

Higher-Order Averaging, Formal Series and Numerical Integration I: B-series

P. Chartier · A. Murua · J.M. Sanz-Serna

Received: 13 October 2009 / Revised: 7 May 2010 / Accepted: 20 May 2010 /
Published online: 22 June 2010
© SFoCM 2010

Abstract We show how B-series may be used to derive in a systematic way the analytical expressions of the high-order stroboscopic averaged equations that approximate the slow dynamics of highly oscillatory systems. For first-order systems we give explicitly the form of the averaged systems with $\mathcal{O}(\epsilon^j)$ errors, $j = 1, 2, 3$ ($2\pi\epsilon$ denotes the period of the fast oscillations). For second-order systems with large $\mathcal{O}(\epsilon^{-1})$ forces, we give the explicit form of the averaged systems with $\mathcal{O}(\epsilon^j)$ errors, $j = 1, 2$. A variant of the Fermi–Pasta–Ulam model and the inverted Kapitza pendulum are used as illustrations. For the former it is shown that our approach establishes the adiabatic invariance of the oscillatory energy. Finally we use B-series to analyze multiscale numerical integrators that implement the method of averaging. We construct

Dedicated to Ernst Hairer on his 60th birthday.

Communicated by Arieh Iserles.

P. Chartier (✉)

INRIA Rennes, ENS Cachan Bretagne, Campus de Ker-Lann, av. Robert Schuman,
35170 Bruz, France
e-mail: Philippe.Chartier@irisa.fr

A. Murua

Konputazio Zientziak eta A.A. Saila, Informatika Fakultatea, UPV/EHU,
20018 Donostia–San Sebastián, Spain
e-mail: Ander.Murua@ehu.es

J.M. Sanz-Serna

Departamento de Matemática Aplicada, Facultad de Ciencias, Universidad de Valladolid, Valladolid,
Spain
e-mail: sanzsern@mac.uva.es

integrators that are able to approximate not only the simplest, lowest-order averaged equation but also its high-order counterparts.

Keywords Averaging · High-order stroboscopic averaging · Highly oscillatory problems · Hamiltonian problems · Multiscale numerical methods · Numerical integrators · Formal series · B-series · Trees · Fermi–Pasta–Ulam problem · Adiabatic invariants · Inverted Kapitza’s pendulum

Mathematics Subject Classification (2000) 34C29 · 65L06 · 34D20 · 70H05 · 79K65

1 Introduction

The aim of this series of papers is to relate the method of averaging to the formal series expansions that are nowadays used as a powerful tool in the analysis of numerical integrators of time-dependent problems. The present work is restricted to B-series and systems with a single fast frequency; Part II will deal with other types of formal series and with quasiperiodic problems.

The method of averaging [14, 26, 27] has a long history that goes back to the work in celestial mechanics of Gauss and Laplace [2]. The aim is to study the long-time behavior of highly oscillatory systems by constructing an averaged system that approximately captures that behavior and ignores the details of the oscillations with small periods.

In 1974 Hairer and Wanner [17] introduced the concept of B-series. They associated with each numerical integrator (within a very broad class) its B-series, a formal series in powers of the step-length, and showed that the properties of the integrator are easily studied by manipulating the corresponding B-series. The importance of the notion of B-series and other similar formal series [25] has grown steadily in recent years [19], in particular since the discovery of their relevance in connection to symplectic integration [4] and modified equations [15] (see also the recent contribution [9]). Even though B-series are formal series, they may lead to rigorous error bounds when suitably truncated [19].

In Sects. 2–4 of this paper we show how B-series may be used to derive in a systematic way the analytical expressions of high-order averaged equations. This is probably the first example of the application of B-series outside the field of numerical ordinary differential equations. Section 2 describes the general idea and Sects. 3 and 4 provide the details for large classes of first-order and second-order differential systems respectively. For first-order systems we give explicitly the form of the averaged systems with $\mathcal{O}(\epsilon^j)$ errors, $j = 1, 2, 3$ ($2\pi\epsilon$ denotes the period of the fast oscillations). When the original oscillatory problem is Hamiltonian, so are the averaged systems and explicit expressions for the corresponding Hamiltonian functions are also provided. For second-order systems with large, $\mathcal{O}(\epsilon^{-1})$ forces, we give the explicit form of the averaged systems with $\mathcal{O}(\epsilon^j)$ errors, $j = 1, 2$. A family of Hamiltonian problems that includes a variant of the Fermi–Pasta–Ulam model is employed to illustrate the material in Sect. 3. We show that our methodology, in addition to determining averaged equations of high order, yields as a byproduct the adiabatic invariance of the oscillatory energy. The material in Sect. 4 is exemplified in the case of Kapitza’s inverted pendulum.

The final Sect. 5 deals with the analysis of numerical multiscale methods such as those introduced by E and Engquist [10, 11]. While the scope of application of those methods widely exceeds the oscillatory problems considered in the present article, they are of particular relevance to our research because, when applied to the classes of systems studied here, they may be considered as numerical implementations of the idea of averaging. We shall construct multiscale methods that are able to approximate not only the simplest, lowest-order averaged equation but also its high-order counterparts.

We end this introduction by pointing out that there are of course analogies between the B-series approach considered here and both the modulated Fourier expansion methodology pioneered by Hairer and Lubich [16] and the WKB and Magnus expansion techniques [20].

2 A Modified Equation Approach to Averaging

This section relates well-known ideas from the theory of averaging [14, 26, 27] to modified equations as used in the analysis of numerical integrators [19, 30].

We are concerned with initial value problems for differential systems of the form

$$\frac{d}{dt}y = f\left(y, \frac{t}{\epsilon}; \epsilon\right), \tag{1}$$

where y is a D -dimensional real vector, ϵ is a small parameter and the indefinitely differentiable function f is assumed to depend 2π -periodically on the variable t/ϵ . Our interest is in situations where, as $\epsilon \rightarrow 0$, the solutions or some of their derivatives with respect to t become *unbounded*. Such is the case, for instance, when f and its partial derivatives are $\mathcal{O}(1)$ and the partial derivative of f with respect to t/ϵ does not vanish (so that (1) is effectively non-autonomous); then differentiation in (1) shows that $(d^2/dt^2)y = \mathcal{O}(1/\epsilon)$. But even in cases where (1) is autonomous, so that there is no effective dependence of f on the fast time t/ϵ , the derivatives of y will become unbounded as $\epsilon \rightarrow 0$ if f itself behaves like a negative power of ϵ .

For the analysis it is useful to rewrite the system in terms of the scaled (non-dimensional) time $\tau = t/\epsilon$:

$$\frac{d}{d\tau}y = \epsilon f(y, \tau; \epsilon). \tag{2}$$

If we denote by $\varphi_{\tau_0, \tau; \epsilon} : \mathcal{R}^D \rightarrow \mathcal{R}^D$ the solution operator of (2), so that

$$y(\tau) = \varphi_{\tau_0, \tau; \epsilon}(y_0)$$

is the solution that satisfies the initial condition $y(\tau_0) = y_0$, then a simple but important observation is that the map $\Psi_{\tau_0; \epsilon} = \varphi_{\tau_0, \tau_0+2\pi; \epsilon}$ depends on τ_0 in a 2π -periodic manner; this follows from the fact that both $\varphi_{\tau_0, \tau; \epsilon}(y_0)$ and $\varphi_{\tau_0+2\pi, \tau+2\pi; \epsilon}(y_0)$ satisfy the initial value problem

$$\begin{cases} \frac{d}{d\tau}y(\tau) = \epsilon f(y(\tau), \tau; \epsilon) = \epsilon f(y(\tau), \tau + 2\pi; \epsilon), \\ y(\tau_0) = y_0. \end{cases}$$

From this observation, it follows that, at the stroboscopic times $\tau_n = \tau_0 + 2\pi n$, $n = 0, \pm 1, \pm 2, \dots$,

$$y(\tau_n) = \varphi_{\tau_0, \tau_n; \epsilon}(y_0) = \varphi_{\tau_{n-1}, \tau_n; \epsilon}(\varphi_{\tau_0, \tau_{n-1}; \epsilon}(y_0)) = \varphi_{\tau_0, \tau_0+2\pi; \epsilon}(\varphi_{\tau_0, \tau_{n-1}; \epsilon}(y_0))$$

and, hence, we arrive at the fundamental formula:

$$y(\tau_n) = (\Psi_{\tau_0; \epsilon})^n(y_0), \quad n = 0, \pm 1, \pm 2, \dots \tag{3}$$

In Sects. 3 and 4 we shall describe general situations where the exact solution of (2) with initial value $y(\tau_0) = y_0$, sampled at $\tau_n = \tau_0 + 2\pi n$, admits a formal expansion of the form

$$y(\tau_n) = y_0 + \sum_{j=1}^{\infty} \epsilon^j \sum_{l=1}^j (\tau_n - \tau_0)^l G_{j,l}(y_0), \tag{4}$$

with suitable indefinitely differentiable maps $G_{j,l} : \mathcal{R}^D \rightarrow \mathcal{R}^D$ independent of ϵ . In particular, this implies that

$$\Psi_{\tau_0; \epsilon}(y_0) = y_0 + \sum_{j=1}^{\infty} \epsilon^j \sum_{l=1}^j (2\pi)^l G_{j,l}(y_0),$$

and thus $\Psi_{\tau_0; \epsilon}$ is a smooth near-to-identity map. Standard backward error analysis [19, 30] then shows the existence of an *autonomous* system (the modified system of $\Psi_{\tau_0; \epsilon}$)

$$\frac{d}{d\tau} Y = F(Y; \epsilon) = \epsilon F_1(Y) + \epsilon^2 F_2(Y) + \epsilon^3 F_3(Y) + \dots \tag{5}$$

or

$$\frac{d}{dt} Y = \frac{1}{\epsilon} F(Y; \epsilon) = F_1(Y) + \epsilon F_2(Y) + \epsilon^2 F_3(Y) + \dots \tag{6}$$

(F and the F_j depend on τ_0 , but this has not been incorporated into the notation) whose (formal) solutions satisfy that $Y(\tau_n) = \Psi_{\tau_0; \epsilon}(Y(\tau_{n-1}))$ for $n = 0, \pm 1, \pm 2, \dots$ so that

$$Y(\tau_n) = (\Psi_{\tau_0; \epsilon})^n(Y_0) \quad n = 0, \pm 1, \pm 2, \dots \tag{7}$$

We conclude from (3) and (7) that, if one chooses $Y(\tau_0) = y(\tau_0)$, then $Y(\tau)$ exactly coincides with $y(\tau)$ at the stroboscopic times $\tau_n = \tau_0 + 2\pi n$. In this way it is possible in principle to find $y(\tau_n)$ by solving the system (5) or (6), where all t -derivatives of Y remain bounded as $\epsilon \rightarrow 0$. When τ does not coincide with one of the stroboscopic times, we note that obviously

$$y(\tau) = (\varphi_{\tau_n, \tau; \epsilon} \circ \Phi_{\tau_n - \tau; \epsilon})Y(\tau),$$

where τ_n is the largest stroboscopic time $\leq \tau$ and $\Phi_{\cdot; \epsilon}$ denotes the flow of (5). In this way, y is ‘enslaved’ to Y through the mapping $\varphi_{\tau_n, \tau; \epsilon} \circ \Phi_{\tau_n - \tau; \epsilon}$ whose dependence on τ is easily seen to be 2π -periodic.

It is well known that the series (5) does not converge in general, and in order to get rigorous results one has to consider truncated versions

$$\frac{d}{d\tau} Y = \epsilon F_1(Y) + \epsilon^2 F_2(Y) + \epsilon^3 F_3(Y) + \dots + \epsilon^J F_J(Y), \tag{8}$$

whose solutions satisfy that $Y(\tau_n) - \Psi_{\tau_0, \epsilon}(Y(\tau_{n-1})) = \mathcal{O}(\epsilon^{J+1})$. If Y solves (8) with $Y(\tau_0) = y(\tau_0)$, then $Y(\tau_n)$ and $y(\tau_n)$ differ by an $\mathcal{O}(\epsilon^J)$ amount, where the constant implied in the \mathcal{O} notation is uniform as the stroboscopic time τ_n ranges in an interval $\tau_0 \leq \tau_n \leq \tau_n + T/\epsilon$, with $T = \mathcal{O}(1)$ as $\epsilon \rightarrow 0$.¹

The process of obtaining the autonomous system (8) from the original system (2) is referred to in the averaging literature [14, 26, 27] as high-order stroboscopic averaging. In this paper we shall show how a number of techniques and results currently used in the analysis of numerical integrators may be applied to the task of constructing explicitly the functions F_j that define the averaged equations (8).

Let us briefly discuss how to find the modified system (5) corresponding to $\Psi_{\tau_0, \epsilon}$. There are several techniques that may be used to accomplish this task; here we shall apply Theorem 1 in [25]. It follows from (4) that the formal solution $Y(\tau)$ of the autonomous system (5) with initial value $Y(\tau_0) = y_0$ can be written as

$$Y(\tau) = y_0 + \sum_{j=1}^{\infty} \epsilon^j \sum_{l=1}^j (\tau - \tau_0)^l G_{j,l}(y_0). \tag{9}$$

(This is proved as follows. The difference $\Delta(\tau)$ between the formal solution Y and the right-hand side of (9) vanishes when τ coincides with one of the stroboscopic times τ_n . Therefore the polynomial $P(\tau)$ with degree $\leq N$ that interpolates $\Delta(\tau)$ at τ_0, \dots, τ_N , vanishes identically. On the other hand the interpolation error $\Delta - P = \Delta$ is easily seen to be $\mathcal{O}(\epsilon^{N+1})$ and, since N is arbitrary, $\Delta \equiv 0$.) The vector field F is then retrieved by differentiating the solution flow:²

$$\left. \frac{d}{d\tau} Y(\tau) \right|_{\tau=\tau_0} = F(y_0; \epsilon).$$

The following proposition lists some useful properties.

¹The size of $Y(\tau) - y(\tau)$ at non-stroboscopic times depends of course on the behavior of $\varphi_{\tau_n, \tau; \epsilon}$. If accurate approximations of $y(\tau)$ are required, they may be obtained by first finding $Y(\tau_n)$, where τ_n is the largest stroboscopic time $\leq \tau$, and then integrating the original oscillatory system from τ_n to τ .

²This procedure is of course classical: old differential equations textbooks used to show that any given family of sufficiently differentiable functions $Q = \mathcal{E}(\tau, C)$ with values in \mathcal{R}^D and depending on a vector parameter C in \mathcal{R}^D was the ‘general solution’ of the differential equation obtained by eliminating the parameters from the equations $Q = \mathcal{E}$, $dQ/d\tau = \partial \mathcal{E} / \partial \tau$. In general the resulting differential equation is non-autonomous. If it is known beforehand—as is the case here—that the differential equation will turn out to be autonomous, then it is sufficient to consider $Q = \mathcal{E}$, $dQ/d\tau = \partial \mathcal{E} / \partial \tau$ at $\tau = 0$. In Sect. 4 we shall use the corresponding idea for second-order systems; these have as ‘general solution’ families $Q = \mathcal{E}(\tau, C_1, C_2)$ with C_1, C_2 in \mathcal{R}^D , and these are retrieved from \mathcal{E} by eliminating the parameters from the relations $Q = \mathcal{E}$, $(d/d\tau)Q = \partial \mathcal{E} / \partial \tau$, $(d^2/d\tau^2)Q = \partial^2 \mathcal{E} / \partial \tau^2$.

Proposition 1

- (1) *If the real function $I(y, \tau)$ depends 2π -periodically on τ and is conserved by all solutions of the original system (2) ($I(y(\tau), \tau) = I(y(\tau_0), \tau_0)$), then $I(Y, \tau_0)$ is a (formal) conserved quantity for the averaged system (5).*
- (2) *If (2) is Hamiltonian, then so is (5). More precisely, there is a formal Hamiltonian*

$$\mathcal{H}(Y; \epsilon) = \epsilon \mathcal{H}_1(Y) + \epsilon^2 \mathcal{H}_2(Y) + \epsilon^3 \mathcal{H}_3(Y) + \dots$$

such that $F(Y; \epsilon) = \mathcal{J}^{-1} \nabla \mathcal{H}(Y; \epsilon)$ and $F_j(Y) = \mathcal{J}^{-1} \nabla \mathcal{H}_j(Y)$ for each j (\mathcal{J} is the canonical symplectic matrix).

- (3) *If (2) is autonomous, then (5) is independent of τ_0 .*
- (4) *If (2) is autonomous and possesses a first integral $I(y)$, then $I(Y)$ is a formal first integral of (5).*
- (5) *If (2) is autonomous and Hamiltonian with Hamiltonian function H , then H is a formal first integral of (5). Furthermore, the Hamiltonian \mathcal{H} of (5) is a formal first integral of (2).*

Proof For (1) (which obviously implies (4)) note that $S(\tau) = I(Y(\tau), \tau_0) - I(Y(\tau_0), \tau_0)$ vanishes when τ coincides with one of the stroboscopic times τ_n . Then the interpolation argument used to prove (9) implies $S \equiv 0$.

If the original system is Hamiltonian (possibly non-autonomous), then $\Psi_{\tau; \epsilon}$ is a symplectic map and therefore its modified equation (5) will be Hamiltonian [19, 30]. This proves (2).

The property (3) is true because for an autonomous system $\varphi_{\tau_0, \tau; \epsilon}$ only depends on the difference $\tau - \tau_0$ and $\Psi_{\tau_0; \epsilon} = \varphi_{\tau_0, \tau_0 + 2\pi; \epsilon}$ is independent of τ_0 .

Under the hypotheses of (5), H is a first integral of the original system and, by (4), also of its averaged counterpart. But then H and \mathcal{H} are in involution (their Poisson bracket vanishes) [3] and therefore \mathcal{H} is a formal invariant of the original (2). □

Let us finish this section with a remark. If the initial value τ_0 of τ is changed to τ'_0 , then, as noted before, the averaged system (5) changes. The relation

$$\varphi_{\tau'_0, \tau'_0 + 2\pi; \epsilon} = \varphi_{\tau_0 + 2\pi, \tau'_0 + 2\pi; \epsilon} \circ \varphi_{\tau_0, \tau_0 + 2\pi; \epsilon} \circ \varphi_{\tau'_0, \tau_0; \epsilon},$$

i.e.

$$\Psi_{\tau'_0; \epsilon} = (\varphi_{\tau'_0, \tau_0; \epsilon})^{-1} \circ \Psi_{\tau_0; \epsilon} \circ \varphi_{\tau'_0, \tau_0; \epsilon},$$

shows that the new $\Psi_{\tau'_0; \epsilon}$ is conjugate to the old $\Psi_{\tau_0; \epsilon}$ by means of the change of variables $\varphi_{\tau'_0, \tau_0; \epsilon}$. Therefore the different averaged equations arising from different values of τ_0 are also conjugate to each other.

3 First-Order Differential Equations

3.1 Framework

Throughout this section we assume that the function f in (1) possesses an expansion in powers of ϵ of the form

$$f(y, \tau; \epsilon) = \frac{1}{\epsilon} \sum_{j=1}^{\infty} \epsilon^j f_j(y, \tau). \tag{10}$$

If

$$f_j(y, \tau) = \sum_{k=-\infty}^{\infty} \exp(ik\tau) f_{j,k}(y) \tag{11}$$

is the Fourier expansion of $f_j(y, \tau)$, then the Fourier coefficients $f_{j,k}(y)$ are in general complex vectors, but, in order to have a real system, we assume that, for each j and k , $f_{j,k} \equiv f_{j,-k}^*$ (* denotes complex conjugate). The system (1) is supplemented by the initial condition $y(0) = y_0$, where y_0 is a given vector of $\mathcal{O}(1)$ magnitude. Note that there is no loss of generality in assuming that the initial value of t is 0; the general case may be reduced to this by shifting t and redefining f .

In terms of the scaled time $\tau = t/\epsilon$, the differential system (2) being solved is then

$$\frac{d}{d\tau} y = \epsilon f(y, \tau; \epsilon) = \sum_{j=1}^{\infty} \epsilon^j f_j(y, \tau). \tag{12}$$

3.2 Trees

As is customary in the analysis of numerical methods for ordinary differential equations, the solution y of (12) will be expanded in an appropriate *B-series*: a formal series whose terms are indexed by (rooted) trees. In this subsection we describe the trees that will be required in the expansion of y which will be carried out in the next.

It is well known, see e.g. [18], that in the standard analysis of numerical integrators for systems $dy/d\tau = f(y)$, the vertices of the trees correspond to f and its derivatives. In view of the structure (10)–(11) of the right-hand side of (12), the vertices of the trees to be used here correspond to the functions $f_{j,k}$ and their derivatives; to keep track of this correspondence, each vertex possesses here a label, i.e. a pair of indices (j, k) , $j = 1, 2, \dots, k = 0, \pm 1, \pm 2, \dots$. Now the set \mathcal{U} of all trees may be defined recursively by the following two rules: (i) for each label (j, k) , the corresponding tree with one vertex, $\bullet jk00$, belongs to \mathcal{U} , (ii) if $u_1, \dots, u_n \in \mathcal{U}$, then, the result

$$u = [u_1, \dots, u_n]_{j,k} \tag{13}$$

of grafting their roots to a new root with label (j, k) belongs to \mathcal{U} .

The expansion of y will be graded according to powers of ϵ (see (14) below) and in this connection we introduce the notion of *weight*. The weight of a vertex with label (j, k) is defined to be j and the weight $|u|$ of $u \in \mathcal{U}$ is the sum of the weights of its vertices.

3.3 The Expansion of y

The solution y of (12) possesses an expansion

$$y(\tau) = y_0 + \sum_{u \in \mathcal{U}} \epsilon^{|u|} \frac{\alpha_u(\tau)}{\sigma_u} \mathcal{F}_u(y_0) = y_0 + \sum_{j=1} \epsilon^j \sum_{|u|=j} \frac{\alpha_u(\tau)}{\sigma_u} \mathcal{F}_u(y_0). \quad (14)$$

Here \mathcal{F}_u , σ_u and α_u are, respectively, the *elementary differential*, the *symmetry* and the *elementary coefficient* of the tree u . These will be described presently.

The elementary differential \mathcal{F}_u is, for each $u \in \mathcal{U}$, a vector-valued function of a vector argument y , constructed in terms of the Fourier coefficients $f_{j,k}$. Recursively, $\mathcal{F}(\bullet j, k00) = f_{j,k}$ and, for u in (13),

$$\mathcal{F}_u(y) = f_{j,k}^{(n)}(y)[F_{u_1}(y), \dots, F_{u_n}(y)],$$

where $f_{j,k}^{(n)}$ is the n th-order Fréchet derivative of $f_{j,k}$.

The integer σ_u counts the number of symmetries of the tree $u \in \mathcal{U}$ and may be defined recursively as follows. For trees with one vertex, $\sigma(\bullet j, k00) = 1$ and, for u in (13),

$$\sigma_u = r_1! \cdots r_m! \sigma_{u_1} \cdots \sigma_{u_m},$$

where u_μ , $\mu = 1, \dots, m$, are the pairwise distinct u_ν , $\nu = 1, \dots, n$, and r_μ counts the number of times that u_μ features among the u_ν .

For $u \in \mathcal{U}$, the elementary coefficient α_u is a complex-valued *function* of the real variable τ .³ Rewriting the differential equation (12) as an integral equation

$$y(\tau) = y_0 + \epsilon \int_0^\tau f(y(\tau'), \tau'; \epsilon) d\tau'$$

and expressing in both sides y by means of its expansion (14), leads to a Picard-like procedure for the recursive computation of the elementary coefficients. For the tree $\bullet j, k00$,

$$\alpha_u(\tau) = \int_0^\tau \exp(ik\tau') d\tau'$$

and, for u in (13),

$$\alpha_u(\tau) = \int_0^\tau \exp(ik\tau') \alpha_{u_1}(\tau') \cdots \alpha_{u_n}(\tau') d\tau'. \quad (15)$$

Two important remarks follow. First, note that the elementary coefficients in the expansion of y are universal in the sense that they are independent of the particular $f_{j,k}$ in the differential equation. This is one of the main advantages of the B-series

³Comparing the general term $\alpha_u/\sigma_u \mathcal{F}_u$ of the B-series used here with that of the standard B-series, see [19], Chap. III, one sees that our elementary coefficients are the counterpart of the product *elementary weight* $\times \tau^m$, where m is the number of vertices in the tree.

approach in numerical differential equations: with B-series there is a clear separation between, on the one hand, quantities depending on the particular differential equation being integrated but not on the specific integration method⁴ and, on the other hand, method-dependent quantities independent of the differential equation. The second observation is that the elementary coefficients only depend on the second component (wave number) of the label (j, k) .

3.4 Computing the Elementary Coefficients

It does not seem possible to obtain a closed-form expression for the coefficients α_u recursively defined in (15).⁵ Nevertheless the use of this recursion may be systematized by considering the functions

$$\phi_{r,k}(\tau) = \tau^r \exp(ik\tau), \quad r = 0, 1, 2, \dots, k = 0, \pm 1, \pm 2, \dots$$

and the linear space Φ spanned by them, i.e. the set of all complex-valued functions ψ of a real variable that may be written as linear combinations

$$\psi(\tau) = \sum_{r,k} a_{r,k} \phi_{r,k}(\tau),$$

where the $a_{r,k}$'s are complex constants.

Clearly, for each $\psi \in \Phi$, the integral

$$\mathcal{I}(\psi)(\tau) = \int_0^\tau \psi(\tau') d\tau'$$

will be given by the corresponding linear combination of the functions $\kappa_{r,k} = \mathcal{I}(\phi_{r,k})$. These are in turn members of the space Φ , because, obviously

$$\kappa_{r,0}(\tau) = \frac{\tau^{r+1}}{r+1}, \tag{16}$$

and, for $k \neq 0$, integration by parts leads to the recursion

$$\kappa_{r,k} = \frac{ir}{k} \kappa_{r-1,k} - \frac{i}{k} \phi_{r,k}, \quad r = 1, 2, \dots \tag{17}$$

For future reference we list here the first few $\kappa_{r,k}$ with $k \neq 0$

$$\kappa_{0,k} = \frac{i}{k} - \frac{i}{k} \exp(ik\tau), \tag{18}$$

$$\kappa_{1,k} = -\frac{1}{k^2} + \frac{1}{k^2} \exp(ik\tau) - \frac{i}{k} \tau \exp(ik\tau), \tag{19}$$

$$\kappa_{2,k} = -\frac{2i}{k^3} + \frac{2i}{k^3} \exp(ik\tau) + \frac{2}{k^2} \tau \exp(ik\tau) - \frac{i}{k} \tau^2 \exp(ik\tau). \tag{20}$$

⁴In the theory of B-series the exact solution counts as a particular instance of numerical solution.

⁵The structure of the elementary coefficients will be analyzed further in part II of the present work.

Since

$$\phi_{r,k}\phi_{r',k'} = \phi_{r+r',k+k'}, \tag{21}$$

the space Φ is also closed under multiplication (i.e. $\psi_1\psi_2 \in \Phi$, if $\psi_1 \in \Phi$ and $\psi_2 \in \Phi$) and, by induction, we conclude from (15) that, for each $u \in \mathcal{U}$, the elementary coefficient α_u belongs to Φ and accordingly there exist complex constants $a_{u,r,k}$ such that

$$\alpha_u(\tau) = \sum_{r,k} a_{u,r,k}\phi_{r,k}(\tau) = \sum_{r,k} a_{u,r,k}\tau^r \exp(ik\tau). \tag{22}$$

With the help of (16), (17) and (21), it is easy to implement a recursive procedure to compute the elementary coefficients α_u . Once the constants $a_{u,r,k}$ corresponding to the trees in the right-hand side of (15) are known, we write, with the help of (21), the integrand as a linear combination of the $\phi_{r,k}$,

$$\exp(ik\tau')\alpha_{u_1}(\tau') \cdots \alpha_{u_n}(\tau') = \sum_{r,k} b_{u,r,k}\phi_{r,k}(\tau),$$

and obtain

$$\alpha_u(\tau) = \sum_{r,k} b_{u,r,k}\kappa_{r,k}(\tau). \tag{23}$$

Then the constants $a_{u,r,k}$ are found by expressing the functions $\kappa_{r,k}$ as linear combinations of the $\phi_{r,k}$ (cf. (18)–(20)).

After computing the elementary coefficients of all trees of weight $j = 1, 2$, we find the first terms in (14):

$$\begin{aligned} y(\tau) = & y_0 + \epsilon \sum_k \kappa_{0,k}(\tau) f_{1,k}(y_0) + \epsilon^2 \sum_k \kappa_{0,k}(\tau) f_{2,k}(y_0) \\ & + \epsilon^2 \sum_{k,\ell} c_{k,\ell}(\tau) f'_{1,k}(y_0) f_{1,\ell}(y_0) + \mathcal{O}(\epsilon^3), \end{aligned} \tag{24}$$

where

$$c_{k,\ell}(\tau) = \frac{i}{\ell}\kappa_{0,k}(\tau) - \frac{i}{\ell}\kappa_{0,k+\ell}(\tau), \quad \ell \neq 0,$$

and

$$c_{k,0}(\tau) = \kappa_{1,k}(\tau),$$

so that, from (16) and (18)–(19),

$$\begin{aligned} c_{0,0}(\tau) &= \frac{\tau^2}{2}, \\ c_{0,\ell}(\tau) &= \frac{1}{\ell^2} - \frac{1}{\ell^2} \exp(i\ell\tau) + \frac{i}{\ell}\tau, \quad \ell \neq 0, \end{aligned}$$

$$\begin{aligned}
 c_{k,0}(\tau) &= -\frac{1}{k^2} + \frac{1}{k^2} \exp(ik\tau) - \frac{i}{k} \tau \exp(ik\tau), \quad k \neq 0, \\
 c_{k,-k}(\tau) &= \frac{1}{k^2} - \frac{1}{k^2} \exp(ik\tau) + \frac{i}{k} \tau, \quad k \neq 0, \\
 c_{k,\ell}(\tau) &= -\frac{1}{k(k+\ell)} + \frac{1}{k\ell} \exp(ik\tau) \\
 &\quad - \frac{1}{\ell(k+\ell)} \exp(i(k+\ell)\tau), \quad k, \ell, k+\ell \neq 0.
 \end{aligned}$$

3.5 The Expansion of the Averaged Solution

In view of (22), the functions (polynomials)

$$\bar{\alpha}_u(\tau) = \sum_{r,k} a_{u,r,k} \tau^r \tag{25}$$

interpolate the values of $\alpha_u(\tau)$ when τ is an integer multiple of 2π . Therefore, the formula

$$Y(\tau) = y_0 + \sum_{j=1} \epsilon^j \sum_{|u|=j} \frac{\bar{\alpha}_u(\tau)}{\sigma_u} \mathcal{F}_u(y_0), \tag{26}$$

provides the B-series expansion of a function $Y(\tau)$ that interpolates the solution $y(\tau)$ with expansion (14) and is *smooth*, in the sense that all derivatives $(d^k/dt^k)Y$ with respect to the original time t remain bounded as $\epsilon \rightarrow 0$ (cf. the bounded derivative principle [22]).

The leading terms of the B-series for y computed in (24) thus yield

$$\begin{aligned}
 Y(\tau) &= y_0 + \epsilon \tau f_{1,0}(y_0) + \epsilon^2 \frac{\tau^2}{2} f'_{1,0}(y_0) f_{1,0}(y_0) + \epsilon^2 \tau f_{2,0}(y_0) \\
 &\quad + \epsilon^2 \tau \sum_{k \neq 0} \frac{i}{k} (f'_{1,0}(y_0) f_{1,k}(y_0) - f'_{1,k}(y_0) f_{1,0}(y_0)) \\
 &\quad + \epsilon^2 \tau \sum_{k \neq 0} \frac{i}{k} f'_{1,k}(y_0) f_{1,-k}(y_0) + \mathcal{O}(\epsilon^3). \tag{27}
 \end{aligned}$$

3.6 The Averaged Differential Equation

As explained in Sect. 2, from the series (26) for the smooth interpolant $Y(\tau)$, we derive, for each integer $J \geq 1$, a differential equation satisfied by Y up to an $\mathcal{O}(\epsilon^{J+1})$ remainder. (This corresponds to an $\mathcal{O}(\epsilon^J)$ remainder if the differential equation is written in terms of the original independent variable t .)

Let us first consider the cases $J = 1, 2$. Differentiation with respect to τ in (27) yields

$$\left. \frac{d}{d\tau} Y \right|_{\tau=0} = \epsilon f_{1,0}(y_0) + \epsilon^2 f_{2,0}(y_0)$$

$$\begin{aligned}
& + \epsilon^2 \sum_{k \neq 0} \frac{i}{k} (f'_{1,0}(y_0) f_{1,k}(y_0) - f'_{1,k}(y_0) f_{1,0}(y_0)) \\
& + \epsilon^2 \sum_{k \neq 0} \frac{i}{k} f'_{1,k}(y_0) f_{1,-k}(y_0) + \mathcal{O}(\epsilon^3)
\end{aligned}$$

and, since $Y(0) = y_0$, we conclude that the lowest-order ($J = 1$) averaged equation is

$$\frac{d}{d\tau} Y = \epsilon f_{1,0}(Y)$$

(a result certainly expected) and that at the next, $J = 2$, order

$$\begin{aligned}
\frac{d}{d\tau} Y & = \epsilon f_{1,0}(Y) + \epsilon^2 f_{2,0}(Y) \\
& + \epsilon^2 \sum_{k \neq 0} \frac{i}{k} (f'_{1,0}(Y) f_{1,k}(Y) - f'_{1,k}(Y) f_{1,0}(Y)) \\
& + \epsilon^2 \sum_{k \neq 0} \frac{i}{k} f'_{1,k}(Y) f_{1,-k}(Y). \tag{28}
\end{aligned}$$

Averaged equations of higher orders may in principle be found following the same methodology used for the cases $J = 1, 2$. The only difficulty stems from the need to determine the elementary coefficients α_u corresponding to trees of weights $|u| = 3, 4, \dots$ that are required to compute, via (25), the coefficients $\bar{\alpha}_u$ in the expansion of Y . To conclude this subsection we find explicitly the averaged equation with remainder $\mathcal{O}(\epsilon^4)$. This is probably as far as one can reasonably go for the *general* format (12), but, of course, in particular cases it is possible to explicitly construct averaged equations with remainders of higher order.

Our task may be simplified by observing that, since we aim at computing the terms with weight $j = 1, 2, 3$ in the expansion of $dY/d\tau$ at $\tau = 0$, it is not necessary to gain full knowledge of the polynomials $\bar{\alpha}_u$, $|u| = 1, 2, 3$, but only of the coefficient of the first power τ^1 in such polynomials. Now each α_u , $|u| = 1, 2, 3$, is of the form (23), with the summation index r restricted to the values 0, 1, 2 and the formulas (16) and (18)–(20) reveal that in the corresponding functions $\kappa_{r,k}(\tau)$, $r = 0, 1, 2$, $k = 0, \pm 1, \pm 2, \dots$, the power τ^1 is only present in $\kappa_{0,0}$ and $\kappa_{1,k}$, $\kappa_{2,k}$, $k \neq 0$. This observation lowers the number of coefficients $b_{u,r,k}$ to be determined.

In Fig. 1 we have depicted the eight ‘families’ of trees u with weight ≤ 3 (each family has infinitely many trees corresponding to different values of the wave numbers). Accordingly, we write the averaged equation with $J = 3$ in the form

$$\frac{d}{d\tau} Y = \sum_{\mu=1}^{\text{VIII}} F_{\mu}(Y), \tag{29}$$

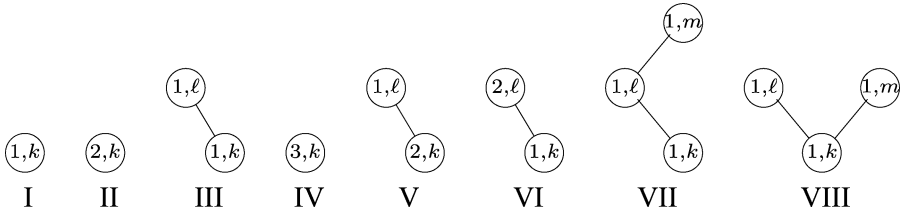


Fig. 1 Families of trees in \mathcal{U} with weight ≤ 3

where F_μ stands for the contribution of the μ th family. Using the methodology outlined above we find that in (29) (all functions evaluated at Y):

$$\begin{aligned}
 F_I &= \epsilon f_{1,0}, & F_{II} &= \epsilon^2 f_{2,0}, & F_{IV} &= \epsilon^3 f_{3,0}, \\
 F_{III} &= \epsilon^2 \left(\sum_{k \neq 0} \frac{i}{k} (f'_{1,0} f_{1,k} - f'_{1,k} f_{1,0}) + \sum_{k \neq 0} \frac{i}{k} f'_{1,k} f_{1,-k} \right), \\
 F_V &= \epsilon^3 \left(\sum_{k \neq 0} \frac{i}{k} (f'_{2,0} f_{1,k} - f'_{2,k} f_{1,0}) + \sum_{k \neq 0} \frac{i}{k} f'_{2,k} f_{1,-k} \right), \\
 F_{VI} &= \epsilon^3 \left(\sum_{k \neq 0} \frac{i}{k} (f'_{1,0} f_{2,k} - f'_{1,k} f_{2,0}) + \sum_{k \neq 0} \frac{i}{k} f'_{1,k} f_{2,-k} \right), \\
 F_{VII} &= \epsilon^3 \left(\sum_{k \neq 0} \frac{1}{k^2} f'_{1,k} f'_{1,0} f_{1,0} - \sum_{\ell \neq 0} \frac{2}{\ell^2} f'_{1,0} f'_{1,\ell} f_{1,0} + \sum_{m \neq 0} \frac{1}{m^2} f'_{1,0} f'_{1,0} f_{1,m} \right. \\
 &\quad + \sum_{\ell \neq 0} \frac{1}{\ell^2} f'_{1,0} f'_{1,-\ell} f_{1,\ell} - \sum_{k \neq 0} \frac{2}{k^2} f'_{1,k} f'_{1,0} f_{1,-k} + \sum_{k \neq 0} \frac{1}{k^2} f'_{1,k} f'_{1,-k} f_{1,0} \\
 &\quad - \sum_{k \neq 0, m \neq 0} \frac{1}{km} (f'_{1,k} f'_{1,-k} f_{1,m} + f'_{1,k} f'_{1,-m} f_{1,m}) \\
 &\quad - \sum_{\substack{k \neq 0, \ell \neq 0 \\ k \neq -\ell}} \frac{1}{(k + \ell)\ell} f'_{1,k} f'_{1,\ell} f_{1,0} - \sum_{\substack{\ell \neq 0, m \neq 0 \\ \ell \neq -m}} \frac{1}{\ell(\ell + m)} f'_{1,0} f'_{1,\ell} f_{1,m} \\
 &\quad \left. + \sum_{\substack{k \neq 0, m \neq 0 \\ k \neq -m}} \frac{1}{km} f'_{1,k} f'_{1,-k-m} f_{1,m} + \sum_{\substack{k \neq 0, m \neq 0 \\ k \neq -m}} \frac{1}{km} f'_{1,k} f'_{1,0} f_{1,m} \right), \\
 F_{VIII} &= \epsilon^3 \left(\sum_{k \neq 0} \frac{1}{k^2} f''_{1,k} [f_{1,0} f_{1,0}] - \sum_{\ell \neq 0} \frac{1}{\ell^2} f''_{1,0} [f_{1,\ell} f_{1,0}] \right. \\
 &\quad \left. - \sum_{\ell \neq 0} \frac{1}{\ell^2} f''_{1,-\ell} [f_{1,\ell}, f_{1,0}] + \sum_{k \neq 0} \frac{1}{k^2} f''_{1,0} [f_{1,\ell}, f_{1,-\ell}] \right)
 \end{aligned}$$

$$\begin{aligned}
& - \sum_{k \neq 0, m \neq 0} \frac{1}{km} f''_{1,k} [f_{1,-k} f_{1,m}] + \sum_{\substack{k \neq 0, \ell \neq 0 \\ k \neq -\ell}} \frac{1}{k(k+\ell)} f''_{1,k} [f_{1,\ell} f_{1,0}] \\
& - \left(\sum_{\substack{\ell \neq 0, m \neq 0 \\ \ell \neq -m}} \frac{1}{2\ell m} f''_{1,-\ell-m} [f_{1,\ell} f_{1,m}] - \sum_{\substack{\ell \neq 0, m \neq 0 \\ \ell \neq -m}} \frac{1}{2\ell m} f''_{1,0} [f_{1,\ell}, f_{1,m}] \right).
\end{aligned}$$

Note that a comparison of the expressions for $F_I - F_{II} - F_{IV}$ or $F_{III} - F_V - F_{VI}$ bears out the independence—mentioned before—of the elementary coefficients of the first index j in the labels (j, k) .

3.7 The Hamiltonian Case

Assume that $H(y, \tau; \epsilon)$ is a Hamiltonian function with d degrees of freedom, 2π -periodic with respect to τ and possessing an expansion (cf. (10))

$$H(y, \tau; \epsilon) = \frac{1}{\epsilon} \sum_{j=1}^{\infty} \epsilon^j H_j(y, \tau)$$

(here y is a $2d$ -dimensional vector). If $f = \mathcal{J}^{-1} \nabla H$, then (1) is a non-autonomous Hamiltonian system whose right-hand side is of the form (10) with $f_j = \mathcal{J}^{-1} \nabla H_j$. Moreover, if (cf. (11))

$$H_j(y, \tau) = \sum_{k=-\infty}^{\infty} \exp(ik\tau) H_{j,k}(y)$$

is the Fourier expansion of H_j , then the Fourier coefficients $f_{j,k}$ in (11) satisfy $f_{j,k} = \mathcal{J}^{-1} \nabla H_{j,k}$.

From Sect. 2 we know that the averaged equations are also Hamiltonian and furthermore, in the present section, we have expressed the averaged vector fields as combinations of elementary differentials. The general theory of symplectic integrators [19, 30] shows that the averaged Hamiltonians will be combinations of the corresponding elementary Hamiltonians. For instance the Hamiltonian for (29) is

$$\begin{aligned}
& \epsilon H_{1,0} + \epsilon^2 H_{2,0} + \epsilon^2 \sum_{k \neq 0} \frac{i}{k} \nabla H_{1,0}^T \mathcal{J}^{-1} \nabla H_{1,k} + \epsilon^2 \sum_{k \neq 0} \frac{i}{2k} \nabla H_{1,k}^T \mathcal{J}^{-1} \nabla H_{1,-k} \\
& + \epsilon^3 H_{3,0} + \epsilon^3 \sum_{k \neq 0} \frac{i}{k} (\nabla H_{2,0}^T \mathcal{J}^{-1} \nabla H_{1,k} + \nabla H_{1,0}^T \mathcal{J}^{-1} \nabla H_{2,k}) \\
& + \epsilon^3 \sum_{k \neq 0} \frac{i}{k} \nabla H_{2,k}^T \mathcal{J}^{-1} \nabla H_{1,-k} \\
& + \epsilon^3 \left(\sum_{k \neq 0} \frac{1}{k^2} \nabla^2 H_{1,k} [\mathcal{J}^{-1} \nabla H_{1,0}, \mathcal{J}^{-1} \nabla H_{1,0}] \right)
\end{aligned}$$

$$\begin{aligned}
 & - \sum_{\ell \neq 0} \frac{1}{\ell^2} \nabla^2 H_{1,0} [\mathcal{J}^{-1} \nabla H_{1,\ell}, \mathcal{J}^{-1} \nabla H_{1,0}] \\
 & - \sum_{\ell \neq 0} \frac{1}{\ell^2} \nabla^2 H_{1,-\ell} [\mathcal{J}^{-1} \nabla H_{1,\ell}, \mathcal{J}^{-1} \nabla H_{1,0}] \\
 & + \sum_{k \neq 0} \frac{1}{k^2} \nabla^2 H_{1,0} [\mathcal{J}^{-1} \nabla H_{1,\ell}, \mathcal{J}^{-1} \nabla H_{1,-\ell}] \\
 & - \sum_{k \neq 0, m \neq 0} \frac{1}{km} \nabla^2 H_{1,k} [\mathcal{J}^{-1} \nabla H_{1,-k}, \mathcal{J}^{-1} \nabla H_{1,m}] \\
 & + \sum_{\substack{k \neq 0, \ell \neq 0 \\ k \neq -\ell}} \frac{1}{k(k+\ell)} \nabla^2 H_{1,k} [\mathcal{J}^{-1} \nabla H_{1,\ell}, \mathcal{J}^{-1} \nabla H_{1,0}] \\
 & - \sum_{\substack{\ell \neq 0, m \neq 0 \\ \ell \neq -m}} \frac{1}{2\ell m} \nabla^2 H_{1,-\ell-m} [\mathcal{J}^{-1} \nabla H_{1,\ell}, \mathcal{J}^{-1} \nabla H_{1,m}] \\
 & - \sum_{\substack{\ell \neq 0, m \neq 0 \\ \ell \neq -m}} \frac{1}{2\ell m} \nabla^2 H_{1,0} [\mathcal{J}^{-1} \nabla H_{1,\ell}, \mathcal{J}^{-1} \nabla H_{1,m}] \Big).
 \end{aligned}$$

3.8 A Class of Autonomous Highly Oscillatory Hamiltonian Systems

We end this section by studying autonomous Hamiltonians of the form

$$\frac{1}{2} p_1^T p_1 + \frac{1}{2} p_2^T p_2 + \frac{1}{2} q_2^T K q_2 + U(q_1, q_2; \epsilon), \tag{30}$$

where p_1, q_1 are d_1 -vectors, p_2, q_2 are d_2 -vectors, U is a real-valued potential that may be expanded in non-negative powers of ϵ and K is a $d_1 \times d_1$ symmetric positive definite matrix with eigenvalues of the form k^2/ϵ^2 (k an integer that may change from eigenvalue to eigenvalue). Observe that when U is independent of q_2 , the system can be decoupled into d_2 harmonic oscillators with large frequencies $k_1/\epsilon, \dots, k_{d_2}/\epsilon$ and a Hamiltonian system with d_1 degrees of freedom and Hamiltonian $(1/2)p_1^T p_1 + U(q_1; \epsilon)$. When $\tau = t/\epsilon$ is used as new independent variable, the Hamiltonian becomes

$$H(p_1, p_2, q_1, q_2; \epsilon) = \frac{\epsilon}{2} p_1^T p_1 + \frac{1}{2} \left(\epsilon p_2^T p_2 + \frac{1}{\epsilon} q_2^T \Omega^2 q_2 \right) + \epsilon U(q_1, q_2; \epsilon), \tag{31}$$

where Ω is a symmetric positive definite matrix with integer eigenvalues, and the equations of motion are then

$$\begin{aligned}
 \frac{d}{d\tau} p_1 &= -\epsilon \nabla_1 U(q_1, q_2; \epsilon), \\
 \frac{d}{d\tau} p_2 &= -\frac{1}{\epsilon} \Omega^2 q_2 - \epsilon \nabla_2 U(q_1, q_2; \epsilon), \\
 \frac{d}{d\tau} q_1 &= \epsilon p_1, \\
 \frac{d}{d\tau} q_2 &= \epsilon p_2.
 \end{aligned} \tag{32}$$

While this system does not fit into the general family (12) considered so far in this section, its solutions $y(\tau) = (p_1(\tau), p_2(\tau), q_1(\tau), q_2(\tau))$ possess at stroboscopic times an expansion of the form (4) and are therefore amenable to our approach. In fact, let us introduce the time-dependent change of variables

$$\begin{aligned}
 \hat{p}_1 &= p_1, & \hat{p}_2 &= \cos(\tau \Omega) p_2 + \epsilon^{-1} \Omega \sin(\tau \Omega) q_2, \\
 \hat{q}_1 &= q_1, & \hat{q}_2 &= -\epsilon \Omega^{-1} \sin(\tau \Omega) p_2 + \cos(\tau \Omega) q_2,
 \end{aligned}$$

that transforms (32) into a non-autonomous Hamiltonian system with Hamiltonian function

$$\begin{aligned}
 \hat{H}(\hat{p}_1, \hat{p}_2, \hat{q}_1, \hat{q}_2, \tau; \epsilon) &= \frac{\epsilon}{2} \hat{p}_1^T \hat{p}_1 + \epsilon U(\hat{q}_1, \cos(\tau \Omega) \hat{q}_2 \\
 &\quad + \epsilon \Omega^{-1} \sin(\tau \Omega) \hat{p}_2; \epsilon).
 \end{aligned} \tag{33}$$

Since the eigenvalues of Ω are positive integers, the change of variables is 2π -periodic in τ , and thus, it reduces to the identity map at stroboscopic times $\tau_n = 2\pi n$. Hence, for any solution $y(\tau) = (p_1(\tau), p_2(\tau), q_1(\tau), q_2(\tau))$ of (32), there exists a solution $\hat{y}(\tau) = (\hat{p}_1(\tau), \hat{p}_2(\tau), \hat{q}_1(\tau), \hat{q}_2(\tau))$ of the system corresponding to (33) such that $y(\tau_n) = \hat{y}(\tau_n)$ for all $\tau_n = 2\pi n$. We conclude that when trying to approximate $y(\tau_n)$ by stroboscopic averaging we may pretend that the solutions being treated are those of the Hamiltonian system with Hamiltonian (33), a system that does fit into the general framework of this section. This implies the existence of an averaged system for (32) that, furthermore, according to Sect. 2, will be Hamiltonian for a Hamiltonian function $\mathcal{H}(Y; \epsilon)$.

From our earlier formulas applied to (33), we find that (capital P 's and Q 's denote the components of the averaged Y) $\mathcal{H}(Y; \epsilon)$ satisfies

$$\mathcal{H}(Y; \epsilon) = \epsilon \left(\frac{1}{2} P_1^T \bar{P}_1 + \frac{1}{2\pi} \int_0^{2\pi} U(Q_1, \cos(\Omega \tau) Q_2; 0) d\tau \right) + \mathcal{O}(\epsilon^2).$$

The proposition in Sect. 2 implies that $\mathcal{H}(y)$ is a conserved quantity of the original system (32) and therefore the function

$$\begin{aligned}
 \epsilon^{-1} (H(y; \epsilon) - \mathcal{H}(y; \epsilon)) &= \frac{1}{2} \left(p_2^T p_2 + \frac{1}{\epsilon^2} q_2^T \Omega^2 q_2 \right) \\
 &\quad - \frac{1}{2\pi} \int_0^{2\pi} (U(q_1, q_2; 0) - U(q_1, \cos(\Omega \tau) q_2; 0)) d\tau \\
 &\quad + \mathcal{O}(\epsilon)
 \end{aligned}$$

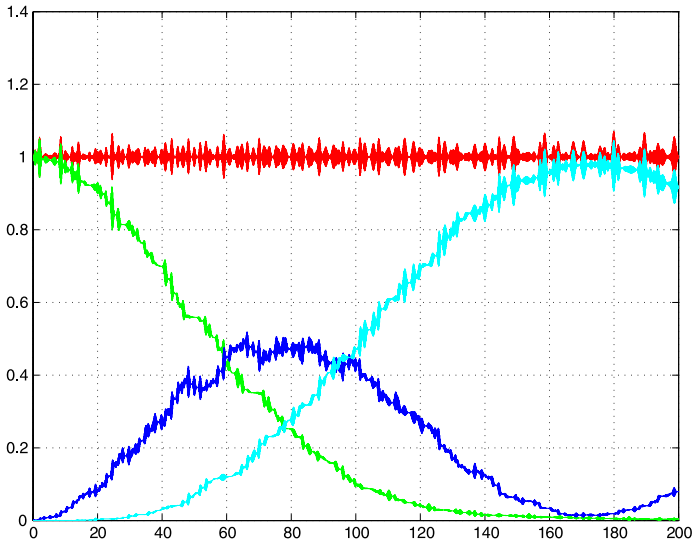


Fig. 2 FPU problem. Oscillatory energies in each of the three stiff springs and total oscillatory energy. Exact solution

will also be conserved. We have established in this way the adiabatic invariance [19] of the oscillatory energy

$$\frac{1}{2} \left(p_2^T p_2 + \frac{1}{\epsilon^2} q_2^T \Omega^2 q_2 \right)$$

of the harmonic oscillators for solutions satisfying $q_2 = \mathcal{O}(\epsilon)$ (in particular for solutions where the total energy (30) remains bounded as $\epsilon \rightarrow 0$).

A well-known example of the family of Hamiltonians considered in this subsection is provided by the variant of the Fermi–Pasta–Ulam problem studied in [19] where q_1 and q_2 are m -dimensional and the potential is given by

$$U(q_1, q_2) = \frac{1}{4} \left((q_{1,1} - q_{2,1})^4 + (-q_{1,m} - q_{2,m})^4 + \sum_{j=1}^{m-1} (q_{1,j+1} - q_{2,j+1} - q_{1,j} - q_{2,j})^4 \right).$$

We have integrated on $0 \leq t \leq 200$ the FPU problem with $m = 3$ (the initial conditions were taken from [19] and we used the implicit midpoint rule with the—very small—time step 0.001). Figure 2 shows the exchange of energy among the three stiff springs, a phenomenon that manifests itself in scales of time $t \sim 1/\epsilon$. Figure 3 shows the results for the corresponding integration of the averaged ($J = 3$) system (29). The averaged system, in spite of not following the oscillations with $\mathcal{O}(\epsilon)$ period in the elongations of the stiff springs, reproduces rather well the exchange between the associated energies.

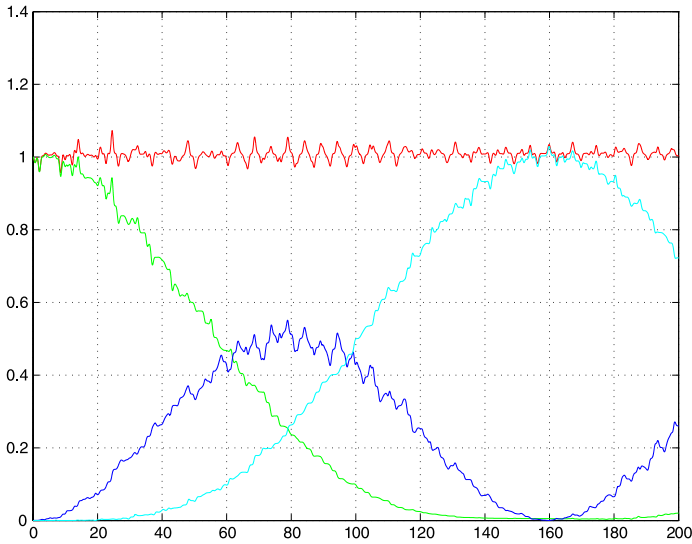


Fig. 3 FPU problem. Oscillatory energies in each of the three stiff springs and total oscillatory energy. Solution of averaged problem

4 Second-Order Differential Equations

4.1 Framework

In this section we study highly oscillatory second-order differential equations

$$\frac{d^2}{dt^2}q = f\left(q, \frac{t}{\epsilon}; \epsilon\right), \tag{34}$$

where $q \in \mathcal{R}^d$ and

$$f(q, \tau; \epsilon) = \frac{1}{\epsilon} \sum_{j=0}^{\infty} \epsilon^j f_j(q, \tau) \tag{35}$$

with

$$f_j(q, \tau) = \sum_{k=-\infty}^{\infty} \exp(ik\tau) f_{j,k}(q).$$

For each j and k , $f_{j,k} \equiv f_{j,-k}^*$. Note that here f is of size $\mathcal{O}(1/\epsilon)$ due to the $j = 0$ term in the sum (35) (cf. (10)).⁶ We assume throughout that $f_{0,0} \equiv 0$: in this way all the leading $\mathcal{O}(1/\epsilon)$ components of the force f in (35) are oscillatory and average out to zero, making it possible for q to undergo variations of size $\mathcal{O}(1)$ (rather

⁶Of course, it would have been possible to include in Sect. 3 first-order differential equations with $\mathcal{O}(1/\epsilon)$ right-hand sides. However, we shall see in this section that the terms with $j = 0$ introduce a number of significant complications that would have hindered the presentation there.

than $\mathcal{O}(1/\epsilon)$ on time-intervals $0 \leq t \leq T$ of length $\mathcal{O}(1)$. The system (34) is supplemented by initial conditions $q(0) = q_0$, $p(0) = p_0$, where $p(t) = dq(t)/dt$ and q_0 , p_0 are given vectors of $\mathcal{O}(1)$ magnitude.

In terms of the scaled time $\tau = t/\epsilon$, we have

$$\frac{d^2}{d\tau^2}q = \epsilon^2 f(q, \tau; \epsilon) = \sum_{j=0}^{\infty} \epsilon^{j+1} f_j(q, \tau); \tag{36}$$

and, initially, $dq(\tau)/d\tau$ has the value ϵp_0 .

4.2 The Expansion of q

The solution $q(\tau)$ of (36) possesses an expansion

$$q(\tau) = q_0 + \sum_{u \in \mathcal{UN}} \epsilon^{|u|} \frac{\alpha_u(\tau)}{\sigma_u} \mathcal{F}_u(p_0, q_0) = q_0 + \sum_{j=1} \epsilon^j \sum_{|u|=j} \frac{\alpha_u(\tau)}{\sigma_u} \mathcal{F}_u(p_0, q_0). \tag{37}$$

Let us briefly point out the most significant differences with (14). Due to the fact that we are now dealing with a second-order differential system (36) where the force $\epsilon^2 f$ does not depend on $dq/d\tau$, we use *special Nyström trees* [18], i.e. trees whose vertices are of two types, meagre and fat, in such a way that the root is meagre, a meagre vertex has at most one son, which is fat, and fat vertices have only meagre sons. Only fat vertices are labeled and the set of labels (j, k) includes the possibility $j = 0$. The notation \mathcal{UN} refers to the set of all special Nyström trees with labeled vertices. The weight $|u|$ of a tree $u \in \mathcal{UN}$ is again the sum of the weights of its vertices; a meagre vertex has unit weight and a fat vertex with label (j, k) has weight j . Note that there exist vertices with zero weight and that, as a consequence, the number of vertices in a tree may be larger than its weight.

The elementary differential \mathcal{F}_u is, for each $u \in \mathcal{UN}$, a mapping from the space $\mathcal{R}^d \times \mathcal{R}^d$ of the variables (p, q) into \mathcal{R}^d constructed in terms of the Fourier coefficients $f_{j,k}$. A fat vertex with label (j, k) and r (meager) sons is associated with the Fréchet derivative of order r of the function $f_{j,k}$ and a terminal meagre son gives rise to a term p .

We rewrite (36) as an integral equation

$$q(\tau) = q_0 + \epsilon p_0 + \epsilon^2 \int_0^\tau (\tau - \tau') f(q(\tau'), \tau'; \epsilon) d\tau',$$

to obtain a procedure for the recursive computation of the elementary coefficients. The recursion starts from the tree consisting only of the (meagre) root, for which $\alpha_u(\tau) = \tau$. For a tree $u \in \mathcal{UN}$ with two or more vertices, let (j, k) be the label of the son of the root and denote by $u_\nu \in \mathcal{UN}$ the trees obtained by removing from u the root and its son; then

$$\alpha_u(\tau) = \int_0^\tau (\tau - \tau') \exp(ik\tau') \prod_\nu \alpha_{u_\nu}(\tau') d\tau'.$$

Let us now turn to the effective computation of the α_n . The role played by the integral operator \mathcal{I} in first-order differential systems is taken here by the operator \mathcal{I}_2 that maps each smooth complex-valued function ψ of a real variable into the function

$$\int_0^\tau (\tau - \tau')\psi(\tau') d\tau'.$$

The elementary coefficients are now linear combinations of the functions (cf. (23)) $\chi_{r,k} = \mathcal{I}_2(\phi_{r,k})$, $r = 0, 1, 2, \dots$, $k = 0, \pm 1, \pm 2, \dots$ (Note that $(d^2/d\tau^2)\chi_{r,k}(\tau) = \phi_{r,k}$.) Clearly

$$\chi_{r,0}(\tau) = \frac{\tau^{r+2}}{(r+2)(r+1)}$$

and furthermore, from

$$\chi_{r,k}(\tau) = \tau\kappa_{r,k}(\tau) - \kappa_{r+1,k}(\tau) \quad k \neq 0,$$

and (18)–(20), we find that, for $k \neq 0$,

$$\chi_{0,k} = \frac{1}{k^2} + \frac{i}{k}\tau - \frac{1}{k^2}\exp(ik\tau), \tag{38}$$

$$\chi_{1,k} = \frac{2i}{k^3} - \frac{1}{k^2}\tau - \frac{2i}{k^3}\exp(ik\tau) - \frac{1}{k^2}\tau\exp(ik\tau), \tag{39}$$

$$\chi_{2,k} = -\frac{6}{k^4} - \frac{2i}{k^3}\tau + \frac{6}{k^4}\exp(ik\tau) - \frac{4i}{k^3}\tau\exp(ik\tau) - \frac{1}{k^2}\tau^2\exp(ik\tau). \tag{40}$$

We have calculated the elementary coefficients of all trees of weight $j = 1, 2$, to find the first terms in (37):

$$\begin{aligned} q(\tau) &= q_0 + \epsilon\tau p_0 + \epsilon \sum_{k \neq 0} \chi_{0,k}(\tau) f_{0,k}(q_0) + \epsilon^2 \chi_{0,0}(\tau) f_{1,0}(q_0) \\ &\quad + \epsilon^2 \sum_{k \neq 0} \chi_{0,k}(\tau) f_{1,k}(q_0) + \epsilon^2 \sum_{k \neq 0} \chi_{1,k}(\tau) f'_{0,k}(q_0) p_0 \\ &\quad + \epsilon^2 \sum_{k \neq 0, \ell \neq 0} c_{k,\ell}(\tau) f'_{0,k}(q_0) f_{0,\ell}(q_0) + \mathcal{O}(\epsilon^3). \end{aligned} \tag{41}$$

Here

$$c_{k,\ell}(\tau) = \frac{1}{\ell^2} \chi_{0,k}(\tau) + \frac{i}{\ell} \chi_{1,k}(\tau) - \frac{1}{\ell^2} \chi_{0,k+\ell}(\tau)$$

or, after replacing the χ 's by their values found above,

$$\begin{aligned} c_{k,\ell}(\tau) &= \frac{1}{k^2\ell^2} - \frac{2}{k^3\ell} - \frac{1}{\ell^2(k+\ell)^2} + i \left(\frac{1}{k\ell^2} - \frac{1}{k^2\ell} - \frac{1}{\ell^2(k+\ell)} \right) \tau \\ &\quad - \left(\frac{1}{k^2\ell^2} - \frac{2}{k^3\ell} - \frac{1}{\ell^2(k+\ell)^2} \right) \exp(ik\tau) - \frac{i}{k^2\ell} \tau \exp(ik\tau), \end{aligned}$$

for $k \neq -\ell$, while

$$c_{k,-k}(\tau) = \frac{3}{k^4} + \frac{2i}{k^3}\tau - \frac{3}{k^4} \exp(ik\tau) + \frac{i}{k^3}\tau \exp(ik\tau) - \frac{1}{2k^2}\tau^2.$$

4.3 The Expansion of the Averaged Solution

Just as in Sect. 3, the expansion of the averaged solution is obtained by interpolating the elementary coefficients. The leading terms of the B-series for q provided in (41) thus yield

$$\begin{aligned} Q(\tau) = & q_0 + \epsilon\tau p_0 + \epsilon \sum_{k \neq 0} \frac{i}{k} \tau f_{0,k}(q_0) + \epsilon^2 \frac{\tau^2}{2} f_{1,0}(q_0) \\ & + \epsilon^2 \sum_{k \neq 0} \frac{i}{k} \tau f_{1,k}(q_0) - \epsilon^2 \sum_{k \neq 0} \frac{2}{k^2} \tau f'_{0,k}(q_0) p_0 \\ & + \epsilon^2 \sum_{k \neq 0, \ell \neq 0} \bar{c}_{k,\ell}(\tau) f'_{0,k}(q_0) f_{0,\ell}(q_0) + \mathcal{O}(\epsilon^3), \end{aligned} \tag{42}$$

where, for $k \neq -\ell$,

$$\bar{c}_{k,\ell}(\tau) = i \left(\frac{1}{k\ell^2} - \frac{2}{k^2\ell} - \frac{1}{\ell^2(k+\ell)} \right) \tau$$

while

$$\bar{c}_{k,-k}(\tau) = \frac{3i}{k^3}\tau - \frac{1}{2k^2}\tau^2.$$

4.4 Averaged Differential Equations

Differentiation with respect to τ of (42) leads to

$$\left. \frac{d^2}{d\tau^2} Q \right|_{\tau=0} = \epsilon^2 \left(f_{1,0}(q_0) - \sum_{k \neq 0} \frac{1}{k^2} f'_{0,k}(q_0) f_{0,-k}(q_0) \right) + \mathcal{O}(\epsilon^3),$$

and therefore the lowest-order ($J = 1$) averaged equation is

$$\frac{d^2}{d\tau^2} Q = \epsilon^2 \left(f_{1,0}(Q) - \sum_{k \neq 0} \frac{1}{k^2} f'_{0,k}(Q) f_{0,-k}(Q) \right) \tag{43}$$

an expression not easily guessed from the mere inspection of (36).

It is important to observe that, from (42),

$$\left. \frac{d}{d\tau} Q \right|_{\tau=0} = \epsilon p_0 + \epsilon \sum_{k \neq 0} \frac{i}{k} f_{0,k}(q_0) + \mathcal{O}(\epsilon^2), \tag{44}$$

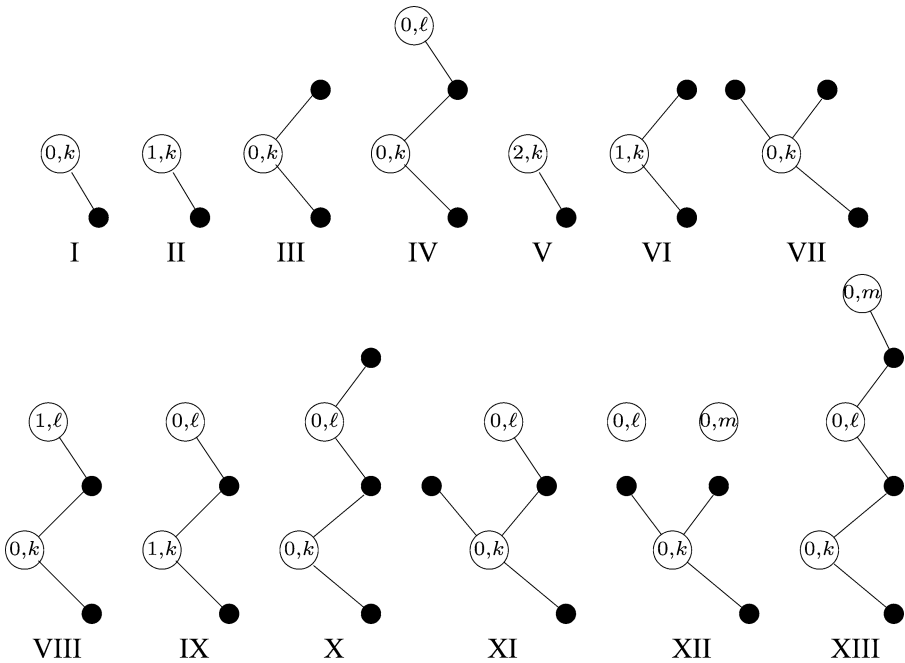


Fig. 4 Families of trees in \mathcal{UN} with weight ≤ 3

so that there is an $\mathcal{O}(\epsilon)$ discrepancy between the initial value ϵp_0 (itself of size $\mathcal{O}(\epsilon)$) of $dq/d\tau$ and the initial value of $dQ/d\tau$. As a consequence, *the (smooth) time derivative $dQ/d\tau$ of the smooth interpolant of q is not an approximation to the smooth interpolant of $p(t) = dq/d\tau$.*

If the force f in the original oscillatory problem (35) is the gradient of a scalar potential, then the averaged (43) possesses also the conservative format $(d^2/d\tau^2)Q = \epsilon^2 \nabla W$, with the potential W given by (cf. [23])

$$W(Q) = V_{1,0}(Q) - \frac{1}{2} \sum_{k \neq 0} \frac{1}{k^2} |\nabla V_{0,k}(Q)|^2.$$

Here $V_{1,0}$ and $V_{0,k}$ are the potentials for $f_{1,0}$ and $f_{0,k}$ respectively, and $|\cdot|$ represents the Euclidean norm.

We conclude this subsection by computing, as in Sect. 3, an averaged equation of higher accuracy. Here we have to face the extra difficulty brought in by the fact that there are low weight trees with a large number of nodes. In Fig. 4 we have listed the thirteen ‘families’ of trees u in \mathcal{UN} with more than one vertex and weight ≤ 3 . Accordingly, we write

$$\frac{d^2}{d\tau^2} Q \Big|_{\tau=0} = \sum_{\mu=I}^{XIII} F_\mu + \mathcal{O}(\epsilon^4) \tag{45}$$

where F_μ stands for the contribution of the μ th family. By using a methodology similar to that employed in Sect. 3 to reduce the number of coefficients $b_{u,r,k}$ to be computed, we find:

$$\begin{aligned}
 F_I &= F_{III} = F_{VI} = 0, \\
 F_{II} &= \epsilon^2 f_{1,0}, \\
 F_{IV} &= -\epsilon^2 \sum_{k \neq 0} \frac{1}{k^2} f'_{0,k} f_{0,-k}, \\
 F_V &= \epsilon^3 f_{2,0}, \\
 F_{VII} &= -\epsilon^3 \sum_{k \neq 0} \frac{1}{k^2} f''_{0,k} [p_0, p_0], \\
 F_{VIII} &= -\epsilon^3 \left(\sum_{k \neq 0} \frac{1}{k^2} f'_{0,k} f_{1,-k} + \sum_{k \neq 0} \frac{1}{k^2} f'_{0,k} f_{1,0} \right), \\
 F_{IX} &= \epsilon^3 \left(-\sum_{k \neq 0} \frac{1}{k^2} f'_{1,k} f_{0,-k} + \sum_{k \neq 0} \frac{1}{k^2} f'_{1,0} f_{0,k} \right), \\
 F_X &= \epsilon^3 \sum_{k \neq 0} \frac{2i}{k^3} f'_{0,k} f'_{0,-k} p_0, \\
 F_{XI} &= -\epsilon^3 \sum_{k \neq 0, \ell \neq 0} \frac{2i}{k^2 \ell} f''_{0,k} [p_0, f_{0,\ell}], \\
 F_{XII} &= \epsilon^3 \left(-\sum_{k \neq 0, \ell \neq 0} \frac{1}{k^2 \ell^2} f''_{0,k} [f_{0,\ell}, f_{0,-k}] \right. \\
 &\quad - \sum_{k \neq 0, \ell \neq 0, m \neq 0} \frac{1}{k^2 \ell m} f''_{0,k} [f_{0,\ell}, f_{0,m}] \\
 &\quad \left. + \sum_{\substack{k \neq 0, \ell \neq 0 \\ k \neq -\ell}} \frac{1}{2k^2 \ell^2 (k + \ell)^2} f''_{0,k} [f_{0,\ell}, f_{0,-(k+\ell)}] \right), \\
 F_{XIII} &= \epsilon^3 \left(-\sum_{k \neq 0, m \neq 0} \frac{2}{k^3 m} f'_{0,k} f'_{0,-k} f_{0,m} \right. \\
 &\quad - \sum_{k \neq 0, m \neq 0} \frac{1}{k^2 m^2} f'_{0,k} f'_{0,-k} f_{0,m} \\
 &\quad + \sum_{k \neq 0, \ell \neq 0} \frac{1}{k^2 \ell^2} f'_{0,k} f'_{0,\ell} f_{0,-\ell} \\
 &\quad \left. + \sum_{\substack{k \neq 0, \ell \neq 0 \\ k \neq -\ell}} \frac{1}{k^2 (k + \ell)^2} f'_{0,k} f'_{0,\ell} f_{0,-(k+\ell)} \right).
 \end{aligned}$$

In (45) it is understood that the functions $f_{j,k}$ and their Fréchet derivatives $f'_{j,k}$, $f''_{j,k}$ that feature in the F_μ 's are evaluated at q_0 .

The averaged differential equation we are seeking is now obtained by discarding the $\mathcal{O}(\epsilon^4)$ remainder in (45) and eliminating q_0 and p_0 with the help of (44), i.e.

$$\frac{d^2}{d\tau^2} Q = \sum_{\mu=1}^{XIII} F_\mu \quad (46)$$

where the functions F_μ are to be evaluated now at Q and

$$P = \omega \frac{d}{d\tau} Q - \sum_{k \neq 0} \frac{i}{k} f_{0,k}(Q)$$

rather than at q_0 and p_0 . Note that the right-hand sides (forces) of the original oscillatory system (36) and the averaged system (43) are independent of the corresponding velocities $dq/d\tau$, $dQ/d\tau$; this is not the case for the system (46).

4.5 An Example: The Kapitza Pendulum

As an illustration of the preceding material, we consider the differential equation

$$\frac{d^2}{dt^2} q = \left(\frac{g}{L} + \frac{1}{\epsilon} \frac{v_{\max}}{L} \cos \left(\frac{t}{\epsilon} + \theta_0 \right) \right) \sin q \quad (47)$$

that describes the motion of a pendulum whose suspension point is subjected to a fast vertical vibration. In (47), q is the angle between the pendulum rod and the upward vertical axis, L the length of the rod, g the acceleration of gravity, $1/\epsilon$ the (large) angular frequency of the vibration of the suspension point, $v_{\max} > 0$ the $\mathcal{O}(1)$ maximum vertical velocity of the suspension point and θ_0 a parameter that governs the initial phase of the vibration. The mechanical system described by (47) is sometimes referred to by the name of Kapitza's pendulum [21] and attracts much interest in physics as an example of the possibility of stabilization by vibration that led to Paul's 1989 Nobel Prize [6, 23, 28, 31].

At the lowest order in ϵ , the averaged equation (43) turns out to be

$$\frac{d^2}{dt^2} Q = \left(\frac{g}{L} - \frac{v_{\max}^2}{2L^2} \cos Q \right) \sin Q; \quad (48)$$

the term with $-v_{\max}^2/(2L^2)$ opposes the gravity term g/L and, for v_{\max} suitable large, turns into a stable equilibrium the usually unstable configuration $Q = 0$, where the pendulum rod is *above* the suspension point.

At the next order in ϵ the averaged equation in formula (46) reads

$$\begin{aligned} \frac{d^2}{dt^2} Q &= \left(\frac{g}{L} - \frac{v_{\max}^2}{2L^2} \cos Q \right) \sin Q \\ &+ \epsilon \left(\frac{v_{\max}}{L} \cos(\theta_0) P^2 \sin Q - \frac{v_{\max}^2}{L^2} \sin(2\theta_0) P \sin^2 Q \right. \\ &\left. - \frac{v_{\max}^3}{4L^3} (\cos(3\theta_0) - 3 \cos(\theta_0)) \sin^3 Q \right), \end{aligned} \tag{49}$$

where

$$P = \frac{d}{dt} Q + \frac{v_{\max}}{L} \sin \theta_0 \sin Q.$$

The $\mathcal{O}(\epsilon)$ terms in (49) stem from the families of trees VII, XI and XII respectively. (The contributions from the families VIII and IX cancel each other and those from families V and XIII vanish.) If $\theta_0 = \pm\pi/2$ (modulo 2π), then (49) reduces to (48), a fact that could have been anticipated by noting that, for those values of θ_0 , the equation (47) is not altered when ϵ is changed into $-\epsilon$.

5 Application to the Analysis of Multiscale Methods

5.1 Heterogeneous Multiscale Methods

For the sake of brevity, we shall only study here the case of second-order problems (34); all our considerations may easily be extended to first-order systems of the form (1).

The numerical integration of the highly oscillatory system (34) may be a very difficult task: standard explicit algorithms suffer from stability restrictions that limit the step-length to size $\mathcal{O}(\epsilon)$ and more sophisticated methods are often hindered by the phenomenon of order reduction [5]. Heterogeneous multiscale methods (HMM) [1, 10–13, 29, 31] (cf. [24]) avoid these difficulties by aiming at finding only the slowly varying components of the solution q without keeping track of the rapidly oscillatory components. A full discussion of the various HMMs is completely out of the scope of this paper and we shall limit ourselves to the asynchronous approach suggested in [7].

From the preceding section we know that the solution q of (34) with initial conditions $q(0) = q_0$, $p(0) = p_0$ differs in terms of size $\mathcal{O}(\epsilon)$ from the solution of the averaged problem

$$\frac{d^2}{dt^2} Q = F(Q), \tag{50}$$

$$F(Q) = f_{1,0}(Q) - \sum_{k \neq 0} \frac{1}{k^2} f'_{0,k}(Q) f_{0,-k}(Q), \tag{51}$$

with initial conditions (see (44))

$$Q(0) = q_0, \quad \left. \frac{d}{dt} Q \right|_{t=0} = p_0 + \sum_{k \neq 0} \frac{i}{k} f_{0,k}(q_0).$$

(Note that, if the $f_{0,k}$ are real, then $(d/dt)Q(0) = p_0$.) HMMs obtain approximations to q by numerically integrating (50) without using the explicit analytic expression given in (51) for the averaged force F . Each time that the numerical integrator (*macro-integrator*) used to integrate (50) requires the value of the function F at a known value Q^* of its argument, $F(Q^*)$ is approximated numerically by an average \hat{F} of values of the force f of the originally given oscillatory problem (34). More precisely:

$$\hat{F}(Q^*) = \frac{2}{\eta} \int_{-\eta/2}^{\eta/2} K\left(\frac{2t}{\eta}\right) f\left(q^*(t), \frac{t}{\epsilon}; \epsilon\right) dt, \quad (52)$$

where η is a (small) scaling parameter, the kernel or weight function K is an even, $K(\xi) = K(-\xi)$, real-valued function of the real variable ξ , $-1 \leq \xi \leq 1$, with unit mass,

$$\int_{-1}^1 K(\xi) d\xi = 1, \quad (53)$$

and $q^*(t)$ is the solution of the original (34) with initial conditions

$$q^*(0) = q_0, \quad \left. \frac{d}{dt} q^* \right|_{t=0} = 0. \quad (54)$$

In practice, the integral in (52) is approximated by a quadrature rule and the required values of q^* are found by numerically integrating (34) with the initial data (54) in the interval $-\eta/2 \leq t \leq \eta/2$. The numerical method used to perform this task is referred to as *micro-integrator*; the step-length δ used by the micro-integrator will typically suffer from a stability/accuracy restriction $\delta = \mathcal{O}(\epsilon)$, but this is acceptable because micro-integrations are only performed over small windows of length η rather than over the whole integration range $0 \leq t \leq T$. (In fact, due to symmetry considerations, it is enough to micro-integrate over the interval $0 \leq t \leq \eta/2$, see [7].)

The error in a HMM consists [29] of the *averaging error* introduced by approximating the solution q by its averaged counterpart Q and the *numerical error* resulting from computing Q numerically. The numerical error may be analyzed by standard techniques after noting that it is nothing but the error of the algorithm used as macro-integrator when applied to (50) with approximate values of the force F . Therefore the numerical error has two origins: (a) the truncation error of the macro-integrator and (b) the use of inexact force values. In turn, the contribution (b) is the result of (b1) the use of the approximation (52) rather than the true averaged force in F in (51), (b2) the use of numerical quadrature to compute the integral in (52), (b3) the use of inexact values of q^* obtained by numerical micro-integration. The analysis of (b2) and (b3) is standard, and accordingly we are only interested in the analysis of (b1), a task that is performed in the subsections that follow.

5.2 The Filtering Error

In view of (34), the formal series expansion for the function $f(q^*(t), t/\epsilon; \epsilon)$ that appears in (52) may be obtained by differentiating twice with respect to t the B-series for $q^*(t)$, which is given in (41) with $q_0 = Q^*$, $p_0 = 0$ (see (54)). After recalling that $(d^2/d\tau^2)\chi_{r,k}(\tau) = \phi_{r,k}(\tau)$, we find in this way

$$\begin{aligned}
 & f\left(q^*(t), \frac{t}{\epsilon}; \epsilon\right) \\
 &= \frac{1}{\epsilon} \sum_{k \neq 0} \phi_{0,k}\left(\frac{t}{\epsilon}\right) f_{0,k}(Q^*) + \phi_{0,0}\left(\frac{t}{\epsilon}\right) f_{1,0}(Q^*) + \sum_{k \neq 0} \phi_{0,k}\left(\frac{t}{\epsilon}\right) f_{1,k}(Q^*) \\
 &+ \sum_{k \neq 0, \ell \neq 0} \frac{1}{\ell^2} \left(\phi_{0,k}\left(\frac{t}{\epsilon}\right) + \frac{i}{\ell} \phi_{1,k}\left(\frac{t}{\epsilon}\right) - \frac{1}{\ell^2} \phi_{0,k+\ell}\left(\frac{t}{\epsilon}\right) \right) f'_{0,k}(Q^*) f_{0,\ell}(Q^*) \\
 &+ R,
 \end{aligned} \tag{55}$$

where R stands for the series consisting of the higher-order terms. We may write $R = R_1(t) + R_2(t)$, where R_1 and R_2 comprise respectively the non-oscillatory and the oscillatory terms. It is easily checked that $R_2 = \mathcal{O}(\epsilon)$, while R_1 is a power series in t^2, t^4, \dots whose coefficients are $\mathcal{O}(1)$.

Next, substitution of (55) into (52) leads to

$$\begin{aligned}
 \hat{F}(Q^*) &= \frac{1}{\epsilon} \sum_{k \neq 0} \hat{\phi}_{0,k} f_{0,k}(Q^*) + \hat{\phi}_{0,0} f_{1,0}(Q^*) + \sum_{k \neq 0} \hat{\phi}_{0,k} f_{1,k}(Q^*) \\
 &+ \sum_{k \neq 0, \ell \neq 0} \frac{1}{\ell^2} \left(\hat{\phi}_{0,k} + \frac{i}{\ell} \hat{\phi}_{1,k} - \frac{1}{\ell^2} \hat{\phi}_{0,k+\ell} \right) f'_{0,k}(Q^*) f_{0,\ell}(Q^*) \\
 &+ \hat{R}_1 + \hat{R}_2,
 \end{aligned} \tag{56}$$

where $\hat{\phi}_{r,k}$ denotes the result of filtering the function $\phi_{r,k}$,

$$\hat{\phi}_{r,k} = \frac{2}{\eta} \int_{-\eta/2}^{\eta/2} K\left(\frac{2t}{\eta}\right) \phi_{r,k}\left(\frac{t}{\epsilon}\right) dt,$$

and \hat{R}_j stands for the result of filtering the remainders R_j , $j = 1, 2$,

$$\hat{R}_j = \frac{2}{\eta} \int_{-\eta/2}^{\eta/2} K\left(\frac{2t}{\eta}\right) R_j(t) dt.$$

Clearly $\hat{R}_1 = \mathcal{O}(\eta^2)$ and $\hat{R}_2 = \mathcal{O}(\epsilon)$. This means that, for accuracy, the value η should be so small that the filtering procedure in (52) retrieves satisfactorily the non-oscillatory components of the micro-force f being filtered.

The normalization (53) implies that $\hat{\phi}_{0,0} = 1$, and then a comparison of (56) with (51) reveals that for an ‘ideal’ filter K for which

$$\hat{\phi}_{r,k} = 0, \quad r = 0, 1, \quad k \neq 0, \tag{57}$$

we have $\hat{F}(Q^*) - F(Q^*) = \hat{R}_1 + \hat{R}_2 = \mathcal{O}(\eta^2 + \epsilon)$.

Reference [13] considers the use of filter functions $K(\xi)$ of class C^ν that vanish with their derivatives at $\xi = \pm 1$ (the exponential filter used in the experiments of [7, 29, 31] satisfies this requirement for any value of $\nu = 1, 2, \dots$). Such filters⁷ only satisfy the target (57) in an approximate manner. In fact, since

$$\hat{\phi}_{r,k} = \left(\frac{\eta}{2\epsilon}\right)^r \int_{-1}^1 K(\xi) \xi^r \exp\left(\frac{ik\eta}{2\epsilon}\xi\right) d\xi,$$

ν integrations by parts reveal that

$$\hat{\phi}_{r,k} = \mathcal{O}\left(k^{-\nu} \left(\frac{\epsilon}{\eta}\right)^{\nu-r}\right), \quad k \neq 0. \quad (58)$$

It follows that

$$\hat{F}(Q^*) - F(Q^*) = \hat{R}_1 + \hat{R}_2 = \mathcal{O}\left(\eta^2 + \epsilon + \frac{1}{\epsilon} \left(\frac{\epsilon}{\eta}\right)^\nu\right)$$

(an expression found in [29] by a different approach). This shows that, for filters in this family, the window length η should be chosen large with respect to the period $2\pi\epsilon$ of the fast oscillations. Since at the same time, as pointed out above, η has to be small with respect to the characteristic time of the smooth, non-oscillatory components of q , it follows that algorithms based on smooth filters cannot operate unless ϵ is much smaller than such a characteristic time, i.e. unless the problem has well-separated time scales.

5.3 Simple Filtering

A simple alternative to the smooth filters we have just discussed was suggested in [7]. It has $K \equiv 1/2$ and $\eta = 2\pi\epsilon$, i.e. it filters through the familiar mean value

$$\hat{F}(Q^*) = \frac{1}{2\pi\epsilon} \int_{-\pi\epsilon}^{\pi\epsilon} f\left(q^*(t), \frac{t}{\epsilon}; \epsilon\right) dt. \quad (59)$$

Since here $\eta = \mathcal{O}(\epsilon)$ while smooth filters require $\eta \gg \epsilon$, the micro-integrations are now much cheaper. However it should be emphasized that this alternative technique exploits the fact that for the problems under consideration the period of the fast oscillations is known beforehand; the scope of the use of smooth filters does not suffer from such a limitation.

As is the case for the smooth filters, the simple filtering formula (59) does not satisfy the conditions (57). In fact, while it is clear that $\hat{\phi}_{0,k} = 0$, $k \neq 0$, the quantities $\hat{\phi}_{1,k}$ are non-zero. However $\hat{\phi}_{1,-k} = -\hat{\phi}_{1,k}$ and therefore under the additional hypothesis $f_{0,k} \equiv f_{0,-k}$ it will be true that, after comparing (51) and (56),

⁷These filters will be referred to as ‘smooth’ filters in view of the fact that, after defining $K(\xi) = 0$ for $|\xi| > 1$, they become C^ν functions in the whole real line.

$$\hat{F}(Q^*) - F(Q^*) = \hat{R}_1 + \hat{R}_2 \text{ or}$$

$$\hat{F}(Q^*) - F(Q^*) = \mathcal{O}(\epsilon).$$

Since $f_{0,k} \equiv f_{0,-k}^*$, the hypothesis $f_{0,k} \equiv f_{0,-k}$ is equivalent to the requirement that the Fourier coefficient $f_{0,k}$ be real or, in other words, that $f_0(q, t/\epsilon)$ be an *even function* of t . When $f_0(q, t/\epsilon)$ is not even, the filtering error $\hat{F}(Q^*) - F(Q^*)$ is $\mathcal{O}(1)$ and simple filtering is inconsistent as already noted in [7].

5.4 Higher-Order Averaged Equations

The HMMs considered so far in the literature have been based on the lowest-order averaged equation (50)–(51) with $\mathcal{O}(\epsilon)$ error at stroboscopic times.⁸ We now illustrate the possibility of devising multiscale numerical methods based on the averaged equation (46), with $\mathcal{O}(\epsilon^2)$ error. We write (46) compactly as

$$\frac{d^2}{dt^2} Q = F_2 \left(Q, \frac{dQ}{dt} \right). \tag{60}$$

(F_2 depends also on ϵ , but this is not reflected in the notation.) To simplify matters we assume that (i) $f(q, t/\epsilon; \epsilon)$ is an even function of t and (ii) $p_0 = 0$. Then the solution Q of (60) with initial condition $Q(0) = q_0$, $(d/dt)Q(0) = 0$ is an $\mathcal{O}(\epsilon^2)$ approximation to the smooth interpolant of the solution q of (34) with $q(0) = q_0$, $p(0) = 0$.⁹ Our aim is then to construct HMM algorithms that integrate (60) in such a way that the filtering errors perpetrated when computing F_2 are also of size $\mathcal{O}(\epsilon^2)$.

When performing the macro-integration of (60) it is necessary to evaluate approximately the force F_2 at known values Q^* , \dot{Q}^* of its arguments, a task that we suggest may be performed with the formula

$$\hat{F}_2(Q^*, \dot{Q}^*) = \frac{2}{\eta} \int_{-\eta/2}^{\eta/2} K \left(\frac{2t}{\eta} \right) f \left(q^*(t), \frac{t}{\epsilon}; \epsilon \right) dt, \tag{61}$$

where, following the ideas in [7], we take $q^*(t)$ to be the solution of (34) with initial conditions

$$q^*(0) = q_0, \quad \left. \frac{d}{dt} q^* \right|_{t=0} = \dot{Q}^*.$$

In the terminology of [7], this is an *asynchronous* approach because the micro-integrations are always performed on the interval $-\eta/2 \leq t \leq \eta/2$, regardless of the current value of t in the macro-integration.

The algorithm just suggested may be analyzed in a way similar to that employed above. Now it is necessary to consider the $\mathcal{O}(\epsilon^3)$ terms in the series (41); the elementary coefficients of those terms are linear combinations of the functions $\chi_{0,k}$,

⁸At non-stroboscopic times $Q - q$ is still $\mathcal{O}(\epsilon)$.

⁹The hypotheses (i) and (ii) are not essential, but without them the presentation is encumbered by the formula that relates $(d/dt)Q(0)$ and p_0 through differentiation of (42), cf. (44).

$\chi_{1,k}$ (as are the $\mathcal{O}(\epsilon^2)$ terms) and $\chi_{2,k}$. A function $\chi_{2,k}(t/\epsilon)$, $k \neq 0$ in the series for $q(t)$ gives rise to a function $\epsilon^{-2}\phi_{2,k}(t/\epsilon)$ in the series for $f = (d^2/dt^2)q$ and, after filtering, to a coefficient $\epsilon^{-2}\hat{\phi}_{2,k}$ in the series for \hat{F}_2 . On the other hand, the contribution of $\chi_{2,k}(t/\epsilon)$ to the exact averaged force is obtained, as discussed in Sect. 4, by first replacing in (40) the factor $\exp(ikt/\epsilon)$ by 1 and then computing the second derivative with respect to t at $t = 0$. It is then obvious that such a contribution equals $-\epsilon^{-2}(2/k^2)$. Therefore, the conditions (57) for an ‘ideal’ filter have now to be supplemented by the additional requirement

$$\hat{\phi}_{2,k} = -\frac{2}{k^2}, \quad k \neq 0. \quad (62)$$

The bound in (58) suggests that it is not possible for the smooth filters considered above to satisfy (62). For the simple filter $K \equiv 1/2$, a trite computation yields

$$\hat{\phi}_{2,k} = \frac{2}{k^2} \exp(ik\pi), \quad k \neq 0,$$

and (62) holds only for *odd* k . Hence, the simpler filter will not achieve the $\mathcal{O}(\epsilon^2)$ error we aim at, unless $f_{0,k}$ vanishes for even k , i.e. unless the leading component $f_0(q, t/\epsilon)$ in the force (35) consists only of Fourier modes $\exp(ikt/\epsilon)f_{0,k}(q)$ for which the filtering interval $-\pi\epsilon \leq t \leq \pi\epsilon$ comprises an odd number of periods, a conclusion not easily guessed in advance!

We have integrated numerically on the interval $0 \leq t \leq \pi/4$ the inverted pendulum problem (47) for which the leading component f_0 of the force f consists only of the odd wave numbers $k = \pm 1$. The values of the parameters were $g = 9.8$, $L = 0.2$, $v_{\max} = 4$, $\theta_0 = 0$, $q(0) = 0.5$, $p(0) = 0.5$, with ϵ ranging from 10^{-2} to 10^{-5} . The ‘classical’ fourth-order Runge–Kutta method¹⁰ was used both as a macro- and micro-integrator and the macro-force was taken from (61) with the simple, constant filter function $K \equiv 1/2$. The macro-step-size Δ was taken from the sequence $(\pi/20)/2^\nu$, $\nu = 0, 1, 2, \dots$ and the micro-step was determined from the formula $\delta = (\pi\epsilon/4)/2^\nu$ (see [7]). Since each micro-integration covers the interval $-\pi\epsilon \leq t \leq \pi\epsilon$, the number of micro-steps in a single micro-integration equals $8 \times 2^\nu$ (independently of ϵ). There are four evaluations of \hat{F}_2 /micro-integrations per macro-step and therefore the total number of micro-steps in $0 \leq t \leq \pi/4$ equals $4 \times 5 \times 8 \times 4^\nu$. The simulations were performed only for those combinations of Δ and ϵ for which the macro-step-length Δ exceeds the vibration period $2\pi\epsilon$; for $\Delta \leq 2\pi\epsilon$, HMMs do not make much sense, and it is better to use a conventional numerical method.

In Table 1 we give the maximum over $0 \leq t \leq \pi/4$ of the absolute value of the difference between the computed Q and the exact solution of (60). Consider first the right-most column of the table ($\epsilon = 10^{-5}$). For the coarser values of Δ , the error behaves like $\mathcal{O}(\Delta^4)$ as corresponds to the fourth-order method used to macro/micro-integrate. For $\Delta = \pi/160$ the algorithm yields errors of roughly one part in ten thousand: a remarkable achievement since the macro-step-length Δ is more than 300

¹⁰More precisely, we used the Runge–Kutta–Nyström method for second-order differential equations implied by the classical Runge–Kutta formula, see [18], Chap. II.14.

Table 1 Errors in Q for the inverted pendulum, simple filtering

Δ	Micro-steps	$\epsilon = 10^{-2}$	$\epsilon = 10^{-3}$	$\epsilon = 10^{-4}$	$\epsilon = 10^{-5}$
$\pi/20$	160	1.94(-1)	2.39(-1)	2.46(-1)	2.47(-1)
$\pi/40$	640	3.24(-2)	1.74(-2)	1.75(-2)	1.75(-2)
$\pi/80$	2,560	***	1.28(-3)	1.02(-3)	1.01(-3)
$\pi/160$	10,240	***	3.15(-4)	7.17(-5)	6.70(-5)
$\pi/320$	40,960	***	2.34(-4)	7.31(-6)	4.47(-6)
$\pi/640$	163,840	***	***	2.67(-6)	3.22(-7)
$\pi/1280$	655,360	***	***	2.26(-6)	4.48(-8)
$\pi/2560$	2,621,440	***	***	2.21(-6)	2.45(-8)

Table 2 Errors in Q for the inverted pendulum, simple filtering on a window of length $4\pi\epsilon$ (two vibrational periods)

Δ	Micro-steps	$\epsilon = 10^{-2}$	$\epsilon = 10^{-3}$	$\epsilon = 10^{-4}$	$\epsilon = 10^{-5}$
$\pi/20$	2×160	4.92(-1)	2.52(-1)	2.47(-1)	2.47(-1)
$\pi/40$	2×640	5.22(-1)	6.76(-2)	1.35(-2)	1.71(-2)
$\pi/80$	$2 \times 2,560$	***	7.57(-2)	7.01(-3)	5.28(-4)
$\pi/160$	$2 \times 10,240$	***	7.65(-2)	7.79(-3)	7.23(-4)
$\pi/320$	$2 \times 40,960$	***	7.69(-2)	7.84(-3)	7.82(-4)

times larger than the vibrational period $2\pi \times 10^{-5}$. However, towards the bottom of the column the error saturates: there the discretization error in the macro/micro-integration is dominated by the filtering error. A comparison of the saturated errors at the bottom of the columns bears out neatly the $\mathcal{O}(\epsilon^2)$ behavior of the filtering error predicted by our analysis. Note also that in the top rows of the table, where the filtering error is negligible, the errors are independent of ϵ .

We integrated again the *same* problem with the *same* parameter values, but now regarding the function $\cos(t/\epsilon)$ in (47) as having period $4\pi\epsilon$, so that micro-integrations are performed in the interval $-2\pi\epsilon \leq t \leq 2\pi\epsilon$. With respect to this new period the wave numbers involved are $k = \pm 2$ and (62) is not satisfied. The numerical results are given in Table 2: now the saturated errors at the bottom of the table are $\mathcal{O}(\epsilon)$ in agreement with our analysis. A comparison between both tables shows that the violation of the condition (62) results in a degraded overall performance of the algorithm.

Finally, we integrated once more the same problem but now regarding $\cos(t/\epsilon)$ in (47) as having period $6\pi/\epsilon$, so that the relevant wave numbers are ± 3 and (62) holds once more. The results of the experiment (not included here) show that the filtering error is restored to being $\mathcal{O}(\epsilon^2)$. This dissipates any possible suspicion that the degraded performance in Table 2 was due to the wider micro-integration window, rather than to the fact that averaging was performed over an even number of cycles of the pendulum vibration.

5.5 Discussion

While, as pointed out before, the material in this section is not meant to discuss exhaustively HMMs, we hope it has illustrated two points:

1. The B-series/modified equation approach suggested in Sects. 2–4 may be advantageously used to analyze the behavior of the error in a variety of HMMs.
2. It is possible to devise multiscale methods that attain, for small macro-step-lengths, errors of size $\mathcal{O}(\epsilon^2)$ rather than merely $\mathcal{O}(\epsilon)$. In fact, we have suggested above an asynchronous algorithm (based on the simple filter function) that achieves $\mathcal{O}(\epsilon^2)$ errors under the hypotheses that the force f is an even function of t and consists only of Fourier modes with odd wave numbers. Presented in [8] (by M.P. Calvo and the present authors) is an algorithm based on finite-difference techniques that results in $\mathcal{O}(\epsilon^\nu)$ errors $\nu = 1, 2, \dots$ in a wide range of oscillatory problems with a single high frequency.

Acknowledgements This research has been supported by ‘Acción Integrada entre España y Francia’ HF2008–0105. A. Murua is also supported by projects MTM2007–61572 (Ministerio de Educación, España) and EHU08/43 (Universidad del País Vasco/Euskal Herriko Unibertsitatea). J.M. Sanz-Serna is also supported by project MTM2007–663257 (Ministerio de Educación, España).

References

1. G. Ariel, B. Engquist, R. Tsai, A multiscale method for highly oscillatory ordinary differential equations with resonance, *Math. Comput.* **78**, 929–956 (2009).
2. V.I. Arnold, *Geometrical Methods in the Theory of Ordinary Differential Equations*, 2nd edn. (Springer, New York, 1988).
3. V.I. Arnold, *Mathematical Methods of Classical Mechanics*, 2nd edn. (Springer, New York, 1989).
4. M.P. Calvo, J.M. Sanz-Serna, Canonical B-series, *Numer. Math.* **67**, 161–175 (1994).
5. M.P. Calvo, J.M. Sanz-Serna, Instabilities and inaccuracies in the integration of highly oscillatory problems, *SIAM J. Sci. Comput.* **31**, 1653–1677 (2009).
6. M.P. Calvo, J.M. Sanz-Serna, Carrying an inverted pendulum on a bumpy road, DCDS B (to appear).
7. M.P. Calvo, J.M. Sanz-Serna, Heterogeneous multiscale methods for mechanical systems with vibrations (submitted).
8. M.P. Calvo, Ph. Chartier, A. Murua, J.M. Sanz-Serna, A stroboscopic numerical method for highly oscillatory problems (submitted).
9. P. Chartier, E. Hairer, G. Vilmart, Algebraic structures of B-series, *Found. Comput. Math.* (to appear).
10. W. E, Analysis of the heterogeneous multiscale method for ordinary differential equations, *Commun. Math. Sci.* **1**, 423–436 (2003).
11. W. E, B. Engquist, The heterogeneous multiscale methods, *Commun. Math. Sci.* **1**, 87–132 (2003).
12. W. E, B. Engquist, X. Li, W. Ren, E. Vanden-Eijnden, Heterogeneous multiscale methods: a review, *Commun. Comput. Phys.* **2**, 367–450 (2007).
13. B. Engquist, R. Tsai, Heterogeneous multiscale methods for stiff ordinary differential equations, *Math. Comput.* **74**, 1707–1742 (2005).
14. R.P. Fedorenko, Derivation and justification of equations in slow time (the stroboscopic method), *Comput. Math. Math. Phys.* **14**, 81–118 (1974).
15. E. Hairer, Backward error analysis of numerical integrators and symplectic methods, *Ann. Numer. Math.* **1**, 107–132 (1994).
16. E. Hairer, Ch. Lubich, Long-time energy conservation of numerical methods for oscillatory differential equations, *SIAM J. Numer. Anal.* **38**, 414–441 (2000).
17. E. Hairer, G. Wanner, On the Butcher group and general multi-value methods, *Computing* **13**, 287–303 (1974).

18. E. Hairer, S.P. Nørsett, G. Wanner, *Solving Ordinary Differential Equations I, Nonstiff Problems*, 2nd edn. (Springer, Berlin, 1993).
19. E. Hairer, Ch. Lubich, G. Wanner, *Geometric Numerical Integration*, 2nd edn. (Springer, Berlin, 2006).
20. A. Iserles, On the global error of discretization methods for highly-oscillatory ordinary differential equations, *BIT Numer. Math.* **42**, 561–599 (2002).
21. P.L. Kapitsa, Dynamical stability of a pendulum when its point of suspension vibrates, in *Collected Papers*, vol. II, ed. by P.L. Kapitsa (Pergamon, Oxford, 1965), pp. 714–725.
22. H.-O. Kreiss, J. Lorenz, Manifolds of slow solutions for highly oscillatory problems, *Indiana Univ. Math. J.* **42**, 1169–1191 (1993).
23. M. Levi, Geometry and physics of averaging with applications, *Physica D* **132**, 150–164 (1999).
24. J. Li, P.G. Kevrekidis, C.W. Gear, I.G. Kevrekidis, Deciding the nature of the coarse equation through microscopic simulations: the baby-bathwater scheme, *SIAM Rev.* **49**, 469–487 (2007).
25. A. Murua, Formal series and numerical integrators. Part I: systems of ODEs and symplectic integrators, *Appl. Numer. Math.* **29**, 221–251 (1999).
26. L.M. Perko, Higher order averaging and related methods for perturbed periodic and quasi-periodic systems, *SIAM J. Appl. Math.* **17**, 698–724 (1968).
27. J.A. Sanders, F. Verhulst, J. Murdock, *Averaging Methods in Nonlinear Dynamical Systems*, 2nd edn. (Springer, New York, 2007).
28. J.M. Sanz-Serna, Stabilizing with a hammer, *Stoch. Dyn.* **8**, 45–57 (2008).
29. J.M. Sanz-Serna, Modulated Fourier expansions and heterogeneous multiscale methods, *IMA J. Numer. Anal.* **29**, 595–605 (2009).
30. J.M. Sanz-Serna, M.P. Calvo, *Numerical Hamiltonian Problems* (Chapman and Hall, London, 1994).
31. R. Sharp, Y.-H. Tsai, B. Engquist, Multiple time scale numerical methods for the inverted pendulum problem, in *Multiscale Methods in Science and Engineering*, ed. by B. Engquist, P. Lötsdtedt, O. Runborg. Lect. Notes Comput. Sci. Eng., vol. 44 (Springer, Berlin, 2005), pp. 241–261.