

Conservation properties of numerical integrators for highly oscillatory Hamiltonian systems

DAVID COHEN[†]

Département de Mathématiques, Université de Genève, CH-1211 Genève 4, Switzerland

[Received on 17 January 2005; revised on 3 March 2005]

Modulated Fourier expansion is used to show long-time near-conservation of the total and oscillatory energies of numerical methods for Hamiltonian systems with highly oscillatory solutions. The numerical methods considered are an extension of the trigonometric methods. A brief discussion of conservation properties in the continuous problem and in the multi-frequency case is also given.

Keywords: trigonometric methods; Hamiltonian systems; modulated Fourier expansion; energy conservation; oscillatory solutions.

1. Introduction

We consider Hamiltonian systems

$$\begin{aligned}\dot{p} &= -\nabla_q H(p, q), \\ \dot{q} &= \nabla_p H(p, q),\end{aligned}\tag{1.1}$$

with the Hamiltonian function

$$H(p, q) = K(p_1, q) + \frac{1}{2} p_2^T p_2 + \frac{\omega^2}{2} q_2^T q_2,\tag{1.2}$$

where the vectors $p = (p_1, p_2)$ and $q = (q_1, q_2)$ are partitioned according to the partition of the square matrix

$$\Omega = \begin{pmatrix} 0 & 0 \\ 0 & \omega I \end{pmatrix},$$

with blocks of arbitrary dimension and where ω is a large positive parameter.

We assume that the initial values satisfy

$$\frac{1}{2} \|p(0)\|^2 + \frac{1}{2} \|\Omega q(0)\|^2 \leq E,\tag{1.3}$$

where E is independent of ω .

Our attention will particularly focus on the near-conservation of the *oscillatory energy*

$$I(p, q) = \frac{1}{2} (p_2^T p_2 + \omega^2 q_2^T q_2),\tag{1.4}$$

over long time intervals.

[†]Email: David.Cohen@math.unige.ch. Present address: Mathematisches Institut, Universität Tübingen, D-72076 Tübingen, Germany (cohen@na.uni-tuebingen.de).

By taking the function K in (1.2) to be $\frac{1}{2}p_1^T p_1 + U(q)$, we recover the Hamiltonian function considered by Hairer & Lubich (2000) (see also Hairer *et al.*, 2002, Chapter XIII). Our aim, in this article, is to extend the results of Hairer & Lubich (2000) to the more general Hamiltonian function (1.2).

In particular, it is possible to consider coupling between the position q and the momentum p_1 , such as $K(p_1, q) = \frac{1}{2}p_1^T M(q)^{-1} p_1$, where $M(q)$ is a mass matrix. Simple examples described by such a Hamiltonian are the stiff spring pendulum (see Ascher & Reich, 1999a) or the diatomic molecule (see Ascher & Reich, 1999b). More complicated examples can be found in physics, molecular dynamics or astronomy (as we will see below).

EXAMPLE 1.1 As a concrete example, we consider the motion of a planar elastic dumbbell spacecraft acting under a central gravitational field. Such a satellite is composed of two equal masses, m , connected by a stiff spring with a stiffness constant $k \gg 1$. As in Sanyal *et al.* (2003), we place the origin at the centre of the central body, the radial distance from the origin to the satellite is denoted by r and the distance of each mass particle from the centre of mass of the spacecraft is q . We denote by ϕ the angular position of the dumbbell and by θ the attitude angle. This is shown in Fig. 1.

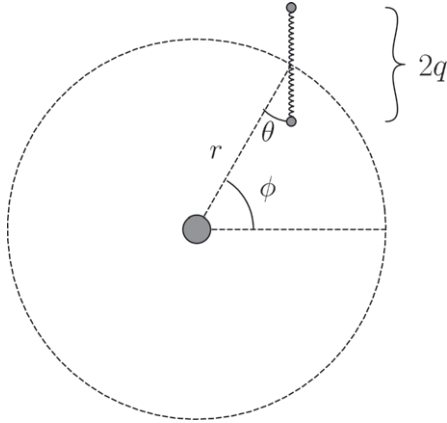


FIG. 1. Planar dumbbell spacecraft.

For this problem, the Lagrangian reads as

$$\begin{aligned} L(\dot{r}, \dot{\phi}, \dot{\theta}, \dot{q}, r, \phi, \theta, q) = & m(\dot{r}^2 + \dot{q}^2 + q^2\dot{\theta}^2 + 2q^2\dot{\theta}\dot{\phi} + (r^2 + q^2)\dot{\phi}^2) \\ & - V_g(r, \theta, q) - 2k(q - l)^2, \end{aligned} \quad (1.5)$$

where l is half the unstretched length of the spring, and

$$V_g(r, \theta, q) = -\frac{\mu m}{r} \left(2 - \frac{q^2}{r^2} (1 - 3 \cos^2(\theta)) \right)$$

is the gravitational potential.

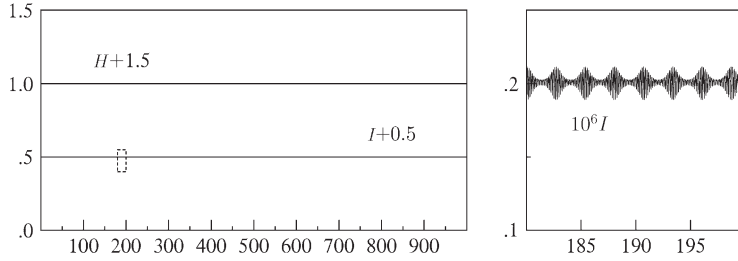


FIG. 2. Scaled total and oscillatory energies for Hamiltonian problem with (1.6), with a zoom of I .

After a change of coordinate (for details see Appendix and Sanyal *et al.*, 2005), we obtain the following Hamiltonian function

$$H(p_\rho, p_\phi, p_\theta, p_\sigma, \rho, \phi, \theta, \sigma) = \frac{1}{2} \left(p_\rho^2 + \frac{1}{\rho^2} (p_\phi - p_\theta)^2 + \frac{1}{(\sigma + \varepsilon)^2} p_\theta^2 + p_\sigma^2 - \frac{2}{\rho} + \frac{(\sigma + \varepsilon)^2}{\rho^3} (1 - 3 \cos^2(\theta)) + \omega^2 \sigma^2 \right), \quad (1.6)$$

where the values for the parameters ε and ω are taken from Sanyal *et al.* (2005) and are given by $\varepsilon = 7.5 \times 10^{-5}$ and $\omega = \sqrt{1800}$. This Hamiltonian function is of the type (1.2) with *slow components* (i.e. q_1) ρ, ϕ, θ and a *fast component* (i.e. q_2) σ .

Let us use a very precise numerical method (namely *DOP853*, for a definition, see Hairer *et al.*, 1993), and plot (see Fig. 2) the different energies involved in this problem for the initial values taken from Sanyal *et al.* (2005): $\rho(0) = 1$, $\theta(0) = \pi/2$, $\sigma(0) = 0.2\varepsilon$, $p_\phi(0) = 0.999958 + (\sigma(0) + \varepsilon)^2$ ($0.07 + 0.999958$), $p_\theta(0) = (\sigma(0) + \varepsilon)^2(0.07 + 0.999958)$ and zero for the remaining ones.

As mentioned above, the oscillatory energy is nearly preserved over long time intervals.

To explain this behaviour, we begin with presenting the *modulated Fourier expansion* of the exact solution (Section 2). Then, we discuss an extension of the numerical methods given in Hairer & Lubich (2000) (Section 3). In Section 4, we apply the approach of the modulated Fourier expansion to the numerical solution and explain its good behaviour. In Section 5, we extend the class of studied problems by adding a small perturbation to the function H of (1.2). In the last section, we briefly discuss the multi-frequency case of the Hamiltonian (1.2).

2. Modulated Fourier expansion of the exact solution

To show the near-conservation of the oscillatory energy for Hamiltonian systems with the Hamiltonian function (1.2), we follow the lines of Hairer *et al.* (2002, Section XIII.5). Here, we state the results omitting the proofs (a detailed version can be found in Cohen, 2004, Chapter 5). In this article, we focus our attention on the proofs for the numerical solution (see Section 4) because the ideas to show the near-conservation of the oscillatory energy for the exact solution are very similar.

To show this geometric property, we assume that the derivatives of the function K in (1.2) are bounded independently of ω . We then give the modulated Fourier expansion of the exact solution.

THEOREM 2.1 If the solution $(p(t), q(t))$ of the Hamiltonian system (1.1) with the Hamiltonian function (1.2) satisfies condition (1.3) and stays in a compact set for $0 \leq t \leq T$, then the solution admits an

expansion of the form

$$\begin{aligned} p(t) &= \sum_{|k| < N} e^{ik\omega t} \eta^k(t) + R_N(t), \\ q(t) &= \sum_{|k| < N} e^{ik\omega t} \zeta^k(t) + S_N(t), \end{aligned} \quad (2.1)$$

for arbitrary $N \geq 2$, where the remainder terms are bounded by

$$R_N(t) = \mathcal{O}(\omega^{-N}), \quad S_N(t) = \mathcal{O}(\omega^{-N}), \quad \text{for } 0 \leq t \leq T.$$

The real functions $\eta = \eta^0 = (\eta_1, \eta_2)$ and $\zeta = \zeta^0 = (\zeta_1, \zeta_2)$ and the complex functions $\eta^k = (\eta_1^k, \eta_2^k)$ and $\zeta^k = (\zeta_1^k, \zeta_2^k)$ are bounded, together with all their derivatives, by

$$\begin{aligned} \zeta_1 &= \mathcal{O}(1), & \eta_1 &= \mathcal{O}(1), & \zeta_2 &= \mathcal{O}(\omega^{-2}), & \eta_2 &= \mathcal{O}(\omega^{-2}), \\ \zeta_1^1 &= \mathcal{O}(\omega^{-2}), & \eta_1^1 &= \mathcal{O}(\omega^{-2}), & \zeta_2^1 &= \mathcal{O}(\omega^{-1}), & \eta_2^1 &= \mathcal{O}(\omega^{-1}), \\ \zeta_1^k &= \mathcal{O}(\omega^{-k-1}), & \eta_1^k &= \mathcal{O}(\omega^{-k-1}), & \zeta_2^k &= \mathcal{O}(\omega^{-k-2}), & \eta_2^k &= \mathcal{O}(\omega^{-k-1}), \end{aligned} \quad (2.2)$$

for $k = 2, \dots, N-1$. Moreover, we have $\eta^{-k} = \overline{\eta^k}$ and $\zeta^{-k} = \overline{\zeta^k}$. These functions are unique up to terms of size $\mathcal{O}(\omega^{-N})$. The constants symbolized by the \mathcal{O} -notation are independent of ω and t with $0 \leq t \leq T$ but depend on N, T and E .

The *modulation functions* η^k and ζ^k have almost-invariants that are related to the total energy H and to the oscillatory energy I . To see this, let us define $\mathbf{p} = (p^{-N+1}, \dots, p^0, \dots, p^{N-1})$ with $p^k = e^{ik\omega t} \eta^k$ (a similar notation is used for \mathbf{q}). To this end, we insert (2.1) into the system (1.1) with the Hamiltonian function (1.2), expand the non-linearity around (p_1^0, q^0) and compare the coefficients of $e^{ik\omega t}$. The modulation functions are then determined to satisfy the following system (for $k = 0, \dots, N-1$)

$$\dot{p}^k + \Omega^2 q^k = -\nabla_{q^{-k}} \mathcal{H}(\mathbf{p}_1, \mathbf{q}) + \mathcal{O}(\omega^{-N}), \quad (2.3)$$

$$\dot{q}^k = \begin{pmatrix} 0 & 0 \\ 0 & I \end{pmatrix} p^k + \nabla_{p^{-k}} \mathcal{H}(\mathbf{p}_1, \mathbf{q}) + \mathcal{O}(\omega^{-N}), \quad (2.4)$$

with

$$\mathcal{H}(\mathbf{p}_1, \mathbf{q}) = K(p_1^0, q^0) + \sum_{s(\alpha)+s(\beta)=0} \frac{1}{m!n!} D_1^m D_2^n K(p_1^0, q^0) (\mathbf{p}_1^\alpha, \mathbf{q}^\beta). \quad (2.5)$$

Here, the sum is over all integers m and n greater than or equal to zero and all multi-indices $\alpha = (\alpha_1, \dots, \alpha_m)$ and $\beta = (\beta_1, \dots, \beta_n)$ with integers $0 < |\alpha_j|, |\beta_j| < N$ which have given sums $s(\alpha)$ and $s(\beta)$, respectively.

Neglecting the $\mathcal{O}(\omega^{-N})$ terms, (2.3)–(2.4) is a Hamiltonian system with

$$\mathcal{H}(\mathbf{p}, \mathbf{q}) = \frac{1}{2} \sum_{|k| < N} \left(q^{-kT} \Omega^2 q^k + p_2^{-kT} p_2^k \right) + \mathcal{H}(\mathbf{p}_1, \mathbf{q}). \quad (2.6)$$

Moreover, this formal invariant is close to the Hamiltonian (1.2): using the bounds of Theorem 2.1 and the fact that $p_2^1 = i\omega q_2^1 + \mathcal{O}(\omega^{-1})$, we see that the dominating terms of (1.2) and (2.6) coincide up to terms of size $\mathcal{O}(\omega^{-1})$.

Besides this formal invariant, system (2.3)–(2.4) has another formal invariant

$$\mathcal{I}(\mathbf{p}, \mathbf{q}) = -i\omega \sum_{0 < |k| < N} k q^{-kT} p^k, \quad (2.7)$$

which turns out to be close to the oscillatory energy (1.4). In fact, like before, the bounds obtained in Theorem 2.1 show that

$$\mathcal{I}(\mathbf{p}, \mathbf{q}) = -i\omega(q_2^{-1})^T p_2^1 + i\omega(q_2^1)^T p_2^{-1} + \mathcal{O}(\omega^{-1}).$$

This implies that (2.7) and (1.4) are equal up to terms of size $\mathcal{O}(\omega^{-1})$.

This permits us to prove the main result of this section, which states that the oscillatory energy (1.4) is nearly conserved over long time intervals.

THEOREM 2.2 If the solution $(p(t), q(t))$ of the Hamiltonian problem (1.1) with the Hamiltonian function (1.2), with initial values satisfying (1.3), stays in a compact set for $0 \leq t \leq \omega^N$, then

$$I(p(t), q(t)) = I(p(0), q(0)) + \mathcal{O}(\omega^{-1}) + \mathcal{O}(t\omega^{-N}).$$

The constants symbolized by \mathcal{O} are independent of ω and t , but depend on E and N .

Benettin *et al.* (1987) studied almost similar Hamiltonian functions and showed, using other techniques, the near-conservation of the oscillatory energy over exponentially long time intervals.

To conclude this section, we want to mention that a finer analysis, similar to the one given in Cohen *et al.* (2003) for the Hamiltonian function $H(p, q) = \frac{1}{2}(p^T p + q^T \Omega^2 q) + U(q)$, should also show the near-conservation of the oscillatory energy over exponentially long time intervals.

3. Numerical methods

In this section, we adapt the trigonometric methods given in Hairer & Lubich (2000) to the case of the Hamiltonian function (1.2). Developing the Hamiltonian system for this Hamiltonian function, we obtain

$$\begin{aligned} \dot{p}_1 &= -\nabla_{q_1} K(p_1, q), \\ \dot{p}_2 &= -\omega^2 q_2 - \nabla_{q_2} K(p_1, q), \\ \dot{q}_1 &= \nabla_{p_1} K(p_1, q), \\ \dot{q}_2 &= p_2. \end{aligned}$$

Treating the second components of p and q with a symmetric trigonometric method and the first components with the Störmer–Verlet method, one gets the following numerical scheme

$$\begin{aligned} p^{n+1/2} &= p^n - \frac{h}{2} \widehat{\Psi} \nabla_q K(p_1^{n+1/2}, \Phi q^n), \\ q_1^{n+1} &= q_1^n + \frac{h}{2} (\nabla_{p_1} K(p_1^{n+1/2}, \Phi q^n) + \nabla_{p_1} K(p_1^{n+1/2}, \Phi q^{n+1})), \\ q_2^{n+1} &= \cos(h\omega) q_2^n + h \operatorname{sinc}(h\omega) p_2^{n+1/2}, \\ \tilde{p}_1^{n+1/2} &= p_1^{n+1/2}, \\ \tilde{p}_2^{n+1/2} &= -\omega \sin(h\omega) q_2^n + \cos(h\omega) p_2^{n+1/2}, \\ p^{n+1} &= \tilde{p}^{n+1/2} - \frac{h}{2} \widehat{\Psi} \nabla_q K(p_1^{n+1/2}, \Phi q^{n+1}), \end{aligned} \tag{3.1}$$

where, here and in the sequel, $\widehat{\Psi} = \widehat{\psi}(h\Omega)$, $\Phi = \phi(h\Omega)$, $\widehat{\Psi}_2 = \widehat{\psi}(h\omega)$ and $\operatorname{sinc}(\zeta) = \sin(\zeta)/\zeta$. The filter functions $\widehat{\psi}$, ϕ are even real-valued functions with $\widehat{\psi}(0) = \phi(0) = 1$.

We remark that the method is explicit if the function $K(p_1, q)$ takes the form $K(p_1, q) = \frac{1}{2} p_1^T M(q_2) p_1 + U(q)$. This was not the case for the dumbbell problem of the first section; however, the special structure of the Hamiltonian (1.6) makes the numerical method explicit for this problem too. We also remark that if $K(p_1, q) = \frac{1}{2} p_1^T p_1 + U(q)$, the numerical method (3.1) reduces to the trigonometric methods analysed in Hairer *et al.* (2002, Chapter XIII). In the next section, we will extend all previous results concerning the trigonometric methods to our numerical method.

As in Hairer *et al.* (2002, Section XIII.2), we can show that if $\widehat{\psi}(\zeta) = \phi(\zeta)$ holds, then method (3.1) is symplectic (for details, see Cohen, 2004, Chapter 5).

EXAMPLE 3.1 Let us return to the dumbbell spacecraft. In Fig. 3, we plot the total energy H and the oscillatory energy I obtained by numerical method (3.1). For the filter functions, we choose $\widehat{\psi}_2(\zeta) = \text{sinc}^2(\zeta/2)/\text{sinc}(\zeta)$ and $\phi_2(\zeta) = \widehat{\psi}_2(\zeta)$. With this choice, the numerical method is symmetric and symplectic.

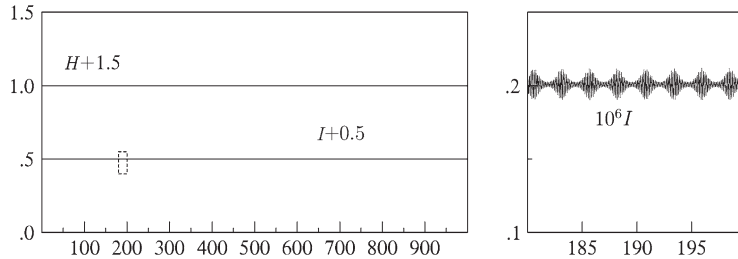


FIG. 3. Energies along the numerical solution of the dumbbell spacecraft problem (see Example 1.1) with step size $h = 0.03$.

We see that this numerical method also approximately conserves H and I . This will be explained in the next section.

4. Modulated Fourier expansion of the numerical solution

In this section, we explain, with the help of the modulated Fourier expansion, the good behaviour of our numerical methods (3.1) applied to the Hamiltonian problems (1.1) with the Hamiltonian function (1.2).

We are interested in the long-time conservation of the total energy H and of the oscillatory energy I along the numerical solution. We make the following assumptions:

- The initial values satisfy

$$\frac{1}{2} \|p^0\|^2 + \frac{1}{2} \|\Omega q^0\|^2 \leq E. \quad (4.1)$$

- The numerical solution stays in a compact set.
- We impose a lower bound on the step size: $h/\omega \geq c_0 > 0$.
- We assume the numerical non-resonance condition:

$$\left| \sin\left(\frac{1}{2}kh\omega\right) \right| \geq c\sqrt{h}, \quad \text{for } k = 1, \dots, N, \quad \text{with } N \geq 2. \quad (4.2)$$

For a given h and ω , this condition imposes a restriction on N . In the sequel, N is a fixed integer such that (4.2) holds.

- For the filter functions, we require the following conditions:

$$\begin{aligned} |\widehat{\psi}(h\omega)| &\leq C_1 \operatorname{sinc}^2\left(\frac{1}{2}h\omega\right), \\ |\widehat{\psi}(h\omega)| &\leq C_2 |\operatorname{sinc}(h\omega)|. \end{aligned} \quad (4.3)$$

- Finally, for

$$\mu(\zeta) = \phi_2(\zeta) \widehat{\psi}_2^{-1}(\zeta), \quad (4.4)$$

we require $\mu(h\omega) \geq c_1 > 0$.

Now, we can state the main result of this section.

THEOREM 4.1 Under the above assumptions, we have, for the numerical solution (3.1),

$$\begin{aligned} H(p^n, q^n) &= H(p^0, q^0) + \mathcal{O}(h), \\ I(p^n, q^n) &= I(p^0, q^0) + \mathcal{O}(h), \end{aligned}$$

for $0 \leq nh \leq h^{-N+1}$.

As in Section 2, we begin to prove that the numerical solution has on a small interval, say $0 \leq t = nh \leq T$, a modulated Fourier expansion.

THEOREM 4.2 Under the assumptions of Theorem 4.1, the numerical solution (p^n, q^n) of (3.1) admits, for $0 \leq t = nh \leq T$, the expansion

$$\begin{aligned} p^n &= \sum_{|k| < N} e^{ik\omega t} \eta_h^k(t) + R_{h,N}(t), \\ q^n &= \sum_{|k| < N} e^{ik\omega t} \zeta_h^k(t) + S_{h,N}(t), \end{aligned} \quad (4.5)$$

where the remainder terms are bounded by

$$R_{h,N}(t) = \mathcal{O}(th^{N-2}), \quad S_{h,N}(t) = \mathcal{O}(th^{N-2}). \quad (4.6)$$

For the modulation functions, we have the following bounds

$$\begin{aligned} \zeta_{h,1} &= \mathcal{O}(1), & \eta_{h,1} &= \mathcal{O}(1), & \zeta_{h,2} &= \mathcal{O}(\omega^{-2}), & \eta_{h,2} &= \mathcal{O}(\omega^{-1}), \\ \zeta_{h,1}^1 &= \mathcal{O}(\omega^{-2}), & \eta_{h,1}^1 &= \mathcal{O}(\omega^{-2}), & \zeta_{h,2}^1 &= \mathcal{O}(\omega^{-1}), & \eta_{h,2}^1 &= \mathcal{O}(\omega^{-1}), \\ \zeta_{h,1}^k &= \mathcal{O}(\omega^{-k-1}), & \eta_{h,1}^k &= \mathcal{O}(\omega^{-k-1}), & \zeta_{h,2}^k &= \mathcal{O}(\omega^{-k-2}), & \eta_{h,2}^k &= \mathcal{O}(\omega^{-k-1}), \end{aligned} \quad (4.7)$$

for $k = 2, \dots, N-1$. Moreover, we have $\eta^{-k} = \overline{\eta^k}$ and $\zeta^{-k} = \overline{\zeta^k}$. The constants symbolized by the \mathcal{O} -notation are independent of ω and h , but depend on E, N, c_0 and T .

To obtain this result, we use ideas very similar to those in the proof of Theorem 5.2 of Hairer *et al.* (2002, Section XIII.5.2). But, in this case, the proof becomes more complicated and more technical.

Proof. We look for two functions of the form

$$\begin{aligned} p_h(t) &= \eta_h(t) + \sum_{0 < |k| < N} e^{ik\omega t} \eta_h^k(t), \\ q_h(t) &= \zeta_h(t) + \sum_{0 < |k| < N} e^{ik\omega t} \zeta_h^k(t), \end{aligned} \quad (4.8)$$

with smooth (in the sense that all their derivatives are bounded independently of h and ω) coefficients ζ_h , ζ_h^k , η_h and η_h^k , which have a small defect when they are inserted into the numerical method (3.1):

$$\begin{aligned} p^n &= p_h(t) + \mathcal{O}(h^{N-1}), \\ q^n &= q_h(t) + \mathcal{O}(h^{N-1}). \end{aligned}$$

Construction of the coefficient functions. To find the coefficient functions η_h^k and ζ_h^k , we insert (4.8) in the numerical method (3.1), expand the non-linearity functions $\nabla_p K$ and $\nabla_q K$ around $(\eta_{h,1}(t), \Phi \zeta_h(t))$ and compare the coefficients of $e^{ik\omega t}$. To motivate the ansatz (4.11) below, we compare the dominant terms appearing when doing these manipulations.

- Looking at the first expression in (3.1), we implicitly define, for $t = nh + \frac{h}{2}$,

$$\widehat{p}_h(t) = p_h\left(t - \frac{h}{2}\right) - \frac{h}{2} \widehat{\Psi} \nabla_q K\left(\widehat{p}_{h,1}(t), \Phi q_h\left(t - \frac{h}{2}\right)\right).$$

As for (4.8), we also define

$$p^{n+1/2} = \widehat{p}_h(t) = \zeta_h(t) + \sum_{0 < |k| < N} e^{ih\omega t} \zeta_h^k(t). \quad (4.9)$$

The coefficient functions of (4.9) satisfy $\zeta_h^k(t) = \eta_h^k(t) + \mathcal{O}(h)$.

- For the second term of the numerical method (3.1), we have

$$\begin{aligned} q_{h,1}\left(t + \frac{h}{2}\right) - q_{h,1}\left(t - \frac{h}{2}\right) &= \frac{h}{2} \left(\nabla_{p_1} K\left(\widehat{p}_{h,1}(t), \Phi q_h\left(t - \frac{h}{2}\right)\right) \right. \\ &\quad \left. + \nabla_{p_1} K\left(\widehat{p}_{h,1}(t), \Phi q_h\left(t + \frac{h}{2}\right)\right) \right). \end{aligned}$$

Using (4.8), we get

$$\begin{aligned} &\sum_{|k| < N} e^{ik\omega(t+h/2)} \zeta_h^k\left(t + \frac{h}{2}\right) - \sum_{|k| < N} e^{ik\omega(t-h/2)} \zeta_h^k\left(t - \frac{h}{2}\right) \\ &= \frac{h}{2} \left(\nabla_{p_1} K\left(\widehat{p}_{h,1}(t), \Phi q_h\left(t - \frac{h}{2}\right)\right) + \nabla_{p_1} K\left(\widehat{p}_{h,1}(t), \Phi q_h\left(t + \frac{h}{2}\right)\right) \right). \end{aligned}$$

Expanding the smooth functions η_h and ζ_h around $h = 0$ and the function $\nabla_{p_1} K$ into their Taylor series and comparing the coefficients of $e^{ik\omega t}$ yields for $k = 0$ (for the sake of clarity, we suppress the argument t in the coefficient functions)

$$\begin{aligned} &\zeta_{h,1} + \frac{h}{2} \dot{\zeta}_{h,1} - \zeta_{h,1} + \frac{h}{2} \dot{\zeta}_{h,1} + \mathcal{O}(h^3) \\ &= h \nabla_{p_1} K(\eta_{h,1}, \Phi \zeta_h) + \frac{h}{2} \sum_{s(\alpha)+s(\beta)=0} \frac{1}{m!n!} D_1^{m+1} D_2^n K(\eta_{h,1}, \Phi \zeta_h)(\eta_{h,1}^\alpha, \Phi \zeta_h^\beta) \\ &\quad + \frac{h}{2} \sum_{s(\alpha)+s(\beta)=0} \frac{1}{m!n!} e^{i\omega h/2(s(\beta)-s(\alpha))} D_1^{m+1} D_2^n K(\eta_{h,1}, \Phi \zeta_h)(\eta_{h,1}^\alpha, \Phi \zeta_h^\beta) + \mathcal{O}(h^2), \end{aligned}$$

where we used the same notations as in Section 2. This yields a relation for $\dot{\zeta}_{h,1}(t)$. Similarly, for $k \neq 0$, we obtain

$$\begin{aligned} \zeta_{h,1}^k &= \frac{h}{4i \sin\left(k\omega\frac{h}{2}\right)} \left(\sum_{s(\alpha)+s(\beta)=k} \frac{1}{m!n!} e^{-ik\omega h/2} D_1^{m+1} D_2^n K(\eta_{h,1}, \Phi_{\zeta_h})(\eta_{h,1}^\alpha, \Phi_{\zeta_h}^\beta) \right. \\ &\quad \left. + \sum_{s(\alpha)+s(\beta)=k} \frac{1}{m!n!} e^{i\omega h/2(s(\beta)-s(\alpha))} D_1^{m+1} D_2^n K(\eta_{h,1}, \Phi_{\zeta_h})(\eta_{h,1}^\alpha, \Phi_{\zeta_h}^\beta) \right) + \mathcal{O}(h^2). \end{aligned}$$

- Similar relations are obtained for the coefficient functions $\zeta_{h,2}(t)$, $\zeta_{h,2}^k(t)$ and $\dot{\eta}_{h,1}$, $\eta_{h,1}^k$.
- For the last formula of (3.1), we use the symmetry of the method, exchanging $n \leftrightarrow n+1$ and $h \leftrightarrow -h$, we get $p_{h,2}^n = \omega \sin(h\omega) q_{h,2}^{n+1} + \cos(h\omega) p_{h,2}^{n+1/2} + \frac{h}{2} \widehat{\psi}_2(h\omega) \nabla_{q_2} K(p_{h,1}^{n+1/2}, \Phi_{q_h}^n)$. Taking $n-1$ in place of n in this last expression and adding this quantity to $p_{h,2}^{n+1}$ yields

$$\begin{aligned} p_2^{n+1} + p_2^{n-1} &= \cos(h\omega)(p_2^{n+1/2} + p_2^{n-1/2}) \\ &\quad + \frac{h}{2} \widehat{\psi}_2(h\omega) (\nabla_{q_2} K(p_1^{n-1/2}, \Phi_{q_h}^{n-1}) - \nabla_{q_2} K(p_1^{n+1/2}, \Phi_{q_h}^{n+1})). \end{aligned}$$

Inserting (4.8) and using the fact that $p_1^{n-1/2} = p_1^{n-1} + \mathcal{O}(h)$ and $p_1^{n+1/2} = p_1^n + \mathcal{O}(h)$, we obtain

$$\begin{aligned} &p_{h,2}\left(t + \frac{h}{2}\right) + p_{h,2}\left(t - \frac{3h}{2}\right) \\ &= 2 \cos(h\omega) p_{h,2}\left(t - \frac{h}{2}\right) + \frac{h}{2} \cos(h\omega) \widehat{\psi}_2(h\omega) \left(\nabla_{q_2} K\left(p_{h,1}\left(t - \frac{3h}{2}\right), \Phi_{q_h}\left(t - \frac{h}{2}\right)\right) \right. \\ &\quad \left. - \nabla_{q_2} K\left(p_{h,1}\left(t - \frac{h}{2}\right), \Phi_{q_h}\left(t - \frac{h}{2}\right)\right) \right) \\ &\quad + \frac{h}{2} \widehat{\psi}_2(h\omega) \left(\nabla_{q_2} K\left(p_{h,1}\left(t - \frac{3h}{2}\right), \Phi_{q_h}\left(t - \frac{3h}{2}\right)\right) \right. \\ &\quad \left. - \nabla_{q_2} K\left(p_{h,1}\left(t - \frac{h}{2}\right), \Phi_{q_h}\left(t + \frac{h}{2}\right)\right) \right) + \mathcal{O}(h^2). \end{aligned} \tag{4.10}$$

This relation is true for every t , so we can exchange t with $t + \frac{h}{2}$. Using the operator

$$\mathcal{L}(hD) = e^{hD} - 2 \cos(h\Omega) + e^{-hD} = 4 \sin\left(\frac{h}{2}h\Omega + \frac{1}{2}ihD\right) \sin\left(\frac{h}{2}h\Omega - \frac{1}{2}ihD\right),$$

defined in Hairer *et al.* (2002, Chapter XIII), we can rewrite formula (4.10) as

$$\begin{aligned} \mathcal{L}(hD)p_{h,2}(t) &= \frac{h}{2} \cos(h\omega) \widehat{\psi}_2(h\omega) (\nabla_{q_2} K(p_{h,1}(t-h), \Phi_{q_h}(t)) - \nabla_{q_2} K(p_{h,1}(t), \Phi_{q_h}(t))) \\ &\quad + \frac{h}{2} \widehat{\psi}_2(h\omega) (\nabla_{q_2} K(p_{h,1}(t-h), \Phi_{q_h}(t-h)) \\ &\quad - \nabla_{q_2} K(p_{h,1}(t), \Phi_{q_h}(t+h))) + \mathcal{O}(h^2). \end{aligned}$$

Now, by the hypothesis (4.2) on N , the dominant terms in the Taylor expansions of $\mathcal{L}(hD)$ and $\mathcal{L}(hD + i h k \omega)$ give the desired first terms for the series of the coefficient functions $\eta_{h,2}^k$. Indeed, we have

$$\begin{aligned}\eta_{h,2}(t) &= \frac{h}{8s_1^2} \widehat{\psi}_2(h\omega) \cos(h\omega)(\dots) + \frac{h}{8s_1^2} \widehat{\psi}_2(h\omega)(\dots) + \mathcal{O}(h^2), \\ \dot{\eta}_{h,2}^1(t) &= \frac{1}{4is_2} \widehat{\psi}_2(h\omega) \cos(h\omega)(\dots) + \frac{1}{4is_2} \widehat{\psi}_2(h\omega)(\dots) + \mathcal{O}(h), \\ \eta_{h,2}^k(t) &= -\frac{h}{8s_{k-1}s_{k+1}} \widehat{\psi}_2(h\omega) \cos(h\omega)(\dots) - \frac{h}{8s_{k-1}s_{k+1}} \widehat{\psi}_2(h\omega)(\dots) + \mathcal{O}(h^2),\end{aligned}$$

where we used the abbreviation $s_k = \sin(\frac{k}{2}h\omega)$, and where the (\dots) terms are big expressions involving sums like those encountered in the formulas for $\zeta_{h,1}^k$ (see above).

This motivates the ansatz

$$\begin{aligned}\dot{\zeta}_{h,1} &= f_{10}(\cdot) + h f_{11}(\cdot) + \dots, \\ \dot{\eta}_{h,1} &= g_{10}(\cdot) + h g_{11}(\cdot) + \dots, \\ \dot{\eta}_{h,2}^1 &= \frac{\widehat{\Psi}_2}{s_2} (g_{20}^1(\cdot) + h g_{21}^1(\cdot) + \dots), \\ \zeta_{h,1}^k &= \frac{h}{s_k} (f_{10}^k(\cdot) + h f_{11}^k(\cdot) + \dots), \\ \eta_{h,1}^k &= \frac{h}{s_k} (g_{10}^k(\cdot) + h g_{11}^k(\cdot) + \dots), \\ \zeta_{h,2}^k &= h f_{21}^k(\cdot) + h^2 f_{22}^k(\cdot) + \dots, \\ \eta_{h,2} &= \frac{h \widehat{\Psi}_2}{s_1^2} (g_{20}(\cdot) + h g_{21}(\cdot) + \dots), \\ \eta_{h,2}^k &= \frac{h \widehat{\Psi}_2}{s_{k-1}s_{k+1}} (g_{20}^k(\cdot) + h g_{21}^k(\cdot) + \dots),\end{aligned}\tag{4.11}$$

where the dots stand for power series in h with coefficient functions f_{mn}^k and g_{mn}^k depending on the variables $\zeta_{h,1}$, $\eta_{h,1}$, $\eta_{h,2}^1$ and $h\omega$. The series present in the ansatz usually diverge, we thus truncate them after the $\mathcal{O}(h^N)$ terms. Inserting this ansatz into the numerical method (3.1) and comparing powers of h yields recurrence relations for the bounded functions f_{mn}^k and g_{mn}^k .

Initial values and bounds (4.7). The conditions $p_{h,1}(0) = p_1(0)$, $p_{h,2}(0) = p_2(0)$, $q_{h,1}(0) = q_1(0)$ and $p_{h,2}(h) = p_2(h)$, give the system

$$\begin{aligned}p_1(0) &= \eta_{h,1}(0) + \mathcal{O}(\omega^{-2}), \\ p_2(0) &= 2 \operatorname{Re}(\eta_{h,2}^1(0)) + \mathcal{O}(\omega^{-1}), \\ q_1(0) &= \zeta_{h,1}(0) + \mathcal{O}(\omega^{-2}), \\ \omega q_2(0) &= 2 \operatorname{Im}(\eta_{h,2}^1(0)) + \mathcal{O}(\omega^{-1}),\end{aligned}$$

which can be solved using the implicit function theorem to yield locally the desired initial values $\eta_{h,1}(0)$, $\zeta_{h,1}(0)$ and $\eta_{h,2}^1(0)$ for the differential equations appearing in the ansatz. The assumption (4.1) on the initial values of the problem and the hypothesis on the filter functions (4.3) and (4.11) imply that $\eta_{h,2}^1(t) = \mathcal{O}(1)$, for $0 \leq t \leq T$.

Using the hypothesis on the filter functions (4.3) and taking a closer look at the functions f_{mn}^k and g_{mn}^k (which contain at least k times the small factors ζ_h^1 or $\eta_{h,1}^1$) give the bounds (4.7) on the coefficient functions of the modulated Fourier expansion.

Defect. Let us define the components of the defect, for $t = nh$,

$$\begin{aligned} d_1(t) &= q_{h,1}(t+h) - q_{h,1}(t) - \frac{h}{2} \left(\nabla_{p_1} K \left(\widehat{p}_{h,1} \left(t + \frac{h}{2} \right), \Phi_{q_h}(t) \right) \right. \\ &\quad \left. + \nabla_{p_1} K \left(\widehat{p}_{h,1} \left(t + \frac{h}{2} \right), \Phi_{q_h}(t+h) \right) \right), \\ d_2(t) &= p_{h,1}(t+h) - p_{h,1}(t) + \frac{h}{2} \left(\nabla_{q_1} K \left(\widehat{p}_{h,1} \left(t + \frac{h}{2} \right), \Phi_{q_h}(t) \right) \right. \\ &\quad \left. + \nabla_{q_1} K \left(\widehat{p}_{h,1} \left(t + \frac{h}{2} \right), \Phi_{q_h}(t+h) \right) \right), \\ d_3(t) &= q_{h,2}(t+h) - \cos(h\omega)q_{h,2}(t) - h \operatorname{sinc}(h\omega)p_{h,2}(t) \\ &\quad + \frac{h^2}{2} \operatorname{sinc}(h\omega) \widehat{\psi}_2(h\omega) \nabla_{q_2} K \left(\widehat{p}_{h,1} \left(t + \frac{h}{2} \right), \Phi_{q_h}(t) \right), \\ d_4(t) &= p_{h,2}(t+h) + \omega \sin(h\omega)q_{h,2}(t) - \cos(h\omega)p_{h,2}(t) \\ &\quad + \frac{h}{2} \operatorname{sinc}(h\omega) \widehat{\psi}_2(h\omega) \nabla_{q_2} K \left(\widehat{p}_{h,1} \left(t + \frac{h}{2} \right), \Phi_{q_h}(t) \right). \end{aligned}$$

By definition of the coefficient functions ζ_h^k and η_h^k , we have $d_1(t) = d_2(t) = d_3(1) = \mathcal{O}(h^N)$. For the fourth component of the defect, we have to use the two-step formulation for $p_{h,2}$, this gives $d_4(t+h) + d_4(t-h) = \mathcal{O}(h^N)$. With our choice for the initial values, the defect at $t = 0$ is $d_4(0) = \mathcal{O}(h^N)$, so that we have $d_4(t) = \mathcal{O}(h^N) + \mathcal{O}(th^{N-1})$.

We still have to estimate the remainders (4.6). To do this, we define $R^n = \|p^n - p_h(t)\|$, $S^n = \|q^n - q_h(t)\|$ and the norm $\|(S_1, R_1, S_2, R_2)\|_* = \|(S_1, R_1, \omega S_2, R_2)\|$. To estimate the remainders (4.6), we first have to estimate the difference $p_1^{n+1/2} - \widehat{p}_{h,1}(t + \frac{h}{2})$. Using the definition (4.9) and the fact that the gradient of the function $K(p_1, q)$ satisfies a Lipschitz condition, we obtain

$$\left\| p_1^{n+1/2} - \widehat{p}_{h,1} \left(t + \frac{h}{2} \right) \right\| \leq \|R_1^n\| + C_1 h \left\| p_1^{n+1/2} - \widehat{p}_{h,1} \left(t + \frac{h}{2} \right) \right\| + C_2 h \|S^n\|,$$

for some constants C_j . This gives

$$\left\| p_1^{n+1/2} - \widehat{p}_{h,1} \left(t + \frac{h}{2} \right) \right\| \leq \alpha, \quad \text{where } \alpha = \frac{1}{1 - C_1 h} (\|R_1^n\| + C_2 h \|S^n\|).$$

Similarly, we obtain

$$\begin{aligned} \|(S_1, R_1, S_2, R_2)^{n+1}\|_* &\leq \|(S_1, R_1, S_2, R_2)^n\|_* + h\kappa_1\alpha + h\kappa_2\|(S_1, R_1, S_2, R_2)^{n+1}\|_* \\ &\quad + h\kappa_3\|(S_1, R_1, S_2, R_2)^n\|_* + \kappa_4h^{N-1}, \end{aligned}$$

for some constants κ_j . Using this relation repeatedly and the fact that $\|(S_1, R_1, S_2, R_2)^0\|_* = \mathcal{O}(h^N)$ (by definition of the initial values), we obtain the following estimate for the remainders

$$\begin{aligned} \|(S_1, R_1, S_2, R_2)^n\|_* &\leq \left(\frac{1+h\kappa_1}{1-h\kappa_2}\right)^n \|(S_1, R_1, S_2, R_2)^0\|_* + \kappa_3(n+1)h^{N-1} \\ &\leq Cnh^{N-1}, \end{aligned}$$

where κ_j and C are some constants. This concludes the proof. \square

Next, we show that the modulation functions of the numerical solution have almost-invariants that are obtained like the one obtained for the exact solution. In the proof of Theorem 4.2, we show that the defects of the functions on the right-hand sides of the equalities in (4.5) inserted into the method (3.1) are small. This implies that the modulated functions satisfy the following system (this is to be compared with (2.3)–(2.4)):

$$\begin{aligned} \widehat{p}_h(t) - p_h\left(t - \frac{h}{2}\right) &= -\frac{h}{2}\widehat{\Psi}\nabla_q K\left(\widehat{p}_{h,1}(t), \Phi q_h\left(t - \frac{h}{2}\right)\right), \\ q_{h,1}\left(t + \frac{h}{2}\right) - q_{h,1}\left(t - \frac{h}{2}\right) &= \frac{h}{2}\left(\nabla_{p_1} K\left(\widehat{p}_{h,1}(t), \Phi q_h\left(t - \frac{h}{2}\right)\right) \right. \\ &\quad \left. + \nabla_{p_1} K\left(\widehat{p}_{h,1}(t), \Phi q_h\left(t + \frac{h}{2}\right)\right)\right) + \mathcal{O}(h^N), \\ p_{h,1}\left(t + \frac{h}{2}\right) - \widehat{p}_{h,1}(t) &= -\frac{h}{2}\nabla_{q_1} K\left(\widehat{p}_{h,1}(t), \Phi q_h\left(t + \frac{h}{2}\right)\right) + \mathcal{O}(h^N), \\ p_{h,2}\left(t + \frac{h}{2}\right) + \omega \sin(h\omega)q_{h,2}\left(t - \frac{h}{2}\right) - \cos(h\omega)\widehat{p}_{h,2}(t) \\ &= -\frac{h}{2}\widehat{\psi}_2(h\omega)\nabla_{q_2} K\left(\widehat{p}_{h,1}(t), \Phi q_h\left(t + \frac{h}{2}\right)\right) + \mathcal{O}(h^N), \\ q_{h,2}\left(t + \frac{h}{2}\right) - \cos(h\omega)q_{h,2}\left(t - \frac{h}{2}\right) &= h \operatorname{sinc}(h\omega)\widehat{p}_{h,2} + \mathcal{O}(h^N), \end{aligned} \tag{4.12}$$

where we recall $q_h(t) = \sum_{|k|<N} q_h^k(t)$, $p_h(t) = \sum_{|k|<N} p_h^k(t)$ and $\widehat{p}_h(t) = \sum_{|k|<N} \widehat{p}_h^k(t)$ with $q_h^k(t) = e^{ik\omega t} \zeta_h^k(t)$, $p_h^k(t) = e^{ik\omega t} \eta_h^k(t)$ and $\widehat{p}_h^k(t) = e^{ik\omega t} \zeta_h^k(t)$. Comparing the coefficients of $e^{ik\omega t}$, we get,

writing the resulting equations in terms of \widehat{p}_h^k , p_h^k and q_h^k ,

$$\begin{aligned}
\widehat{p}_h^k(t) - p_h^k\left(t - \frac{h}{2}\right) &= -\frac{h}{2} \widehat{\Psi} \Phi^{-1} \nabla_{q^{-k}} \mathcal{K}_h\left(\widehat{\mathbf{p}}_1(t), \mathbf{q}\left(t - \frac{h}{2}\right)\right), \\
q_{h,1}^k\left(t + \frac{h}{2}\right) - q_{h,1}^k\left(t - \frac{h}{2}\right) &= \frac{h}{2} \left(\nabla_{p_1^{-k}} \mathcal{K}_h\left(\widehat{\mathbf{p}}_1(t), \mathbf{q}\left(t - \frac{h}{2}\right)\right) + \nabla_{p_1^{-k}} \mathcal{K}_h\left(\widehat{\mathbf{p}}_1(t), \mathbf{q}\left(t + \frac{h}{2}\right)\right) \right) + \mathcal{O}(h^N), \\
p_{h,1}^k\left(t + \frac{h}{2}\right) - \widehat{p}_{h,1}^k(t) &= -\frac{h}{2} \nabla_{q_1^{-k}} \mathcal{K}_h\left(\widehat{\mathbf{p}}_1(t), \mathbf{q}\left(t + \frac{h}{2}\right)\right) + \mathcal{O}(h^N), \\
p_{h,2}^k\left(t + \frac{h}{2}\right) + \omega \sin(h\omega) q_{h,2}^k\left(t - \frac{h}{2}\right) - \cos(h\omega) \widehat{p}_{h,2}^k(t) &= -\frac{h}{2} \widehat{\psi}_2(h\omega) \phi_2^{-1}(h\omega) \nabla_{q_2^{-k}} \mathcal{K}_h\left(\widehat{\mathbf{p}}_1(t), \mathbf{q}\left(t + \frac{h}{2}\right)\right) + \mathcal{O}(h^N), \\
q_{h,2}^k\left(t + \frac{h}{2}\right) - \cos(h\omega) q_{h,2}^k\left(t - \frac{h}{2}\right) &= h \operatorname{sinc}(h\omega) \widehat{p}_{h,2}^k(t) + \mathcal{O}(h^N),
\end{aligned} \tag{4.13}$$

where, similar to (2.5), we define

$$\mathcal{K}_h(\widehat{\mathbf{p}}_1, \mathbf{q}) = K(\widehat{p}_1^0, \Phi q^0) + \sum_{s(\alpha)+s(\beta)=0} \frac{1}{m!n!} D_1^m D_2^n K(\widehat{p}_1^0, \Phi q^0)(\widehat{\mathbf{p}}_1^\alpha, (\Phi \mathbf{q})^\beta), \tag{4.14}$$

for a vector $\widehat{\mathbf{p}}_1 = (\widehat{p}_{h,1}^{-N+1}, \dots, \widehat{p}_{h,1}^0, \dots, \widehat{p}_{h,1}^{N-1})$ and $\widehat{p}_{h,1}^k = e^{ik\omega t} \zeta_{h,1}^k(t)$, where $\zeta_{h,1}^k(t)$ are the modulation functions of $\widehat{p}_{h,1}(t)$. The same notation is used for \mathbf{q} . From here on, we omit the index h in the modulation functions.

As for the exact solution, the modulation functions $\boldsymbol{\eta} = (\eta^{-N+1}, \dots, \eta^{N-1})$ and $\boldsymbol{\zeta} = (\zeta^{-N+1}, \dots, \zeta^{N-1})$ have two formal invariants. We now give the result concerning the first one.

LEMMA 4.1 Under the assumptions of Theorem 4.1, the coefficient functions $\boldsymbol{\eta}$ and $\boldsymbol{\zeta}$ of the modulated Fourier expansion of the numerical solution satisfy

$$\widehat{\mathcal{H}}_h[\boldsymbol{\eta}, \boldsymbol{\zeta}](t) = \widehat{\mathcal{H}}_h[\boldsymbol{\eta}, \boldsymbol{\zeta}](0) + \mathcal{O}(th^N), \tag{4.15}$$

for $0 \leq t \leq T$. Moreover, we have

$$\widehat{\mathcal{H}}_h[\boldsymbol{\eta}, \boldsymbol{\zeta}](t) = 2\omega^2 \mu(h\omega) (\zeta_2^{-1})^T \zeta_2^{-1} + K(\eta_1, \Phi \zeta) + \mathcal{O}(h). \tag{4.16}$$

Proof. The idea of the proof is to multiply the relations in (4.13) by a derivative of some coefficient functions, then we take the sum over all k with $|k| < N$ and show that the resulting formula is in fact a total derivative of a function, say $\widehat{\mathcal{H}}_h[\boldsymbol{\eta}, \boldsymbol{\zeta}](t)$.

After multiplications and summations, we get from (4.13) that

$$\begin{aligned}
 & \sum_{|k| < N} \left\{ -\dot{q}^{-k} \left(t - \frac{h}{2} \right)^{\text{T}} \Phi \widehat{\Psi}^{-1} \left(\widehat{p}^k(t) - p^k \left(t - \frac{h}{2} \right) \right) \right. \\
 & \quad + \widehat{p}_1^{-k}(t)^{\text{T}} \left(q_1^k \left(t + \frac{h}{2} \right) - q_1^k \left(t - \frac{h}{2} \right) \right) - \dot{q}_1^{-k} \left(t + \frac{h}{2} \right)^{\text{T}} \left(p_1^k \left(t + \frac{h}{2} \right) - \widehat{p}_1^k(t) \right) \\
 & \quad \left. - \dot{q}_2^{-k} \left(t + \frac{h}{2} \right)^{\text{T}} \phi_2(h\omega) \widehat{\psi}_2^{-1}(h\omega) \left(p_2^k \left(t + \frac{h}{2} \right) \omega \sin(h\omega) q_2^k \left(t - \frac{h}{2} \right) - \cos(h\omega) \widehat{p}_2^k(t) \right) \right\} \\
 & = \frac{h}{2} \frac{\text{d}}{\text{d}t} \left\{ \mathcal{K}_h \left(\widehat{\mathbf{p}}_1(t), \mathbf{q} \left(t + \frac{h}{2} \right) \right) + \mathcal{K}_h \left(\widehat{\mathbf{p}}_1(t), \mathbf{q} \left(t - \frac{h}{2} \right) \right) \right\} + \mathcal{O}(h^N). \tag{4.17}
 \end{aligned}$$

Expanding the functions $\zeta^k(t \pm \frac{h}{2})$ and $\eta^k(t \pm \frac{h}{2})$ around t and replacing $\widehat{p}_2^k(t)$ by the last formula of (4.13) show that the left-hand side of this equation is a total derivative. Moving the terms from the left-hand side to right-hand side of the equation, we get

$$\frac{\text{d}}{\text{d}t} \widehat{\mathcal{H}}_h[\boldsymbol{\eta}, \boldsymbol{\zeta}](t) = \mathcal{O}(h^N),$$

and an integration yields statement (4.15) of the theorem.

This construction of $\widehat{\mathcal{H}}_h[\boldsymbol{\eta}, \boldsymbol{\zeta}](t)$, the bounds of the modulation functions, hypothesis (4.3) on the filter functions and the fact that we have $\eta_2^1 = i\omega\zeta_2^1 + \mathcal{O}(h^2)$ yield (4.16) and conclude the proof. \square

Concerning the second formal invariant, similar to formula (6.16) in Hairer *et al.* (2002, Section XIII.6), we have the following relation

$$\omega \sum_{0 < |k| < N} ik \left((\widehat{p}^k)^{\text{T}} \nabla_{p^k} \mathcal{K}_h(\widehat{\mathbf{p}}_1, \mathbf{q}) + (q^k)^{\text{T}} \nabla_{q^k} \mathcal{K}_h(\widehat{\mathbf{p}}_1, \mathbf{q}) \right) = 0, \tag{4.18}$$

for $\mathcal{K}_h(\widehat{\mathbf{p}}_1(t), \mathbf{q}(t))$ of (4.14). The same tricks as those used in the proof of the last lemma permit to prove the following lemma.

LEMMA 4.2 Under the assumptions of Theorem 4.1, the coefficient functions $\boldsymbol{\eta}$ and $\boldsymbol{\zeta}$ of the modulated Fourier expansion of the numerical solution satisfy

$$\widehat{\mathcal{F}}_h[\boldsymbol{\eta}, \boldsymbol{\zeta}](t) = \widehat{\mathcal{F}}_h[\boldsymbol{\eta}, \boldsymbol{\zeta}](0) + \mathcal{O}(th^N), \tag{4.19}$$

for $0 \leq t \leq T$. Moreover, we have

$$\widehat{\mathcal{F}}_h[\boldsymbol{\eta}, \boldsymbol{\zeta}](t) = 2\omega^2 \mu(h\omega) (\zeta_2^{-1})^{\text{T}} \zeta_2^{-1} + \mathcal{O}(h^2). \tag{4.20}$$

Proof. This time, we multiply and sum the equations in (4.13) in order to apply the identity (4.18). We get

$$\begin{aligned}
& i\omega \sum_{0 < |k| < N} k \left\{ -q^{-k} \left(t - \frac{h}{2} \right)^T \Phi \widehat{\Psi}^{-1} \left(\widehat{p}^k(t) - p^k \left(t - \frac{h}{2} \right) \right) + \widehat{p}_1^{-k}(t)^T \left(q_1^k \left(t + \frac{h}{2} \right) \right. \right. \\
& \quad \left. \left. - q_1^k \left(t - \frac{h}{2} \right) \right) - q_1^{-k} \left(t + \frac{h}{2} \right)^T \left(p_1^k \left(t + \frac{h}{2} \right) - \widehat{p}_1^k(t) \right) - q_2^{-k} \left(t + \frac{h}{2} \right)^T \phi_2(h\omega) \widehat{\psi}_2^{-1}(h\omega) \right. \\
& \quad \left. \times \left(p_2^k \left(t + \frac{h}{2} \right) + \omega \sin(h\omega) q_2^k \left(t - \frac{h}{2} \right) - \cos(h\omega) \widehat{p}_2^k(t) \right) \right\} \\
& = \frac{h\omega}{2} \sum_{0 < |k| < N} ik \left\{ \widehat{p}^k(t)^T \nabla_{p^k} \mathcal{K}_h \left(\widehat{\mathbf{p}}_1(t), \mathbf{q} \left(t - \frac{h}{2} \right) \right) \right. \\
& \quad \left. + q^k \left(t - \frac{h}{2} \right)^T \nabla_{q^k} \mathcal{K}_h \left(\widehat{\mathbf{p}}_1(t), \mathbf{q} \left(t - \frac{h}{2} \right) \right) \right. \\
& \quad \left. + \widehat{p}^k(t)^T \nabla_{p^k} \mathcal{K}_h \left(\widehat{\mathbf{p}}_1(t), \mathbf{q} \left(t + \frac{h}{2} \right) \right) + q^k \left(t + \frac{h}{2} \right)^T \nabla_{q^k} \mathcal{K}_h \left(\widehat{\mathbf{p}}_1(t), \mathbf{q} \left(t + \frac{h}{2} \right) \right) \right\} + \mathcal{O}(h^N).
\end{aligned}$$

The left-hand side of this equation is again a total derivative. For the right-hand side, we have, using (4.18), $0 + \mathcal{O}(h^N)$. Thus, we get

$$\frac{d}{dt} \widehat{\mathcal{F}}_h[\boldsymbol{\eta}, \boldsymbol{\zeta}](t) = \mathcal{O}(h^N),$$

and an integration from 0 to t yields the result (4.19). As before, statement (4.20) follows from the bounds on the modulation functions. \square

We now return to the proof of Theorem 4.1. We see that for symplectic numerical methods, we have $\mu(h\omega) = 1$ and hence $\widehat{\mathcal{F}}_h[\boldsymbol{\eta}, \boldsymbol{\zeta}](nh) = I(p^n, q^n) + \mathcal{O}(h)$ and $\widehat{\mathcal{H}}_h[\boldsymbol{\eta}, \boldsymbol{\zeta}](nh) = H(p^n, q^n) + \mathcal{O}(h)$. This proves the theorem in the case of symplectic methods. The additional hypothesis on the function μ (see (4.4)) and the arguments given in the proof of Theorem 7.1 in Hairer *et al.* (2002, Section XIII.7) show that the numerical method (3.1) nearly preserves the total energy H and the oscillatory energy I over long time intervals as stated in Theorem 4.1.

5. Further generalization in the kinetic energy

In this section we apply techniques similar to those given above to Hamiltonian problems (1.1) with a small perturbation in the Hamiltonian function (1.2). In fact, we consider the Hamiltonian

$$H(p, q) = K(p_1, q) + \frac{1}{2} p_2^T p_2 + \frac{\omega^2}{2} q_2^T q_2 + S(p, q), \quad (5.1)$$

where $S(p, q)$ is a quadratic function in the variable p and satisfies $S(p_1, p_2, q_1, 0) = 0$ (i.e. it is small). Basically, the only thing that changes is that we have to add (for the notations, see Section 2)

$$\mathcal{S}(\mathbf{p}, \mathbf{q}) = S(p^0, q^0) + \sum_{s(\alpha)+s(\beta)=0} \frac{1}{m!n!} D_1^m D_2^n S(p^0, q^0)(\mathbf{p}^\alpha, \mathbf{q}^\beta),$$

to the function \mathcal{K} in (2.5).

Similarly to Theorem 2.2, we can show that the oscillatory energy (1.4) is nearly preserved along the exact solution of Hamiltonian problems with Hamiltonian function (5.1) over long time intervals.

5.1 Numerical methods

Concerning numerical methods to solve Hamiltonian problems with (5.1), we propose to make a splitting and obtain the following numerical method

$$\Phi_h = (\phi_{h/2}^S)^* \circ \phi_h \circ \phi_{h/2}^S, \quad (5.2)$$

where the $*$ denotes the adjoint method. For the numerical scheme ϕ_h , we take the numerical method described in Section 3, and for $\phi_{h/2}^S$, we take the explicit Euler method, the Störmer–Verlet scheme or the symplectic Euler method (all methods work equally well).

REMARK 5.1 Since the function $S(p, q)$ in (5.1) is small, we do not apply a filter function to the method $\phi_{h/2}^S$. It is however possible to adapt the following proofs to this case.

EXAMPLE 5.1 The motion of a triatomic molecule can be modelled by a Hamiltonian system with the Hamiltonian function (5.1). To describe the motion of such a molecule, we use polar coordinates, as shown in Fig. 4.

The third mass (m_3) is kept fixed, the angle between the other masses is stiff with stiffness constant $\omega/\sqrt{2}$ and the two other springs have a stiffness constant ω . The Hamiltonian reads

$$\begin{aligned} H(p_{r_1}, p_{r_2}, p_{\theta_1}, p_{\theta_2}, r_1, r_2, \theta_1, \theta_2) \\ = \frac{1}{2} \left(p_{r_1}^2 + p_{r_2}^2 + (r_2 + 1)^{-2} (p_{\theta_2} - p_{\theta_1})^2 + (r_1 + 1)^{-2} p_{\theta_1}^2 \right) \\ + \frac{\omega^2}{2} \left(r_1^2 + r_2^2 + \frac{\theta_1^2}{2} \right) + \frac{1}{2} \theta_2^2. \end{aligned} \quad (5.3)$$

The last term is just an external potential to keep the molecule moving. After suitable coordinate changes (see Appendix), this Hamiltonian becomes

$$\begin{aligned} H(p_1, p_{2,1}, p_{2,2}, p_{2,3}, q_1, q_{2,1}, q_{2,2}, q_{2,3}) \\ = \frac{1}{2} (p_1^2 + p_{2,1}^2 + p_{2,2}^2 + p_{2,3}^2) + \frac{\omega^2}{2} (q_{2,1}^2 + q_{2,2}^2 + q_{2,3}^2) \\ + \frac{1}{4} (q_1 - q_{2,3})^2 - \frac{1}{4} \frac{2q_{2,2} + q_{2,2}^2}{(1 + q_{2,2})^2} (p_1 - p_{2,3})^2 - \frac{1}{4} \frac{2q_{2,1} + q_{2,1}^2}{(1 + q_{2,1})^2} (p_1 + p_{2,3})^2, \end{aligned} \quad (5.4)$$

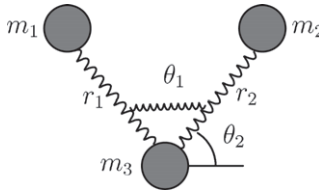


FIG. 4. Triatomic molecule.

which is of the form (5.1) with

$$S(p, q) = -\frac{1}{4} \frac{2q_{2,2} + q_{2,2}^2}{(1 + q_{2,2})^2} (p_1 - p_{2,3})^2 - \frac{1}{4} \frac{2q_{2,1} + q_{2,1}^2}{(1 + q_{2,1})^2} (p_1 + p_{2,3})^2.$$

Let us apply our numerical method to this problem with $\omega = 50$ and initial conditions $p(0) = 1$, $q_1(0) = 0.4$, $q_{2,1}(0) = q_{2,2}(0) = 1/\omega$, $q_{2,3}(0) = 1/(\sqrt{2}\omega)$. In Fig. 5, we plot the Hamiltonian H and the oscillatory energy I obtained by numerical method (5.2) (where we choose for $\phi_{h/2}^S$ the Störmer–Verlet method).

We note that, in this example, our numerical methods are symplectic, symmetric and explicit.

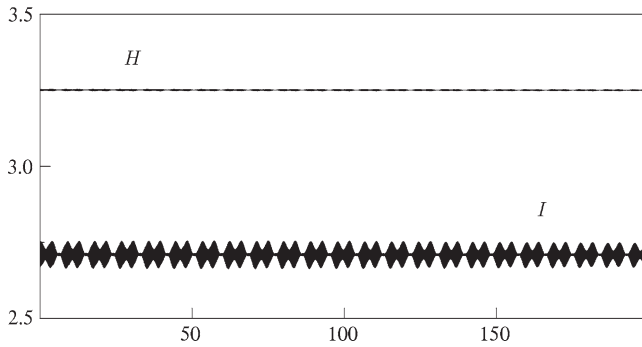


FIG. 5. Numerical solution of Hamiltonian problem with (5.4). For $\phi_{h/2}^S$, we use the Störmer–Verlet method. Step size $h = 0.01$.

The proofs given in the preceding section can be adapted to the numerical method (5.2). However, they become more technical and therefore we only mention some important points of the proofs.

5.2 Numerical energy conservation for the method (5.2)

Here, we prove the analogue of Theorem 4.1 for the numerical method (5.2) where, for the choice of the method $\phi_{h/2}^S$, we take the symplectic Euler method.

THEOREM 5.1 Under the assumptions of Theorem 4.1, we have, for the numerical solution (5.2),

$$\begin{aligned} H(p^n, q^n) &= H(p^0, q^0) + \mathcal{O}(h), \\ I(p^n, q^n) &= I(p^0, q^0) + \mathcal{O}(h), \end{aligned}$$

for $0 \leq nh \leq h^{-N+1}$.

The proof of this conservation result follows the lines of the proof of the aforementioned theorem: we first recall the numerical method and then give the analogues of the Lemmas 4.1 and 4.2. They help us to explain the near-conservation of the total and oscillatory energies for the numerical method (5.2) over long time intervals. Since the proof of the existence of a modulated Fourier expansion for the numerical scheme (5.2) is very similar to the one of Theorem 4.2, we do not give it here.

For our particular choice of $\phi_{h/2}^S$, the numerical scheme (5.2) now reads

$$\begin{aligned}
 \widehat{p}^n &= p^n - \frac{h}{2} \nabla_q S(\widehat{p}^n, q^n), \\
 \widehat{q}^n &= q^n + \frac{h}{2} \nabla_p S(\widehat{p}^n, q^n), \\
 \widetilde{p}^{n+1/2} &= \widehat{p}^n - \frac{h}{2} \widehat{\Psi} \nabla_q K(\widetilde{p}_1^{n+1/2}, \Phi \widehat{q}^n), \\
 \check{q}_1^{n+1} &= \widehat{q}_1^n + \frac{h}{2} (\nabla_{p_1} K(\widetilde{p}_1^{n+1/2}, \Phi \widehat{q}^n) + \nabla_{p_1} K(\widetilde{p}_1^{n+1/2}, \Phi \check{q}^{n+1})), \\
 \check{q}_2^{n+1} &= \cos(h\omega) \widehat{q}_2^n + h \operatorname{sinc}(h\omega) \widetilde{p}_2^{n+1/2}, \\
 \check{p}_1^{n+1} &= \widetilde{p}_1^{n+1/2} - \frac{h}{2} \nabla_{q_1} K(\widetilde{p}_1^{n+1/2}, \Phi \check{q}^{n+1}), \\
 \check{p}_2^{n+1} &= -\omega \sin(h\omega) \widehat{q}_2^n + \cos(h\omega) \widetilde{p}_2^{n+1/2} - \frac{h}{2} \widehat{\Psi}_2(h\omega) \nabla_{q_2} K(\widetilde{p}_1^{n+1/2}, \Phi \check{q}^{n+1}), \\
 p^{n+1} &= \check{p}^{n+1} - \frac{h}{2} \nabla_q S(\check{p}^{n+1}, q^{n+1}), \\
 q^{n+1} &= \check{q}^{n+1} + \frac{h}{2} \nabla_p S(\check{p}^{n+1}, q^{n+1}).
 \end{aligned}$$

As in (4.12), the coefficients of the modulated Fourier expansion of the numerical scheme (5.2) satisfy

$$\begin{aligned}
 \widehat{p}_h(t) &= p_h(t) - \frac{h}{2} \nabla_q S(\widehat{p}_h(t), q_h(t)), \\
 \widehat{q}_h(t) &= q_h(t) + \frac{h}{2} \nabla_p S(\widehat{p}_h(t), q_h(t)), \\
 \widetilde{p}_h(t) - \widehat{p}_h\left(t - \frac{h}{2}\right) &= -\frac{h}{2} \widehat{\Psi} \nabla_q K\left(\widetilde{p}_{h,1}(t), \Phi \widehat{q}_h\left(t - \frac{h}{2}\right)\right), \\
 \check{q}_{h,1}\left(t + \frac{h}{2}\right) - \widehat{q}_{h,1}\left(t - \frac{h}{2}\right) &= \frac{h}{2} \left(\nabla_{p_1} K\left(\widetilde{p}_{h,1}(t), \Phi \widehat{q}_h\left(t - \frac{h}{2}\right)\right) \right. \\
 &\quad \left. + \nabla_{p_1} K\left(\widetilde{p}_{h,1}(t), \Phi \check{q}_h\left(t + \frac{h}{2}\right)\right) \right), \\
 \check{p}_{h,1}\left(t + \frac{h}{2}\right) - \widetilde{p}_{h,1}(t) &= -\frac{h}{2} \nabla_{q_1} K\left(\widetilde{p}_{h,1}(t), \Phi \check{q}_h\left(t + \frac{h}{2}\right)\right), \\
 \check{p}_{h,2}\left(t + \frac{h}{2}\right) + \omega \sin(h\omega) \widehat{q}_{h,2}\left(t - \frac{h}{2}\right) - \cos(h\omega) \widetilde{p}_{h,2}(t) \\
 &= -\frac{h}{2} \widehat{\Psi}_2(h\omega) \nabla_{q_2} K\left(\widetilde{p}_{h,1}(t), \Phi \check{q}_h\left(t + \frac{h}{2}\right)\right), \\
 \check{q}_{h,2}\left(t + \frac{h}{2}\right) - \cos(h\omega) \widehat{q}_{h,2}\left(t - \frac{h}{2}\right) &= h \operatorname{sinc}(h\omega) \widetilde{p}_{h,2}, \\
 p_h(t) &= \check{p}_h(t) - \frac{h}{2} \nabla_q S(\check{p}_h(t), q_h(t)) + \mathcal{O}(h^N), \\
 q_h(t) &= \check{q}_h(t) + \frac{h}{2} \nabla_p S(\check{p}_h(t), q_h(t)) + \mathcal{O}(h^N),
 \end{aligned}$$

where we define $q_h(t) = \sum_{|k| < N} q_h^k(t)$ and $p_h(t) = \sum_{|k| < N} p_h^k(t)$ with $q_h^k(t) = e^{ik\omega t} \zeta_h^k(t)$ and $p_h^k(t) = e^{ik\omega t} \eta_h^k(t)$ (similar notations are used for $\widehat{p}_h(t)$, $\widehat{q}_h(t)$, $\check{p}_h(t)$, $\check{q}_h(t)$ and $\widetilde{p}_h(t)$). Comparing the coefficients of $e^{ik\omega t}$, we get, writing the resulting equations in terms of \widetilde{p}_h^k , p_h^k , q_h^k , \widehat{p}_h^k , \widehat{q}_h^k , \check{p}_h^k and \check{q}_h^k ,

$$\begin{aligned}
\widehat{p}_h^k(t) &= p_h^k(t) - \frac{h}{2} \nabla_{q^{-k}} \mathcal{S}_h(\widehat{\mathbf{p}}(t), \mathbf{q}(t)), \\
\widehat{q}_h^k(t) &= q_h^k(t) + \frac{h}{2} \nabla_{p^{-k}} \mathcal{S}_h(\widehat{\mathbf{p}}(t), \mathbf{q}(t)), \\
\widetilde{p}_h^k(t) - \widehat{p}_h^k\left(t - \frac{h}{2}\right) &= -\frac{h}{2} \widehat{\Psi} \Phi^{-1} \nabla_{q^{-k}} \mathcal{K}_h\left(\widetilde{\mathbf{p}}_1(t), \widehat{\mathbf{q}}\left(t - \frac{h}{2}\right)\right), \\
\check{q}_{h,1}^k\left(t + \frac{h}{2}\right) - \widehat{q}_{h,1}^k\left(t - \frac{h}{2}\right) &= \frac{h}{2} \left(\nabla_{p_1^{-k}} \mathcal{K}_h\left(\widetilde{\mathbf{p}}_1(t), \widehat{\mathbf{q}}\left(t - \frac{h}{2}\right)\right) + \nabla_{p_1^{-k}} \mathcal{K}_h\left(\widetilde{\mathbf{p}}_1(t), \check{\mathbf{q}}\left(t + \frac{h}{2}\right)\right) \right), \\
\check{p}_{h,1}^k\left(t + \frac{h}{2}\right) - \widetilde{p}_{h,1}^k(t) &= -\frac{h}{2} \nabla_{q_1^{-k}} \mathcal{K}_h\left(\widetilde{\mathbf{p}}_1(t), \check{\mathbf{q}}\left(t + \frac{h}{2}\right)\right), \\
\check{p}_{h,2}^k\left(t + \frac{h}{2}\right) + \omega \sin(h\omega) \widehat{q}_{h,2}^k\left(t - \frac{h}{2}\right) - \cos(h\omega) \widetilde{p}_{h,2}^k(t) &= -\frac{h}{2} \widehat{\Psi}_2(h\omega) \phi_2^{-1}(h\omega) \nabla_{q_2^{-k}} \mathcal{K}_h\left(\widetilde{\mathbf{p}}_1(t), \check{\mathbf{q}}\left(t + \frac{h}{2}\right)\right), \\
\check{q}_{h,2}^k\left(t + \frac{h}{2}\right) - \cos(h\omega) \widehat{q}_{h,2}^k\left(t - \frac{h}{2}\right) &= h \operatorname{sinc}(h\omega) \widetilde{p}_{h,2}^k(t), \\
\check{p}_h^k(t) &= \check{p}_h^k(t) - \frac{h}{2} \nabla_{q^{-k}} \mathcal{S}_h(\check{\mathbf{p}}(t), \mathbf{q}(t)) + \mathcal{O}(h^N), \\
q_h^k(t) &= \check{q}_h^k(t) + \frac{h}{2} \nabla_{p^{-k}} \mathcal{S}_h(\check{\mathbf{p}}(t), \mathbf{q}(t)) + \mathcal{O}(h^N),
\end{aligned} \tag{5.5}$$

where, similar to (2.5), we define

$$\mathcal{S}_h(\widehat{\mathbf{p}}, \mathbf{q}) = S(\widehat{p}^0, q^0) + \sum_{s(\alpha)+s(\beta)=0} \frac{1}{m!n!} D_1^m D_2^n K(\widehat{p}^0, q^0)(\widehat{\mathbf{p}}^\alpha, \mathbf{q}^\beta), \tag{5.6}$$

for a vector $\widehat{\mathbf{p}} = (\widehat{p}_h^{-N+1}, \dots, \widehat{p}_h^0, \dots, \widehat{p}_h^{N-1})$ and $\widehat{p}_h^k = e^{ik\omega t} \zeta_h^k(t)$, where $\zeta_h^k(t)$ are the modulation functions of $\widehat{p}_h(t)$. The same notation is used for \mathbf{q} and $\mathcal{S}_h(\check{\mathbf{p}}, \mathbf{q})$. From here on, we do not write the index h in the modulation functions.

As before, the modulation functions $\boldsymbol{\eta} = (\eta^{-N+1}, \dots, \eta^{N-1})$ and $\boldsymbol{\zeta} = (\zeta^{-N+1}, \dots, \zeta^{N-1})$ have two formal invariants. We now give the result concerning the first one.

LEMMA 5.1 Under the assumptions of Theorem 4.1, the coefficient functions $\boldsymbol{\eta}$ and $\boldsymbol{\zeta}$ of the modulated Fourier expansion of the numerical solution satisfy

$$\widehat{\mathcal{H}}_h[\boldsymbol{\eta}, \boldsymbol{\zeta}](t) = \widehat{\mathcal{H}}_h[\boldsymbol{\eta}, \boldsymbol{\zeta}](0) + \mathcal{O}(th^N), \tag{5.7}$$

for $0 \leq t \leq T$. Moreover, we have

$$\widehat{\mathcal{H}}_h[\boldsymbol{\eta}, \boldsymbol{\zeta}](t) = 2\omega^2 \mu(h\omega) (\zeta_2^{-1})^T \zeta_2^{-1} + K(\eta_1, \Phi \zeta) + \mathcal{O}(h). \tag{5.8}$$

Proof. To simplify the following proof, we consider the case $\mu(h\omega) = 1$ (i.e. the numerical method ϕ_h in (5.2) is symplectic).

Multiplying the relations in (5.5) (except those that contain the function $\mathcal{S}_h(\mathbf{p}, \mathbf{q})$, which will be used later) by the same coefficient functions as in (4.17) and summing up, we get

$$\begin{aligned} & \sum_{|k|<N} \left\{ -\hat{q}^{-k} \left(t - \frac{h}{2} \right)^{\text{T}} \left(\tilde{p}^k(t) - \hat{p}^k \left(t - \frac{h}{2} \right) \right) \right. \\ & \quad + \check{p}_1^{-k}(t)^{\text{T}} \left(\check{q}_1^k \left(t + \frac{h}{2} \right) - \hat{q}_1^k \left(t - \frac{h}{2} \right) \right) - \check{q}_1^{-k} \left(t + \frac{h}{2} \right)^{\text{T}} \left(\check{p}_1^k \left(t + \frac{h}{2} \right) - \tilde{p}_1^k(t) \right) \\ & \quad \left. - \left(\check{q}_2^{-k} \left(t + \frac{h}{2} \right) \right)^{\text{T}} \left(\check{p}_2^k \left(t + \frac{h}{2} \right) + \omega \sin(h\omega) \hat{q}_2^k \left(t - \frac{h}{2} \right) - \cos(h\omega) \tilde{p}_2^k(t) \right) \right\} \\ & = \frac{h}{2} \frac{\text{d}}{\text{d}t} \left\{ \mathcal{K}_h \left(\tilde{\mathbf{p}}_1(t), \hat{\mathbf{q}} \left(t - \frac{h}{2} \right) \right) + \mathcal{K}_h \left(\tilde{\mathbf{p}}_1(t), \check{\mathbf{q}} \left(t + \frac{h}{2} \right) \right) \right\}. \end{aligned}$$

Expanding the functions $\zeta^k(t \pm \frac{h}{2})$ and $\eta^k(t \pm \frac{h}{2})$ around t shows that the left-hand side of this equation is a total derivative. In contrast to the proof of Lemma 4.1, we have the following term,

$$\sum_{|k|<N} \left\{ \hat{q}^{-k} \left(t - \frac{h}{2} \right)^{\text{T}} \hat{p}^k \left(t - \frac{h}{2} \right) - \check{q}^{-k} \left(t + \frac{h}{2} \right)^{\text{T}} \check{p}^k \left(t + \frac{h}{2} \right) \right\},$$

which depends on the numerical method $\phi_{h/2}^S$. In order to show that this expression is also a total derivative, we insert the first two and last two formulas of (5.5) into it and get

$$\begin{aligned} & \sum_{|k|<N} \left\{ \dot{q}^{-k} \left(t - \frac{h}{2} \right)^{\text{T}} p^k \left(t - \frac{h}{2} \right) - \dot{q}^{-k} \left(t + \frac{h}{2} \right)^{\text{T}} p^k \left(t + \frac{h}{2} \right) \right. \\ & \quad - \frac{h}{2} \left(\dot{q}^{-k} \left(t - \frac{h}{2} \right)^{\text{T}} \nabla_{q^{-k}} \mathcal{S}_h \left(\hat{\mathbf{p}} \left(t - \frac{h}{2} \right), \mathbf{q} \left(t - \frac{h}{2} \right) \right) \right. \\ & \quad \left. - p^k \left(t - \frac{h}{2} \right)^{\text{T}} \frac{\text{d}}{\text{d}t} \nabla_{p^k} \mathcal{S}_h \left(\hat{\mathbf{p}} \left(t - \frac{h}{2} \right), \mathbf{q} \left(t - \frac{h}{2} \right) \right) \right. \\ & \quad \left. + \dot{q}^{-k} \left(t + \frac{h}{2} \right)^{\text{T}} \nabla_{q^{-k}} \mathcal{S}_h \left(\check{\mathbf{p}} \left(t + \frac{h}{2} \right), \mathbf{q} \left(t + \frac{h}{2} \right) \right) \right. \\ & \quad \left. - p^k \left(t + \frac{h}{2} \right)^{\text{T}} \frac{\text{d}}{\text{d}t} \nabla_{p^k} \mathcal{S}_h \left(\check{\mathbf{p}} \left(t + \frac{h}{2} \right), \mathbf{q} \left(t + \frac{h}{2} \right) \right) \right. \\ & \quad \left. - \frac{h^2}{4} \left(\nabla_{q^{-k}} \mathcal{S}_h \left(\hat{\mathbf{p}} \left(t - \frac{h}{2} \right), \mathbf{q} \left(t - \frac{h}{2} \right) \right) \right)^{\text{T}} \frac{\text{d}}{\text{d}t} \nabla_{p^k} \mathcal{S}_h \left(\hat{\mathbf{p}} \left(t - \frac{h}{2} \right), \mathbf{q} \left(t - \frac{h}{2} \right) \right) \right. \\ & \quad \left. \left. - \nabla_{q^{-k}} \mathcal{S}_h \left(\check{\mathbf{p}} \left(t + \frac{h}{2} \right), \mathbf{q} \left(t + \frac{h}{2} \right) \right) \right)^{\text{T}} \frac{\text{d}}{\text{d}t} \nabla_{p^k} \mathcal{S}_h \left(\check{\mathbf{p}} \left(t + \frac{h}{2} \right), \mathbf{q} \left(t + \frac{h}{2} \right) \right) \right\} + \mathcal{O}(h^N). \end{aligned}$$

The first two terms of this expression are in fact a total derivative. To show that the remaining terms are also a total derivative, we add and subtract

$$\dot{\hat{p}}^{-k} \left(t - \frac{h}{2} \right)^T \nabla_{p^{-k}} \mathcal{S}_h \left(\hat{\mathbf{p}} \left(t - \frac{h}{2} \right), \mathbf{q} \left(t - \frac{h}{2} \right) \right)$$

and

$$\dot{\check{p}}^{-k} \left(t + \frac{h}{2} \right)^T \nabla_{p^{-k}} \mathcal{S}_h \left(\check{\mathbf{p}} \left(t + \frac{h}{2} \right), \mathbf{q} \left(t + \frac{h}{2} \right) \right)$$

to make the total derivatives of $\mathcal{S}_h(\hat{\mathbf{p}}(t - \frac{h}{2}), \mathbf{q}(t - \frac{h}{2}))$, $p^k(t - \frac{h}{2})^T \nabla_{p^k} \mathcal{S}_h(\hat{\mathbf{p}}(t - \frac{h}{2}), \mathbf{q}(t - \frac{h}{2}))$ and $\nabla_{q^{-k}} \mathcal{S}_h(\hat{\mathbf{p}}(t - \frac{h}{2}), \mathbf{q}(t - \frac{h}{2}))^T \nabla_{p^k} \mathcal{S}_h(\hat{\mathbf{p}}(t - \frac{h}{2}), \mathbf{q}(t - \frac{h}{2}))$ appear (as well as the corresponding ones with argument $t + \frac{h}{2}$).

Moving the terms from the left-hand side to the right-hand side of the equation, we get

$$\frac{d}{dt} \widehat{\mathcal{H}}_h[\boldsymbol{\eta}, \boldsymbol{\zeta}](t) = \mathcal{O}(h^N),$$

and an integration yields statement (5.7) of the theorem.

This construction of $\widehat{\mathcal{H}}_h[\boldsymbol{\eta}, \boldsymbol{\zeta}](t)$, the bounds of the modulation functions, hypothesis (4.3) on the filter functions and the fact that we have $\eta_2^1 = i\omega\zeta_2^1 + \mathcal{O}(h^2)$ yield (5.8) and conclude the proof. \square

Concerning the second formal invariant, similarly to formula (4.18), we have the following equality

$$\omega \sum_{0 < |k| < N} ik((p^k)^T \nabla_{p^k} \mathcal{S}_h(\mathbf{p}, \mathbf{q}) + (q^k)^T \nabla_{q^k} \mathcal{S}_h(\mathbf{p}, \mathbf{q})) = 0, \quad (5.9)$$

for $\mathcal{S}_h(\mathbf{p}(t), \mathbf{q}(t))$ of (5.6). Tricks similar to those used in the proof of the last lemma help to prove the following lemma.

LEMMA 5.2 Under the assumptions of Theorem 4.1, the coefficient functions of the modulated Fourier expansion of the numerical solution satisfy

$$\widehat{\mathcal{I}}_h[\boldsymbol{\eta}, \boldsymbol{\zeta}](t) = \widehat{\mathcal{I}}_h[\boldsymbol{\eta}, \boldsymbol{\zeta}](0) + \mathcal{O}(th^N), \quad (5.10)$$

for $0 \leq t \leq T$. Moreover, we have

$$\widehat{\mathcal{I}}_h[\boldsymbol{\eta}, \boldsymbol{\zeta}](t) = 2\omega^2 \mu(h\omega) (\zeta_2^{-1})^T \zeta_2^1 + \mathcal{O}(h^2). \quad (5.11)$$

Proof. Again, for the sake of simplicity, we only give the proof for $\mu(h\omega) = 1$. This time, we multiply and add the equations in (5.5) in order to apply the identities (4.18) and (5.9). We get

$$\begin{aligned} & i\omega \sum_{0 < |k| < N} k \left\{ -\widehat{q}^{-k} \left(t - \frac{h}{2} \right)^T \left(\widetilde{p}^k(t) - \widehat{p}^k \left(t - \frac{h}{2} \right) \right) \right. \\ & \quad + \widetilde{p}_1^{-k}(t)^T \left(\check{q}_1^k \left(t + \frac{h}{2} \right) - \widehat{q}_1^k \left(t - \frac{h}{2} \right) \right) - \check{q}_1^{-k} \left(t + \frac{h}{2} \right)^T \left(\check{p}_1^k \left(t + \frac{h}{2} \right) - \widetilde{p}_1^k(t) \right) \\ & \quad \left. - \check{q}_2^{-k} \left(t + \frac{h}{2} \right)^T \left(\check{p}_2^k \left(t + \frac{h}{2} \right) + \omega \sin(h\omega) \widehat{q}_2^k \left(t - \frac{h}{2} \right) - \cos(h\omega) \widetilde{p}_2^k(t) \right) \right\} \end{aligned}$$

$$\begin{aligned}
 &= \frac{h\omega}{2} \sum_{0 < |k| < N} ik \left\{ \tilde{p}^k(t)^T \nabla_{p^k} \mathcal{K}_h \left(\tilde{\mathbf{p}}_1(t), \widehat{\mathbf{q}} \left(t - \frac{h}{2} \right) \right) \right. \\
 &\quad + \widehat{q}^k \left(t - \frac{h}{2} \right)^T \nabla_{q^k} \mathcal{K}_h \left(\tilde{\mathbf{p}}_1(t), \widehat{\mathbf{q}} \left(t - \frac{h}{2} \right) \right) + \tilde{p}^k(t)^T \nabla_{p^k} \mathcal{K}_h \left(\tilde{\mathbf{p}}_1(t), \check{\mathbf{q}} \left(t + \frac{h}{2} \right) \right) \\
 &\quad \left. + \check{q}^k \left(t + \frac{h}{2} \right)^T \nabla_{q^k} \mathcal{K}_h \left(\tilde{\mathbf{p}}_1(t), \check{\mathbf{q}} \left(t + \frac{h}{2} \right) \right) \right\}.
 \end{aligned}$$

Inserting the definition of the modulation functions corresponding to the symplectic Euler scheme and its adjoint, adding and subtracting

$$\widehat{p}^{-k} \left(t - \frac{h}{2} \right)^T \nabla_{p^{-k}} \mathcal{S}_h \left(\widehat{\mathbf{p}} \left(t - \frac{h}{2} \right), \mathbf{q} \left(t - \frac{h}{2} \right) \right)$$

and

$$\check{p}^{-k} \left(t + \frac{h}{2} \right)^T \nabla_{p^{-k}} \mathcal{S}_h \left(\check{\mathbf{p}} \left(t + \frac{h}{2} \right), \mathbf{q} \left(t + \frac{h}{2} \right) \right),$$

we see that the left-hand side of this equation is again a total derivative. Using (4.18), the right-hand side is zero. Thus, we get

$$\frac{d}{dt} \widehat{\mathcal{F}}_h[\boldsymbol{\eta}, \boldsymbol{\zeta}](t) = \mathcal{O}(h^N),$$

and an integration from 0 to t yields the result (5.10). As before, statement (5.11) follows from the bounds on the modulation functions. \square

These two lemmas explain the long-time conservation of the total and oscillatory energies along the numerical solution of the scheme (5.2), as stated in Theorem 5.1.

Finally, we would like to mention that the proofs given above can also be repeated for the composition

$$\Phi_h = \phi_{h/2}^S \circ \phi_h \circ \phi_{h/2}^S,$$

where $\phi_{h/2}^S$ is still the symplectic Euler method. This leads to conservation properties for a symplectic, non-symmetric numerical scheme.

6. The multi-frequency case

In this final section, we briefly discuss the multi-frequency case. We consider the Hamiltonian function (in accordance with the notations used in Cohen *et al.* (2005) and (5.1))

$$H(p, q) = K(p_1, q) + \frac{1}{2} \sum_{j=2}^l (p_j^T p_j + \omega_j^2 q_j^T q_j) + S(p, q), \quad (6.1)$$

where $q = (q_1, \dots, q_l)$ with $q_j \in \mathbb{R}^{d_j}$ (the same notation is used for p) and $\omega_j = \lambda_j \frac{1}{\varepsilon}$ with $\lambda_j \geq 1$ real distinct numbers and ε a small positive parameter.

Concerning the exact solution of Hamiltonian systems with the Hamiltonian (6.1), in complete analogy to Cohen *et al.* (2005, Theorem 7.1), we have the following result.

THEOREM 6.1 Let N be such that (weak non-resonance condition)

$$|k \cdot \lambda| \geq C\sqrt{\varepsilon} \quad \text{for } k \in \mathbb{Z}^{\ell-1} \setminus \mathcal{M} \text{ with } |k| \leq N,$$

where $k \cdot \lambda = k_2\lambda_2 + \dots + k_\ell\lambda_\ell$, $|k| = |k_2| + \dots + |k_\ell|$ and $\mathcal{M} = \{k \in \mathbb{Z}^{\ell-1} : k_2\lambda_2 + \dots + k_\ell\lambda_\ell = 0\}$. If the initial values satisfy (1.3), then, as long as the exact solution of the system stays in a compact set, we have

$$I_j(p(t), q(t)) = I_j(p(0), q(0)) + \mathcal{O}(\varepsilon) \quad \text{for } 0 \leq t \leq \varepsilon \cdot \min(\varepsilon^{-M+1}, \varepsilon^{-N}),$$

where

$$I_j(p(t), q(t)) = \frac{1}{2}(p_j^\top p_j + \omega_j^2 q_j^\top q_j), \quad (6.2)$$

for $j = 2, \dots, \ell$, and with $M = \min\{|k| : 0 \neq k \in \mathcal{M}\}$.

The idea of the proof is still to write the solution as a modulated Fourier expansion and to construct a system that determines the modulation functions of this expansion. One gets a system similar to (2.3)–(2.4) and finds almost-invariants related to (6.2).

For the multi-frequency case too, similar results have also been shown by Benettin *et al.* (1989).

Concerning the numerical solution, we extend method (5.2) to the multi-frequency case and obtain results concerning the near-conservation properties of the numerical solution similar to those given in Cohen *et al.* (2005). We make the following assumptions (see Section 4):

- The initial values satisfy

$$\frac{1}{2}\|p^0\|^2 + \frac{1}{2}\|\Omega q^0\|^2 \leq E.$$

- The numerical solution stays in a compact set.
- We impose a lower bound on the step size: $h/\omega \geq c_0 > 0$.
- We assume the numerical non-resonance condition:

$$\left| \sin\left(\frac{h}{2\varepsilon}k \cdot \lambda\right) \right| \geq c\sqrt{h} \quad \text{for all } k \in \mathbb{Z}^{\ell-1} \setminus \mathcal{M} \text{ with } |k| \leq N, \text{ with } N \geq 2.$$

- The function ψ satisfies, with $\xi_j = h\omega_j = h\lambda_j/\varepsilon$,

$$|\psi(\xi_j)| \leq C \left| \operatorname{sinc}\left(\frac{1}{2}\xi_j\right) \right| \quad \text{for } j = 2, \dots, \ell.$$

- We finally assume that

$$\begin{aligned} |\psi(\xi_j)| &\leq C \operatorname{sinc}^2\left(\frac{1}{2}\xi_j\right), \\ |\psi(\xi_j)| &\leq C|\phi(\xi_j)| \quad \text{for } j = 2, \dots, \ell. \end{aligned}$$

THEOREM 6.2 Under the above conditions, the numerical solution obtained by the method (5.2) satisfies

$$\begin{aligned} H(p^n, q^n) &= H(p^0, q^0) + \mathcal{O}(h) \quad \text{for } 0 \leq nh \leq \sigma_1 h \cdot \min(\varepsilon^{-M+1}, h^{-N}), \\ I_j(p^n, q^n) &= I_j(p^0, q^0) + \mathcal{O}(h) \quad \text{for } 0 \leq nh \leq \sigma_j h \cdot \min(\varepsilon^{-M+1}, h^{-N}), \end{aligned}$$

for $j = 2, \dots, \ell$. Here, $\sigma_j = |\sigma(\xi_j)|$ and $\sigma_1 = \min\{1, \sigma_2, \dots, \sigma_\ell\}$, where $\sigma(\xi) = \operatorname{sinc}(\xi)\phi(\xi)/\psi(\xi)$. The constants symbolized by \mathcal{O} are independent of n, h, ε and λ_j satisfying the above conditions, but depend on N and the constants in the conditions.

EXAMPLE 6.1 Taking different spring constants in Example 5.1, one can get a simple model of the water molecule. Following Izaguirre *et al.* (1999) (see also <http://amber.scripps.edu>), we take for the bond length constant $\omega_2 = \sqrt{553}$ and for the harmonic bond angle constant $\omega_3 = \sqrt{100}$. For such a molecule, the Hamiltonian (5.4) now reads

$$\begin{aligned} H(p_1, p_{2,1}, p_{2,2}, p_3, q_1, q_{2,1}, q_{2,2}, q_3) &= \frac{1}{2}(p_1^2 + p_{2,1}^2 + p_{2,2}^2 + p_3^2) + \frac{1}{2}(\omega_2^2 q_{2,1}^2 + \omega_2^2 q_{2,2}^2 + 2\omega_3^2 q_3^2) + \frac{1}{4}(q_1 - q_3)^2 \\ &+ \frac{1}{4} \left(\frac{1}{(r_0 + q_{2,2})^2} - 1 \right) (p_1 - p_3)^2 + \frac{1}{4} \left(\frac{1}{(r_0 + q_{2,1})^2} - 1 \right) (p_1 + p_3)^2, \end{aligned} \quad (6.3)$$

where $r_0 = 0.9572$ is the unstretched length of the springs. For initial values $p(0) = 0.5$, $q_1(0) = \sqrt{2}$, $q_{2,1}(0) = 1/\omega_2$, $q_{2,2}(0) = 1/\omega_2$ and $q_3(0) = 1/\omega_3$, we plot, in Fig. 6, the total and oscillatory energies and the first component of I along the numerical solution of the Hamiltonian system with Hamiltonian function (6.1).

As predicted, I is nearly preserved. This is not the case for I_2 , due to the fact that the frequencies ω_2 and ω_3 are not sufficiently large. Indeed, in Fig. 7, we plot the same quantities as in Fig. 6 but with a 10 times larger vector ω .

This time, all these quantities are nearly preserved.

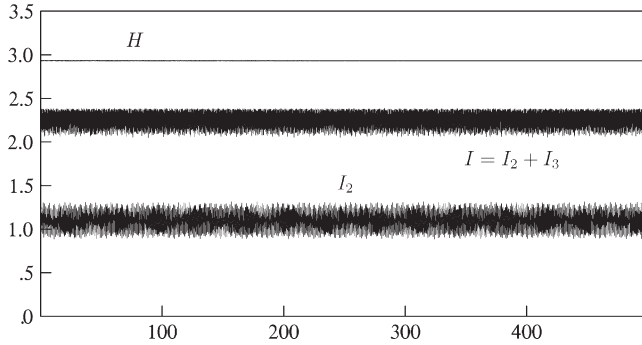


FIG. 6. Energies along the numerical solution of the Hamiltonian problem (6.3) with $h = 0.01$ and using for ϕ_h^S the Störmer–Verlet method.

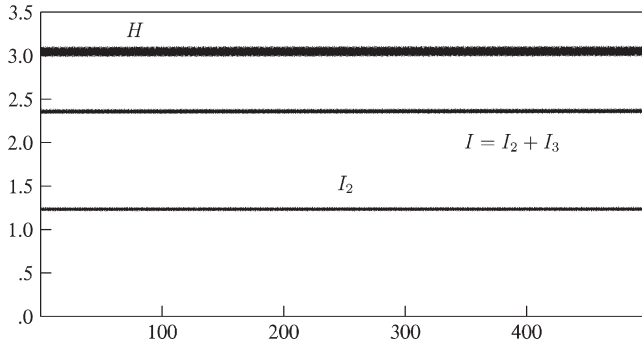


FIG. 7. Same plot as in Fig. 6 but with ω 10 times larger (and $h = 0.01$).

Acknowledgements

I am grateful to Christian Lubich and Ernst Hairer for interesting discussions on the subject. I wish to thank the anonymous referees for their helpful criticisms. This work was supported by the Fonds National Suisse.

REFERENCES

- ASCHER, U. & REICH, S. (1999a) The midpoint scheme and variants for Hamiltonian systems: advantages and pitfalls. *SIAM J. Sci. Comput.*, **21**, 1045–1065 (electronic).
- ASCHER, U. & REICH, S. (1999b) On some difficulties in integrating highly oscillatory Hamiltonian systems. *Computational Molecular Dynamics*. Lecture Notes in Computational Science and Engineering, vol. 4, pp. 281–296. Berlin: Springer.
- BENETTIN, G., GALGANI, L. & GIORGILLI, A. (1987) Realization of holonomic constraints and freezing of high frequency degrees of freedom in the light of classical perturbation theory. Part I. *Commun. Math. Phys.*, **113**, 87–103.
- BENETTIN, G., GALGANI, L. & GIORGILLI, A. (1989) Realization of holonomic constraints and freezing of high frequency degrees of freedom in the light of classical perturbation theory. Part II. *Commun. Math. Phys.*, **121**, 557–601.
- COHEN, D. (2004) Analysis and numerical treatment of highly oscillatory differential equations. *Ph.D. Thesis*, University of Geneva (www.unige.ch/cyberdocuments/theses2004/CohenD/meta.html).
- COHEN, D., HAIRER, E. & LUBICH, C. (2003) Modulated Fourier expansions of highly oscillatory differential equations. *Found. Comput. Math.*, **3**, 327–345.
- COHEN, D., HAIRER, E. & LUBICH, C. (2005) Numerical energy conservation for multi-frequency oscillatory differential equations. *BIT* (to appear).
- HAIRER, E. & LUBICH, C. (2000) Long-time energy conservation of numerical methods for oscillatory differential equations. *SIAM J. Numer. Anal.*, **38**, 414–441.
- HAIRER, E., LUBICH, C. & WANNER, G. (2002) *Geometric Numerical Integration. Structure-Preserving Algorithms for Ordinary Differential Equations*. Springer Series in Computational Mathematics 31. Berlin: Springer.
- HAIRER, E., NØRSETT, S. P. & WANNER, G. (1993) *Solving Ordinary Differential Equations I. Nonstiff Problems*. Springer Series in Computational Mathematics 8. Berlin: Springer.
- IZAGUIRRE, J. A., REICH, S. & SKEEL, R. D. (1999) Longer time steps for molecular dynamics. *J. Chem. Phys.*, **110**, 9853–9864.
- SANYAL, A. K., SHEN, J. & MCCLAMROCH, N. H. (2003) Dynamics and control of an Elastic Dumbbell Spacecraft in a central gravitational field. *Proceedings of the 42nd IEEE Conference on Decision and Control*, Maui, Hawaii, 2798–2803.
- SANYAL, A. K., SHEN, J. & MCCLAMROCH, N. H. (2005) Variational integrators for mechanical systems with configuration dependent inertia. Submitted to *Numerical Methods in Engineering*, under review (<http://math.la.asu.edu/sanyal/research/research.html>).

Appendix

A.1 Coordinate changes in Example 1.1

To obtain the Hamiltonian function (1.6), we first consider, as in Sanyal *et al.* (2003), the Lagrangian

$$L(\dot{r}, \dot{\phi}, \dot{\theta}, \dot{q}, r, \phi, \theta, q) = m(\dot{r}^2 + \dot{q}^2 + q^2\dot{\theta}^2 + 2q^2\dot{\theta}\dot{\phi} + (r^2 + q^2)\dot{\phi}^2) - V_g(r, \theta, q) - 2k(q - l)^2$$

and scale the variables: for $R > 0$, we define $\widehat{\omega} = \sqrt{\frac{\mu}{R^3}}$ and $\tau = \widehat{\omega}t$. We also define the new positions $\rho = \frac{r}{R}$ and $\sigma = \frac{q}{l}$. In the new variables, the Lagrangian function reads

$$L(\dot{\rho}, \dot{\phi}, \dot{\theta}, \dot{\sigma}, \rho, \phi, \theta, \sigma) = m\widehat{\omega}^2 R^2 \left\{ \dot{\rho}^2 + \varepsilon^2 \dot{\sigma}^2 + \varepsilon^2 \sigma^2 \dot{\theta}^2 + 2\varepsilon^2 \sigma^2 \dot{\theta} \dot{\phi} + (\rho^2 + \varepsilon^2 \sigma^2) \dot{\phi}^2 + \frac{1}{\rho} \left(2 - \varepsilon^2 \frac{\sigma^2}{\rho^2} (1 - 3 \cos^2(\theta)) \right) - 2\chi \varepsilon^2 (\sigma - 1)^2 \right\},$$

with $\varepsilon = \frac{l}{R}$ and $\chi = \frac{k}{m\widehat{\omega}^2}$. A last coordinate change, namely $\sigma = \varepsilon(\sigma - 1)$, leads to

$$L(\dot{\rho}, \dot{\phi}, \dot{\theta}, \dot{\sigma}, \rho, \phi, \theta, \sigma) = \frac{1}{2} \left(\dot{\rho}^2 + \rho^2 \dot{\phi}^2 + \dot{\sigma}^2 + (\varepsilon + \sigma)^2 (\dot{\phi} + \dot{\theta})^2 + \frac{2}{\rho} - \frac{(\varepsilon + \sigma)^2}{\rho^3} (1 - 3 \cos^2(\theta)) - 2\chi \sigma^2 \right),$$

where we have chosen the constants so that we obtain a factor $\frac{1}{2}$ in front of the Lagrangian. Finally, calculating the corresponding momentum, one gets the Hamiltonian function (1.6).

A.2 Coordinate changes in Example 5.1

To obtain the Hamiltonian (5.4), we rewrite (5.3) as

$$H(p, q) = \frac{1}{2} p^T M p + \frac{1}{2} q^T A q + \dots,$$

where the dots stand for small terms (i.e. terms containing r_1 or r_2), $q = (r_1, r_2, \theta_1, \theta_2)$ and

$$M = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 2 & -1 \\ 0 & 0 & -1 & 1 \end{pmatrix} \quad \text{and} \quad A = \begin{pmatrix} \omega^2 & 0 & 0 & 0 \\ 0 & \omega^2 & 0 & 0 \\ 0 & 0 & \omega^2/2 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}.$$

We make the symplectic change of coordinates $\widehat{p} = Cp$ and $\widehat{q} = Dq$ with the following matrices

$$C = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & \sqrt{2} & -\sqrt{2}/2 \\ 0 & 0 & 0 & \sqrt{2}/2 \end{pmatrix} \quad \text{and} \quad D = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & \sqrt{2}/2 & 0 \\ 0 & 0 & \sqrt{2}/2 & \sqrt{2} \end{pmatrix}.$$

The Hamiltonian function now reads $H(\widehat{p}, \widehat{q}) = \frac{1}{2} \widehat{p}^T \widehat{p} + \frac{1}{2} \widehat{q}^T \widehat{A} \widehat{q} + \dots$ with

$$\widehat{A} = \begin{pmatrix} \omega^2 & 0 & 0 & 0 \\ 0 & \omega^2 & 0 & 0 \\ 0 & 0 & \omega^2 + 1/2 & -1/2 \\ 0 & 0 & -1/2 & 1/2 \end{pmatrix},$$

and it is of the desired form (5.1).