal{D}$  satisfy a large deviation theorem. To be more specific [39], given any  $c > 0$  there is an  $h(c) > 0$  such that

$$(5.20) \quad \mu_0(\{\mathbf{x}_0 \in \mathcal{D} : |\langle A \rangle_N - \langle A, \mu_0 \rangle| > c\}) \leq e^{-Nh(c)}$$

for all large  $N \geq 1$ .

These results can be proven (see, for example, [36]) by carefully studying the properties of the corresponding Frobenius–Perron operator  $\mathbf{P}_0 : L^1(\mathcal{D}) \rightarrow L^1(\mathcal{D})$  defined by

$$\mathbf{P}_0 \mu := \mu \circ [\mathbf{\Pi}_0]^{-1},$$

$\mu \in L^1(\mathcal{D})$ .

DEFINITION 1. We call a Hamiltonian system (5.18)–(5.19) with the above introduced properties Poincaré hyperbolic. In particular, we assume (i) that the level sets  $\mathcal{M}_E$ ,  $E \in (-\Delta E, +\Delta E)$  of constant energy are compact submanifolds; (ii) that there is a constant  $\gamma_1 > 0$  such that  $\|\nabla_{\mathbf{x}} H(\mathbf{x})\| > \gamma_1$  for all  $\mathbf{x} \in \mathcal{U}$ ; (iii) that a global Poincaré map  $\mathbf{\Pi}$  can be defined on

$$\mathcal{V} = \bigcup_{s \in (-\Delta s, +\Delta s)} \mathcal{S}_s,$$

which is uniformly hyperbolic as a map restricted to  $\mathcal{D} = \mathcal{S}_0 \cap \mathcal{M}_0$ ; (iv) that the Poincaré return times  $t_p(\mathbf{x})$ ,  $\mathbf{x} \in \mathcal{V}$ , are bounded by some constant  $K > 0$ ; and (v) that there is a constant  $\gamma_2 > 0$  such that  $|\{\psi, H\}(\mathbf{x})| > \gamma_2$  on  $\mathcal{V}$ .

Let  $\tilde{H}$  be a perturbation of  $H$  such that

$$|H(\mathbf{x}) - \tilde{H}(\mathbf{x})| + \|\nabla_{\mathbf{x}} H(\mathbf{x}) - \nabla_{\mathbf{x}} \tilde{H}(\mathbf{x})\| \leq \epsilon$$

for all  $\mathbf{x} \in \mathcal{U}$  and some  $\epsilon > 0$ . Then we call  $\tilde{H}$  an  $\epsilon$ -perturbation of  $H$ .

LEMMA 2. The property of being Poincaré hyperbolic is stable under  $\epsilon$ -perturbations of the Hamiltonian  $H$  provided  $\epsilon$  is sufficiently small.

Proof. The assumption  $\|\nabla_{\mathbf{x}} H(\mathbf{x})\| > \gamma_1$  on the level sets  $\mathcal{M}_E$  implies that these sets are persistent under small perturbations. Furthermore, there exists a constant  $\tilde{\gamma}_2 > 0$  such that  $|\{\psi, \tilde{H}\}(\mathbf{x})| > \tilde{\gamma}_2$  for a perturbed Hamiltonian  $\tilde{H}$  and  $\mathbf{x} \in \mathcal{V}$ . Thus a Poincaré map is also defined for the perturbed Hamiltonian  $\tilde{H}$ . Uniform hyperbolicity is also stable under small perturbations of the Poincaré map [2].  $\square$

Let us discretize (5.18)–(5.19) by a symplectic (real analytic) integrator  $\Psi_{\delta t}$  of order  $p \geq 1$ .

ASSUMPTION. We assume that backward error analysis can be applied on a compact subset  $\mathcal{K}$  with  $\mathcal{U} \subset \mathcal{K}$ . The corresponding perturbed Hamiltonian is denoted by  $\tilde{H}(\delta t)$ ; i.e., for all  $\mathbf{x} \in \mathcal{K}$ ,

$$\|\Phi_{\delta t, \tilde{H}}(\mathbf{x}) - \Psi_{\delta t}(\mathbf{x})\| \leq \delta t d_1 e^{-p} e^{-d_2/\delta t},$$

$d_1, d_2 > 0$  are appropriate constants. Let the step size  $\delta t$  be sufficiently small such that the perturbed Hamiltonian system is also Poincaré hyperbolic. For simplicity, we shift the modified Hamiltonian  $\tilde{H}(\delta t)$  such that  $H(\mathbf{x}_0) = \tilde{H}(\mathbf{x}_0; \delta t) = 0$ .

Let us introduce notation for the perturbed system. As for the unperturbed system, we define the compact level sets  $\tilde{\mathcal{M}}_E$  and the open set  $\tilde{\mathcal{U}}$  (replacing  $H$  by  $\tilde{H}$  in the definition). Without loss of generality, we can assume that  $\tilde{\mathcal{U}} \subset \mathcal{K}$ . Furthermore,

$$\tilde{\mathcal{S}}_s := \{\mathbf{x} \in \tilde{\mathcal{U}} : \psi(\mathbf{x}) = s\},$$

$s \in (-\Delta s, +\Delta s)$ . The corresponding sets  $\tilde{\mathcal{V}}$  and  $\tilde{\mathcal{D}}$  are now defined in the obvious way. Finally, the global Poincaré map  $\tilde{\Pi}$  and the reduced Poincaré map  $\tilde{\Pi}_0$  are introduced as for the unperturbed system. Again, without loss of the generality, we can assume that  $\mathcal{D} \subset \tilde{\mathcal{V}}$ .

We extend the discrete time map  $\Psi_{\delta t}$  to a map  $\Psi_t$ ,  $t \in [0, \delta t]$ , by using the exact flow map  $\Phi_{t, \tilde{H}}$  of the modified problem as an interpolation for  $t \in [0, \delta t]$ . The map is then extended to  $t \geq \delta t$  in the obvious way<sup>8</sup> as the composition of  $k$  steps with  $\Psi_{\delta t}$  and one-step with  $\Phi_{dt, \tilde{H}}$  where  $t = k\delta t + dt$ ,  $\delta t > dt \geq 0$ . Thus, in correspondence with the definition of the global Poincaré map

$$\tilde{\Pi}(\mathbf{x}) := \Phi_{\tilde{t}_p(\mathbf{x}), \tilde{H}}(\mathbf{x}),$$

we define

$$\hat{\Pi}(\mathbf{x}) := \Psi_{\tilde{t}_p(\mathbf{x})}(\mathbf{x})$$

for all  $\mathbf{x} \in \tilde{\mathcal{V}}$ . Here the Poincaré return times  $\tilde{t}_p(\mathbf{x})$ ,  $\mathbf{x} \in \tilde{\mathcal{V}}$ , are the same as in the definition of  $\tilde{\Pi}$ . Lemma 2 implies that there is a constant  $\tilde{K} > 0$  such that

$$\sup_{\mathbf{x} \in \tilde{\mathcal{V}}} \tilde{t}_p(\mathbf{x}) \leq \tilde{K}.$$

It follows from backward and forward error analysis [18] that there is a constant  $d_3 > 0$  such that

$$\|\tilde{\Pi}(\mathbf{x}) - \hat{\Pi}(\mathbf{x})\| \leq d_3 e^{-p} e^{-d_2/\delta t}$$

for all  $\mathbf{x} \in \tilde{\mathcal{V}}$  and for all  $\delta t$  sufficiently small. More importantly, let  $\{\mathbf{x}_i\}_{i=1, \dots, N}$  be a “numerically” computed sequence of points with  $\mathbf{x}_{i+1} = \hat{\Pi}(\mathbf{x}_i)$  and let  $\{\tilde{\mathbf{x}}_i\}_{i=1, \dots, N}$  be the corresponding sequence under the map  $\tilde{\Pi}$  with  $\mathbf{x}_0 = \tilde{\mathbf{x}}_0 \in \tilde{\mathcal{D}}$ . Each sequence  $\{\mathbf{x}_i\}$ ,  $\{\tilde{\mathbf{x}}_i\}$ , respectively, generates two sequences of real numbers  $\{E_i\}$  and  $\{s_i\}$ ,  $\{\tilde{E}_i\}$  and  $\{\tilde{s}_i\}$ , respectively, which are defined by  $E_i = \tilde{H}(\mathbf{x}_i)$  and  $s_i = \psi(\mathbf{x}_i)$ , and by  $\tilde{E}_i = \tilde{H}(\tilde{\mathbf{x}}_i)$  and  $\tilde{s}_i = \psi(\tilde{\mathbf{x}}_i)$ , respectively. We obviously have  $\tilde{E}_i = 0$  and  $\tilde{s}_i = 0$ , while

<sup>8</sup>The resulting map is discontinuous at multiples of the step size  $\delta t$ . A smooth interpolation could be defined. But this is not needed in the context of this paper.



the “drift” in the values of  $E_i$  and  $s_i$  per step away from zero is exponentially small and sums up linearly with the number of steps. This energy conserving property of a symplectic method has been discussed by Benettin and Giorgilli [5] and Hairer and Lubich [18]. The same exponentially slow drift follows for the sequence  $\{s_i\}$  from

$$\begin{aligned} |\psi(\mathbf{x}_N) - \psi(\mathbf{x}_0)| &\leq \sum_{i=1}^N |\psi(\mathbf{x}_i) - \psi(\mathbf{x}_{i-1})| \\ &\leq \sum_{i=1}^N |\psi(\widehat{\Pi}(\mathbf{x}_{i-1})) - \psi(\widetilde{\Pi}(\mathbf{x}_{i-1}))| \\ &\leq N \lambda d_3 e^{-p} e^{-d_2/\delta t}, \end{aligned}$$

$\lambda > 0$  the Lipschitz constant of  $\psi$  on  $\widetilde{\mathcal{V}}$ .

In other words, if we start initially on  $\widetilde{\mathcal{D}}$ , then the points computed “numerically” with the Poincaré map  $\widehat{\Pi}$  will stay in an exponentially small neighborhood of  $\mathcal{D}$  over exponentially many iterates of  $\widehat{\Pi}$ . Now, since our numerical method is of order  $p \geq 1$ , the compact manifolds  $\widetilde{\mathcal{M}}_E$  and  $\mathcal{M}_E$  are  $\mathcal{O}(\delta t^p)$  away from each other; i.e., the modified Hamiltonian  $\widetilde{H}$  is an  $\epsilon$ -perturbation of the Hamiltonian  $H$  with  $\epsilon \sim \delta t^p$ . Thus the sequence  $\{\mathbf{x}_i\}$  will also stay in an  $\mathcal{O}(\delta t^p)$  neighborhood of  $\mathcal{D}$  as long as the number of iterates  $N$  satisfies

$$(5.21) \quad N \leq d_4 e^{d_2/(2\delta t)},$$

$d_4 > 0$  an appropriate constant.

Now the *shadowing lemma* [33] is applied to the sequence  $\{\mathbf{x}_i\}_{i=1,\dots,N}$ .

PROPOSITION 1. *There exists an exact trajectory  $\{\widehat{\mathbf{x}}_i\}_{i=1,\dots,N}$  of the Poincaré map  $\Pi_0$  on  $\mathcal{D}$  such that the “numerically” computed sequence  $\{\mathbf{x}_i\}_{i=1,\dots,N}$  stays in a  $\mathcal{O}(\delta t^p)$  neighborhood of the (shadowing) exact trajectory provided the number of iterates  $N$  satisfies (5.21).*

*Proof.* We first project the sequence  $\{\mathbf{x}_i\}_{i=1,\dots,N}$  down onto  $\mathcal{D}$  using a “search direction” orthogonal to the manifold  $\mathcal{D}$ . Denote the result by  $\{\bar{\mathbf{x}}_i\}$ . The projected sequence  $\{\bar{\mathbf{x}}_i\}$  and the sequence  $\{\mathbf{x}_i\}$  are  $\mathcal{O}(\delta t^p)$  close to each other provided  $N$  satisfies (5.21). The “local” error per step between the “exact” Poincaré map  $\Pi$  and the “numerical” Poincaré map  $\widehat{\Pi}$  is also of order  $p$  in the step size  $\delta t$ . This follows from standard forward error analysis. Thus the shadowing lemma [33] for uniformly hyperbolic maps can be applied to the Poincaré map  $\Pi_0 : \mathcal{D} \rightarrow \mathcal{D}$  and the projected sequence  $\{\bar{\mathbf{x}}_i\}$  on  $\mathcal{D}$ . The shadowing distance is  $\mathcal{O}(\delta t^p)$ . This shadowing result also applies to the sequence  $\{\mathbf{x}_i\}$ .  $\square$

Let us now assume that we want to compute the ensemble average of a smooth function  $A : \mathbb{R}^{2n} \rightarrow \mathbb{R}$  up to a certain accuracy  $c > 0$ . The large deviation theorem (5.20) for hyperbolic maps tells us that the probability to obtain the ensemble average in the desired accuracy as the time average along a single trajectory goes to one exponentially fast as the length  $N$  of the trajectory is increased. If we numerically compute an approximative trajectory for the system (5.18)–(5.19), then we know from Proposition 1 that this trajectory is  $\mathcal{O}(\delta t^p)$  close to *some* exact trajectory over exponentially many integration steps  $N$ . Let us denote the time average of  $A$  along this exact trajectory by  $\langle A \rangle_N^e$  and the numerically computed time average by  $\langle A \rangle_N$ ; then

$$(5.22) \quad \langle A \rangle_N - \langle A \rangle_N^e = \mathcal{O}(\delta t^p)$$

for all  $N$  satisfying a bound of type (4.13). Thus we obtain the following.

PROPOSITION 2. *Let (5.18)–(5.19) be a Poincaré hyperbolic (real-analytic) system, which we discretize by a symplectic method of order  $p \geq 1$  in the step size  $\delta t$ . Then the time average  $\langle A \rangle_N$  of an observable  $A$  along a “numerically” computed trajectory  $\{\mathbf{x}_n\}_{n=1, \dots, N}$ ,*

$$\mathbf{x}_{n+1} = \Psi_{\delta t}(\mathbf{x}_n),$$

*satisfies (5.22), where  $\langle A \rangle_N^e$  is the time average along some exact trajectory and the number of steps  $N$  satisfies a bound of type (5.21). Furthermore, assume we want to compute the ensemble average of  $A$  within a given accuracy  $c > 0$ . We assume, for simplicity, that the constant  $c$  is larger than the difference between the time averages (5.22), which is always true for sufficiently small step sizes  $\delta t$ . Then the probability to obtain the average in the desired accuracy as the time average along a numerically computed trajectory goes to one exponentially fast as the number of integration steps  $N$  is increased. Taking the maximum number (5.21) of steps, the probability can be made double exponentially close to one in (5.20) as  $\delta t \rightarrow 0$ .*

*Example.* As a numerical example, we look at the following planar anisotropic Kepler problem [15]:

$$\begin{aligned} \frac{d}{dt} \mathbf{q} &= \mathbf{M}^{-1} \mathbf{p}, \\ \frac{d}{dt} \mathbf{p} &= -\nabla_{\mathbf{q}} V(\mathbf{q}), \end{aligned}$$

$$\mathbf{q} = (q_x, q_y)^T, \mathbf{p} = (p_x, p_y)^T \in \mathbb{R}^2,$$

$$V(q_x, q_y) = \frac{-1}{\sqrt{(q_x)^2 + (q_y)^2}},$$

and

$$\mathbf{M} = \begin{bmatrix} 10 & 0 \\ 0 & 1 \end{bmatrix}.$$

The initial conditions are chosen such that

$$(5.23) \quad H = \frac{\mathbf{p}^T \mathbf{M}^{-1} \mathbf{p}}{2} + V(\mathbf{q}) = -\frac{1}{2}$$

and  $L = p_y q_x - p_x q_y \neq 0$ . Note that angular momentum  $L$  is not conserved.

We define the Poincaré section  $\mathcal{S}_0$  by  $\psi = q_y = 0$  and record the sequence of points  $(q_x, p_x)$ . Conservation of energy implies that the thus-defined sequence is restricted to the subset

$$|q_x| < \frac{2}{1 + (p_x)^2/10},$$

where  $-\infty < p_x < +\infty$ . This subset has an awkward shape. But it can be transformed into a rectangle by means of the area preserving transformation

$$(5.24) \quad X_1 := q_x (1 + (p_x)^2/10),$$

$$(5.25) \quad X_2 := \sqrt{10} \arctg(p_x/\sqrt{10}),$$

where  $|X_1| \leq 2$  and  $|X_2| \leq \sqrt{10}\pi/2$ . The corresponding Poincaré map is hyperbolic; i.e., stable and unstable manifolds intersect transversally, and the dynamics can be encoded in a binary Bernoulli shift [15].

The main computational difficulty consists in the existence of a weak singularity at  $\mathbf{q} = \mathbf{0}$ . To remove this singularity, we have to scale the equations of motion by introducing the time transformation

$$dt = \rho(\mathbf{q}) d\tau, \quad \rho(\mathbf{q}) := (q_x)^2 + (q_y)^2,$$

which implies that, in the new time  $\tau$ , the norm of the vector field remains bounded at  $\mathbf{q} = \mathbf{0}$ . The time transformation has to be introduced such that the transformed equations of motion are still Hamiltonian. A constant step-size symplectic method can then be used to integrate the transformed system. Let us describe the general approach: Assume that a Hamiltonian function  $H(\mathbf{q}, \mathbf{p}) = \mathbf{p}^T \mathbf{M}^{-1} \mathbf{p}/2 + V(\mathbf{q})$  and a scaling function  $\rho(\mathbf{q})$  are given. Following Zare and Szebehely [41], we introduce the modified Hamiltonian function

$$E(\mathbf{q}, \mathbf{p}, t, e) := \rho(\mathbf{q}) [H(\mathbf{q}, \mathbf{p}) - e]$$

with corresponding Hamiltonian equations of motion

$$(5.26) \quad \frac{d}{d\tau} \mathbf{q} = \rho(\mathbf{q}) \mathbf{M}^{-1} \mathbf{p},$$

$$(5.27) \quad \frac{d}{d\tau} \mathbf{p} = -\rho(\mathbf{q}) \nabla_{\mathbf{q}} V(\mathbf{q}) - [H(\mathbf{q}, \mathbf{p}) - e] \nabla_{\mathbf{q}} \rho(\mathbf{q}),$$

$$(5.28) \quad \begin{aligned} \frac{d}{d\tau} t &= \rho(\mathbf{q}), \\ \frac{d}{d\tau} e &= 0 \end{aligned}$$

in extended phase space  $\mathbb{R}^{2n} \times \mathbb{R}^2$ . In particular, let us consider the case  $e = H(\mathbf{q}(0), \mathbf{p}(0))$ . Then (5.26)–(5.28) can be simplified to

$$\begin{aligned} \frac{d}{d\tau} \mathbf{q} &= \rho(\mathbf{q}) \mathbf{M}^{-1} \mathbf{p}, \\ \frac{d}{d\tau} \mathbf{p} &= -\rho(\mathbf{q}) \nabla_{\mathbf{q}} V(\mathbf{q}), \\ \frac{d}{d\tau} t &= \rho(\mathbf{q}) \end{aligned}$$

on the hypersurface of constant energy  $E = 0$ . This is a scaled vector field as desired that is not Hamiltonian anymore. Therefore, as suggested by the author in [27] and independently by Hairer [17], the Hamiltonian equations (5.26)–(5.28) are discretized by a symplectic method and  $e = H(\mathbf{q}_0, \mathbf{p}_0)$ . For example, the equations can be discretized by the symplectic Euler method; i.e.,

$$\begin{aligned} \mathbf{q}_{n+1} &= \mathbf{q}_n + \delta\tau \rho(\mathbf{q}_n) \mathbf{M}^{-1} \mathbf{p}_{n+1}, \\ \mathbf{p}_{n+1} &= \mathbf{p}_n - \delta\tau \rho(\mathbf{q}_n) \nabla_{\mathbf{q}} V(\mathbf{q}_n) - \delta\tau (H(\mathbf{q}_n, \mathbf{p}_{n+1}) - e) \nabla_{\mathbf{q}} \rho(\mathbf{q}_n), \\ t_{n+1} &= t_n + \delta\tau \rho(\mathbf{q}_n). \end{aligned}$$

The method is explicit in the variable  $\mathbf{q}$ . Unfortunately this implies that the method is only first order in  $\delta\tau$ . However, the method is symplectic and, therefore, the Hamiltonian  $E = [H(\mathbf{q}, \mathbf{p}) - e]\rho(\mathbf{q})$  is conserved to  $\mathcal{O}(\delta\tau)$  over exponentially long periods of

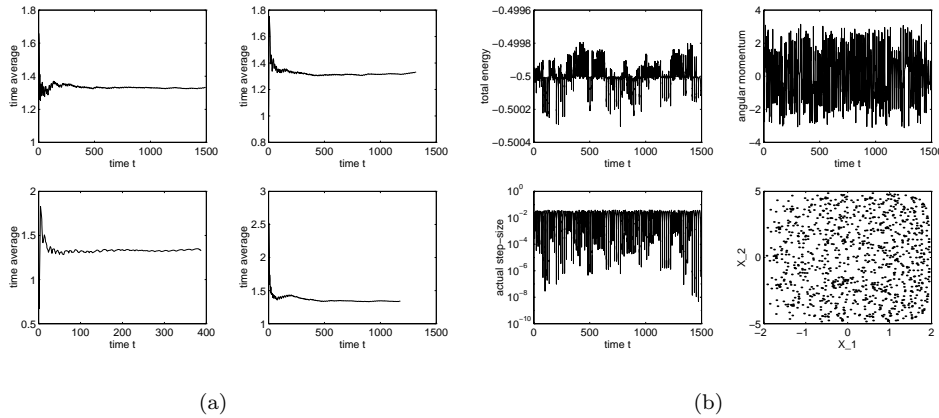


FIG. 1. (a) The time evolution of the average  $\langle r \rangle_n$  (mean distance) is shown for four different initial conditions with equal initial energy  $e = -1/2$ . (b) Time evolution of the error in energy, angular momentum, and the actual step size. The bottom-right figure shows the intersections of the trajectory with the Poincaré section in the  $(X_1, X_2)$  coordinates. One thousand intersections are plotted.

time. A second-order symplectic discretization can be obtained by using the second-order Lobatto IIIa–b partitioned Runge–Kutta formula [34], i.e.,

$$(5.29) \quad \mathbf{p}_{n+1/2} = \mathbf{p}_n - \frac{\delta\tau}{2} [\rho(\mathbf{q}_n) \nabla_{\mathbf{q}} V(\mathbf{q}_n) - [H(\mathbf{q}_n, \mathbf{p}_{n+1/2}) - e] \nabla_{\mathbf{q}} \rho(\mathbf{q}_n)],$$

$$(5.30) \quad \mathbf{q}_{n+1} = \mathbf{q}_n + \frac{\delta\tau}{2} [\rho(\mathbf{q}_{n+1}) + \rho(\mathbf{q}_n)] \mathbf{M}^{-1} \mathbf{p}_{n+1/2},$$

$$(5.31) \quad \mathbf{p}_{n+1} = \mathbf{p}_{n+1/2} - \frac{\delta\tau}{2} [\rho(\mathbf{q}_{n+1}) \nabla_{\mathbf{q}} V(\mathbf{q}_{n+1}) - [H(\mathbf{q}_{n+1}, \mathbf{p}_{n+1/2}) - e] \nabla_{\mathbf{q}} \rho(\mathbf{q}_{n+1})],$$

$$(5.32) \quad t_{n+1} = t_n + \frac{\delta\tau}{2} [\rho(\mathbf{q}_n) + \rho(\mathbf{q}_{n+1})].$$

The resulting scheme is implicit in  $\rho(\mathbf{q})$ .

This approach is applied to the anisotropic Kepler problem with a scaling function  $\rho(\mathbf{q}) = (q_x)^2 + (q_y)^2$ . We chose initial values such that  $e = H = -1/2$  and  $L \neq 0$ . The equations of motion are integrated using the second-order symplectic method (5.29)–(5.32) with  $\delta\tau = 0.05$ . The time average of an observable  $A(\mathbf{q})$  along a trajectory  $\{\mathbf{q}_n\}_{n=1, \dots, M}$  is computed according to the recursive formula

$$\langle A \rangle_n = \frac{1}{t_n} [t_{n-1} \langle A \rangle_{n-1} + \delta t_n A(\mathbf{q}_n)].$$

The time average of  $r = \sqrt{(q_x)^2 + (q_y)^2}$  (mean distance) was computed for four different initial conditions, and the evolution of the corresponding time averages  $\langle r \rangle_n$  can be found in Figure 1(a). The different lengths of the time intervals are due to the fact that the same number of steps with step size  $\delta\tau$  were taken, which leads to different actual step sizes  $\delta t_n$ . Within a tolerance of  $c = 0.04$ , these averages converge to the same value  $\approx 1.33$ . In Figure 1(b), the total energy  $H$ , the angular momentum  $L$ , and the variation in the actual step size  $\delta t = r^2 \delta\tau$  can be found for a particular

trajectory. We also plotted the intersections of the trajectory with the  $(q_x, p_x)$  plane in the  $(X_1, X_2)$  representation (5.24)–(5.25). Theoretically, these points should fill the rectangle in a uniform way (the invariant measure is the Lebesgue measure). With one thousand points plotted, the uniform distribution is satisfied quite well.

**Acknowledgments.** I would like to thank Ernst Hairer, Chus Sanz-Serna, Andrew Stuart, and Claudia Wulff for comments on an earlier version of this paper and the referees for making valuable suggestions.

## REFERENCES

- [1] M. ADAMS, T. RATIU, AND R. SCHMID, *The Lie group structure of diffeomorphism groups and invertible Fourier integrals operators with Applications*, in *Infinite Dimensional Groups with Applications*, V. Kac, ed., Springer-Verlag, New York, 1985.
- [2] V. I. ARNOLD, *Geometrical Methods in the Theory of Ordinary Differential Equations*, 2nd ed., Springer-Verlag, New York, 1987.
- [3] V. I. ARNOLD, *Mathematical Methods of Classical Mechanics*, 2nd ed., Springer-Verlag, New York, 1988.
- [4] S. P. AUERBACH, AND A. FRIEDMAN, *Long-time behaviour of numerically computed orbits: Small and intermediate time-step analysis of one-dimensional systems*, *J. Comput. Phys.*, 93 (1991), pp. 189–223.
- [5] G. BENETTIN AND A. GIORGILLI, *On the Hamiltonian interpolation of near to the identity symplectic mappings with application to symplectic integration algorithms*, *J. Statist. Phys.*, 74 (1994), pp. 1117–1143.
- [6] W.-J. BEYN, *Numerical methods for dynamical systems*, in *Advances in Numerical Analysis*, Vol. I, Clarendon Press, Oxford, 1991.
- [7] M. P. CALVO, A. MURUA, AND J. M. SANZ-SERNA, *Modified Equations for ODEs*, *Contemp. Math.*, 172 (1994), pp. 63–74.
- [8] M. P. CALVO, A. ISERLES, AND A. ZANNA, *Runge–Kutta methods on manifolds*, in *Numerical Analysis: A. R. Mitchell’s 75th Birthday Volume*, G. A. Watson and D. F. Griffiths, eds., World Scientific, Singapore, 1997, pp. 57–70.
- [9] T. EIROLA, *Aspects of backward error analysis of numerical ODEs*, *J. Comput. Appl. Math.*, 45 (1993), pp. 65–73.
- [10] K. FENG, *Formal power series and numerical algorithms for dynamical systems*, in *Proceedings of International Conference on Scientific Computation*, T. Chan and Z.-C. Shi, eds., Ser. Appl. Math. 1, World Scientific, Singapore, 1991, pp. 28–35.
- [11] B. FIEDLER AND J. SCHEURLE, *Discretization of homoclinic orbits, rapid forcing, and “invisible” chaos*, *Mem. Amer. Math. Soc.*, 119 (1996).
- [12] D. F. GRIFFITHS AND J. M. SANZ-SERNA, *On the scope of the modified equations*, *J. Sci. Statist. Comput.*, 7 (1986), pp. 994–1008.
- [13] O. GONZALEZ, D. J. HIGHAM, AND A. STUART, *Qualitative properties of modified equations*, *IMA J. Numer. Anal.*, 19 (1999), pp. 169–190.
- [14] J. GUCKENHEIMER AND P. HOLMES, *Nonlinear Oscillations, Dynamical Systems and Bifurcations of Vector Fields*, Springer-Verlag, New York, 1983.
- [15] M. C. GUTZWILLER, *Chaos in Classical and Quantum Mechanics*, Springer-Verlag, New York, 1990.
- [16] E. HAIRER, *Backward analysis of numerical integrators and symplectic methods*, *Ann. Numer. Math.*, 1 (1994), pp. 107–132.
- [17] E. HAIRER, *Variable time step integration with symplectic methods*, *Appl. Numer. Math.*, 25 (1997), pp. 219–227.
- [18] E. HAIRER AND CH. LUBICH, *The life-span of backward error analysis for numerical integrators*, *Numer. Math.*, 76 (1997), pp. 441–462.
- [19] E. HAIRER, S. P. NØRSETT, AND G. WANNER, *Solving Ordinary Differential Equations, Vol. I*, 2nd ed., Springer-Verlag, Berlin, New York, 1993.
- [20] E. HAIRER AND D. STOFFER, *Reversible long-term integration with variable stepsizes*, *SIAM J. Sci. Comput.*, 18 (1997), pp. 257–269.
- [21] R. HAMILTON, *The inverse function theorem of Nash and Moser*, *Bull. Amer. Math. Soc.*, 7 (1982), pp. 65–222.
- [22] K. JÄNICH, *Funktionentheorie*, 3rd ed., Springer-Verlag, Berlin, New York, 1993.

- [23] F. M. LASAGNI, *Canonical Runge–Kutta methods*, Z. Angew. Math. Phys., 5, 39 (1988), pp. 952–953.
- [24] J. MOSER, *Lectures on Hamiltonian systems*, Mem. Amer. Math. Soc., 81 (1968), pp. 1–60.
- [25] A. I. NEISHTADT, *The separation of motions in systems with rapidly rotating phase*, J. Appl. Math. Mech., 48 (1984), pp. 133–139.
- [26] S. REICH, *Numerical Integration of Generalized Euler Equations*, Tech. Report 93-20, University of British Columbia, Vancouver, British Columbia, Canada, 1993.
- [27] S. REICH, *Backward Error Analysis for Numerical Integrators*, Preprint SC 96-21, Konrad-Zuse-Zentrum, Berlin, 1996.
- [28] S. REICH, *On higher order semi-explicit symplectic partitioned Runge–Kutta methods for constrained Hamiltonian systems*, Numer. Math., 76 (1997), pp. 231–247.
- [29] S. REICH, *Preservation of adiabatic invariants under symplectic discretization*, Appl. Numer. Math., 29 (1999), pp. 45–55.
- [30] S. REICH, *Dynamical Systems, Exponentially Small Estimates, and Numerical Integration*, Habilitationsschrift, FU Berlin, 1998.
- [31] J. M. SANZ-SERNA, *Symplectic integrators for Hamiltonian problems: An overview*, Acta Numer., 1 (1992), pp. 243–286.
- [32] J. M. SANZ-SERNA AND M. P. CALVO, *Numerical Hamiltonian Problems*, Chapman and Hall, London, 1994.
- [33] T. SAUER AND J. A. YORK, *Rigorous verification of trajectories for the computer simulation of dynamical systems*, Nonlinearity, 4 (1994), pp. 961–979.
- [34] G. SUN, *Symplectic partitioned Runge–Kutta methods*, J. Comput. Math., 11 (1993), pp. 365–372.
- [35] V. S. VARADARAJAN, *Lie Groups, Lie Algebras, and Their Representation*, Prentice–Hall, Englewood Cliffs, NJ, 1974.
- [36] M. VIANA, *Stochastic Dynamics of Deterministic Systems*, Instituto de Matemática Pura e Aplicada (IMPA), Rio de Janeiro, 1997.
- [37] P. WALTERS, *Introduction to Ergodicity Theory*, 2nd ed., Springer-Verlag, New York, 1985.
- [38] R. F. WARMING AND B. J. HYETT, *The modified equation approach to the stability and accuracy of finite-difference methods*, J. Comput. Phys., 14 (1974), pp. 159–179.
- [39] L.-S. YOUNG, *Large deviations in dynamical systems*, Trans. Amer. Math. Soc., 318 (1990), pp. 525–543.
- [40] H. YOSHIDA, *Construction of higher order symplectic integrators*, Phys. Lett. A, 150 (1990), pp. 262–268.
- [41] K. ZARE AND V. SZEBEHELY, *Time transformations for the extended phase space*, Celestial Mech., 11 (1975), pp. 469–482.