

## Conservation of energy, momentum and actions in numerical discretizations of non-linear wave equations

COHEN, David, HAIRER, Ernst, LUBICH, Christian

### Abstract

For classes of symplectic and symmetric time-stepping methods - trigonometric integrators and the Störmer-Verlet or leapfrog method - applied to spectral semi-discretizations of semilinear wave equations in a weakly nonlinear setting, it is shown that energy, momentum, and all harmonic actions are approximately preserved over long times. For the case of interest where the CFL number is not a small parameter, such results are outside the reach of standard backward error analysis. Here, they are instead obtained via a modulated Fourier expansion in time.

### Reference

COHEN, David, HAIRER, Ernst, LUBICH, Christian. Conservation of energy, momentum and actions in numerical discretizations of non-linear wave equations. *Numerische Mathematik*, 2008, vol. 110, p. 113-143

DOI : 10.1007/s00211-008-0163-9

Available at:

<http://archive-ouverte.unige.ch/unige:5202>

Disclaimer: layout of this document may differ from the published version.



UNIVERSITÉ  
DE GENÈVE

# Conservation of energy, momentum and actions in numerical discretizations of nonlinear wave equations

David Cohen<sup>1</sup>, Ernst Hairer<sup>2</sup>, Christian Lubich<sup>3</sup>

<sup>1</sup> Mathematisches Institut, Univ. Basel, CH-4051 Basel, Switzerland,  
e-mail: David.Cohen@unibas.ch

<sup>2</sup> Dept. de Mathématiques, Univ. de Genève, CH-1211 Genève 24, Switzerland,  
e-mail: Ernst.Hairer@math.unige.ch

<sup>3</sup> Mathematisches Institut, Univ. Tübingen, D-72076 Tübingen, Germany,  
e-mail: Lubich@na.uni-tuebingen.de

Received: April 2007 / Revised version: date

**Summary** For classes of symplectic and symmetric time-stepping methods — trigonometric integrators and the Störmer–Verlet or leap-frog method — applied to spectral semi-discretizations of semilinear wave equations in a weakly nonlinear setting, it is shown that energy, momentum, and all harmonic actions are approximately preserved over long times. For the case of interest where the CFL number is not a small parameter, such results are outside the reach of standard backward error analysis. Here, they are instead obtained via a modulated Fourier expansion in time.

## 1 Introduction

This paper is concerned with the long-time behaviour of symplectic integrators applied to Hamiltonian nonlinear partial differential equations, such as semilinear wave equations. For symplectic methods applied to Hamiltonian systems of *ordinary* differential equations, the numerically observed long-time near-conservation of the total energy, and of actions in near-integrable systems, can be rigorously proved with the help of backward error analysis. This interprets a step of a symplectic method as the exact flow of a modified Hamiltonian system, up to an error which in the case of an analytic Hamiltonian is exponentially small in  $1/(h\omega)$ , where  $h$  is the small step

size and  $\omega$  represents the largest frequency in a local linearization of the system; see Benettin & Giorgilli [2], Hairer & Lubich [11], Reich [16], and Chapter IX in Hairer, Lubich & Wanner [14]. When the symplectic method is applied to a semi-discretization of a partial differential equation, however, then the product  $h\omega$  corresponds to the CFL number, which in typical computations is not small but of size 1. In this situation, the “exponentially small” remainder terms become of magnitude  $\mathcal{O}(1)$ , and no conclusions on the long-time behaviour of the method can then be drawn from the familiar backward error analysis. Nevertheless, long-time conservation of energy, and of momentum and actions when appropriate, *is* observed in numerical computations with symplectic methods used with reasonable CFL numbers. The present paper presents a proof of such conservation properties in the case of semilinear wave equations in the weakly nonlinear regime, over time scales that go far beyond linear perturbation arguments. To our knowledge, the results of this paper are the first results that rigorously prove the remarkable long-time conservation properties of symplectic integrators on a class of nonlinear partial differential equations.

We consider the one-dimensional nonlinear wave equation

$$u_{tt} - u_{xx} + \rho u + g(u) = 0 \quad (1)$$

for  $t > 0$  and  $-\pi \leq x \leq \pi$  subject to periodic boundary conditions. We assume  $\rho > 0$  and a nonlinearity  $g$  that is a smooth real function with  $g(0) = g'(0) = 0$ . We consider small initial data: in appropriate Sobolev norms, the initial values  $u(\cdot, 0)$  and  $u_t(\cdot, 0)$  are bounded by a small parameter  $\varepsilon$ . Notice that by rescaling  $u$ , this assumption could be rephrased as a  $\mathcal{O}(1)$  initial datum but a small non-linearity.

In Section 2 we recall the exact conservation of energy and momentum and, less obvious, the near-conservation of actions over long times  $t \leq \varepsilon^{-N}$ , where  $N$  only depends on a non-resonance condition on the frequencies, as shown by Bambusi [1] and Bourgain [3]. With the technique of modulated Fourier expansions that is central also to the present paper, the near-conservation of actions along solutions of (1) has been studied in our paper [6], and for spatial semi-discretizations of (1) by spectral methods in [13]. After discussing the semi-discretization in Section 3, we turn to the time discretization in Section 4.

We consider a class of symplectic and symmetric trigonometric integrators discussed in [14, Chap. XIII], and the familiar Störmer–Verlet or leapfrog method. In Section 4 we describe the trigonometric methods and present numerical experiments illustrating their

conservation properties, which appear particularly remarkable when confronted with the behaviour of a standard explicit Runge-Kutta method.

In Section 5 we state the main result of this paper, concerning the long-time near-conservation of energy, momentum and actions along numerical solutions in the full discretization. The result is proved in Sections 6 and 7, using the technique of *modulated Fourier expansions*. This approach was first used for studying long-time conservation properties of numerical methods for highly oscillatory ordinary differential equations with a single high frequency in [12], and later extended to several frequencies in [5]; see also [14, Chap. XIII] and further references given there. The extension of this technique to infinitely many frequencies, as occur in equation (1), was studied for the analytical problem in [6], and our treatment here essentially follows the lines of this previous work, with additional technical complications arising from the discretization.

In Section 8 we give similar long-time conservation results for the Störmer–Verlet/leapfrog method used with step sizes in the linear stability interval. These results follow from the previous ones by interpreting the leapfrog method as a trigonometric method with modified frequencies.

We are aware of two other papers that deal with long-time energy conservation of symplectic integrators for partial differential equations. Cano [4] also considers the nonlinear wave equation and aims at extending the classical backward error analysis to this situation. Long-time conservation properties are obtained under a list of unverified conditions formulated as conjectures. For symplectic splitting methods applied to the *linear* Schrödinger equation with a small potential, results on long-time energy conservation are given by Dujardin & Faou [8].

## 2 The nonlinear wave equation with small data

The semilinear wave equation (1) conserves several quantities along every solution  $(u(x, t), v(x, t))$ , with  $v = \partial_t u$ . The *total energy* or Hamiltonian, defined for  $2\pi$ -periodic functions  $u, v$  as

$$H(u, v) = \frac{1}{2\pi} \int_{-\pi}^{\pi} \left( \frac{1}{2} (v^2 + (\partial_x u)^2 + \rho u^2)(x) + U(u(x)) \right) dx, \quad (2)$$

where the potential  $U(u)$  is such that  $U'(u) = g(u)$ , and the *momentum*

$$K(u, v) = \frac{1}{2\pi} \int_{-\pi}^{\pi} \partial_x u(x) v(x) dx = - \sum_{j=-\infty}^{\infty} i j u_{-j} v_j \quad (3)$$

are exactly conserved along every solution  $(u(\cdot, t), v(\cdot, t))$  of (1). Here,  $u_j = \mathcal{F}_j u$  and  $v_j = \mathcal{F}_j v$  are the Fourier coefficients in the series  $u(x) = \sum_{j=-\infty}^{\infty} u_j e^{ijx}$  and correspondingly  $v(x)$ . Since we consider only real solutions, we note that  $u_{-j} = \bar{u}_j$  and  $v_{-j} = \bar{v}_j$ . In terms of the Fourier coefficients, equation (1) reads

$$\partial_t^2 u_j + \omega_j^2 u_j + \mathcal{F}_j g(u) = 0, \quad j \in \mathbb{Z}, \quad (4)$$

with the frequencies

$$\omega_j = \sqrt{\rho + j^2}.$$

The *harmonic actions*

$$I_j(u, v) = \frac{\omega_j}{2} |u_j|^2 + \frac{1}{2\omega_j} |v_j|^2, \quad (5)$$

for which we note  $I_{-j} = I_j$ , are conserved for the linear wave equation, that is, for  $g(u) \equiv 0$ . In the semilinear equation (1), they turn out to remain constant up to small deviations over long times for almost all values of  $\rho > 0$ , when the initial data are smooth and small [1, 3, 6]. We recall the precise statement of this result, because this will help to understand related assumptions for the numerical discretizations.

We work with the Sobolev space, for  $s \geq 0$ ,

$$H^s = \{v \in L^2(\mathbb{T}) : \|v\|_s < \infty\}, \quad \|v\|_s = \left( \sum_{j=-\infty}^{\infty} \omega_j^{2s} |v_j|^2 \right)^{1/2},$$

where  $v_j$  denote the Fourier coefficients of a  $2\pi$ -periodic function  $v$ . For the initial position and velocity we assume that for suitably large  $s$  and small  $\varepsilon$ ,

$$\left( \|u(\cdot, 0)\|_{s+1}^2 + \|v(\cdot, 0)\|_s^2 \right)^{1/2} \leq \varepsilon. \quad (6)$$

Since the analysis of the near-conservation of actions encounters problems with small denominators, we prepare for the formulation of a non-resonance condition. Consider sequences of integers  $\mathbf{k} = (k_\ell)_{\ell=0}^{\infty}$  with only finitely many  $k_\ell \neq 0$ . We denote  $|\mathbf{k}| = (|k_\ell|)_{\ell=0}^{\infty}$  and let

$$\|\mathbf{k}\| = \sum_{\ell=0}^{\infty} |k_\ell|, \quad \mathbf{k} \cdot \boldsymbol{\omega} = \sum_{\ell=0}^{\infty} k_\ell \omega_\ell, \quad \boldsymbol{\omega}^{\sigma|\mathbf{k}|} = \prod_{\ell=0}^{\infty} \omega_\ell^{\sigma|k_\ell|} \quad (7)$$

for real  $\sigma$ , where we use the notation  $\boldsymbol{\omega} = (\omega_\ell)_{\ell=0}^\infty$ . For  $j \in \mathbb{Z}$ , we write  $\langle j \rangle = (0, \dots, 0, 1, 0, \dots)$  with the only entry at the  $|j|$ -th position.

For an arbitrary fixed integer  $N \geq 1$  and for small  $\varepsilon > 0$ , we consider the set of near-resonant indices

$$\mathcal{R}_\varepsilon = \{(j, \mathbf{k}) : j \in \mathbb{Z}, \mathbf{k} \neq \pm \langle j \rangle, \|\mathbf{k}\| \leq 2N \text{ with } |\omega_j - |\mathbf{k} \cdot \boldsymbol{\omega}|| < \varepsilon^{1/2}\}. \quad (8)$$

We impose the following *non-resonance condition*: there are  $\sigma > 0$  and a constant  $C_0$  such that

$$\sup_{(j, \mathbf{k}) \in \mathcal{R}_\varepsilon} \frac{\omega_j^\sigma}{\boldsymbol{\omega}^{|\mathbf{k}|}} \varepsilon^{\|\mathbf{k}\|/2} \leq C_0 \varepsilon^N. \quad (9)$$

As is shown in [6], condition (9) is implied, for sufficiently large  $\sigma$ , by the non-resonance condition of Bambusi [1], which holds true for almost all (w.r.t. Lebesgue measure)  $\rho$  in any fixed interval of positive numbers.

**Theorem 1** [6, Theorem 1] *Under the non-resonance condition (9) and assumption (6) on the initial data with  $s \geq \sigma + 1$ , the estimate*

$$\sum_{\ell=0}^{\infty} \omega_\ell^{2s+1} \frac{|I_\ell(t) - I_\ell(0)|}{\varepsilon^2} \leq C\varepsilon \quad \text{for } 0 \leq t \leq \varepsilon^{-N+1}$$

with  $I_\ell(t) = I_\ell(u(\cdot, t), v(\cdot, t))$  holds with a constant  $C$  which depends on  $s, N$ , and  $C_0$ , but is independent of  $\varepsilon$  and  $t$ .

The smallness of the initial data, which implies that the non-linearity is small compared to the linear terms, is essential for our analysis. Since we do not impose any further restrictions on the non-linearity, such an assumption permits to avoid blow-up in finite time.

### 3 Spectral semi-discretization in space

For the numerical solution of (1) we first discretize in space (method of lines) and then in time (Section 4). Following [13], we consider pseudo-spectral semi-discretization in space with equidistant collocation points  $x_k = k\pi/M$  (for  $k = -M, \dots, M-1$ ). This yields an approximation in form of real-valued trigonometric polynomials

$$u^M(x, t) = \sum'_{|j| \leq M} q_j(t) e^{ijx}, \quad v^M(x, t) = \sum'_{|j| \leq M} p_j(t) e^{ijx} \quad (10)$$

where the prime indicates that the first and last terms in the sum are taken with the factor  $1/2$ . We have  $p_j(t) = \frac{d}{dt}q_j(t)$ , and the  $2M$ -periodic coefficient vector  $q(t) = (q_j(t))$  is a solution of the  $2M$ -dimensional system of ordinary differential equations

$$\frac{d^2q}{dt^2} + \Omega^2 q = f(q) \quad \text{with} \quad f(q) = -\mathcal{F}_{2M}g(\mathcal{F}_{2M}^{-1}q). \quad (11)$$

The matrix  $\Omega$  is diagonal with entries  $\omega_j$  for  $|j| \leq M$ , and  $\mathcal{F}_{2M}$  denotes the discrete Fourier transform:

$$(\mathcal{F}_{2M}w)_j = \frac{1}{2M} \sum_{k=-M}^{M-1} w_k e^{-ijx_k}.$$

Since the components of the nonlinearity in (11) are of the form

$$f_j(q) = -\frac{\partial}{\partial q_{-j}}V(q) \quad \text{with} \quad V(q) = \frac{1}{2M} \sum_{k=-M}^{M-1} U((\mathcal{F}_{2M}^{-1}q)_k),$$

we are concerned with a finite-dimensional complex Hamiltonian system with energy

$$H_M(q, p) = \frac{1}{2} \sum'_{|j| \leq M} (|p_j|^2 + \omega_j^2 |q_j|^2) + V(q), \quad (12)$$

which is exactly conserved along the solution  $(q(t), p(t))$  of (11) with  $p(t) = dq(t)/dt$ . We further consider the actions (for  $|j| \leq M$ ) and the momentum

$$I_j(q, p) = \frac{\omega_j}{2}|q_j|^2 + \frac{1}{2\omega_j}|p_j|^2, \quad K(q, p) = -\sum''_{|j| \leq M} ij q_{-j} p_j, \quad (13)$$

where the double prime indicates that the first and last terms in the sum are taken with the factor  $1/4$ . The definition of these expressions is motivated by the fact that they agree with the corresponding quantities of Section 2 along the trigonometric polynomials  $u^M, v^M$  (with the exception of  $I_{\pm M}$ , where a factor 4 has been included to get a unified formula). Since we are concerned with real approximations (10), the Fourier coefficients satisfy  $q_{-j} = \bar{q}_j$  and  $p_{-j} = \bar{p}_j$ , so that also  $I_{-j} = I_j$ .

On the space of  $2M$ -periodic sequences  $q = (q_j)$  we consider the weighted norm

$$\|q\|_s = \left( \sum''_{|j| \leq M} \omega_j^{2s} |q_j|^2 \right)^{1/2}, \quad (14)$$

which is defined such that it equals the  $H^s$  norm of the trigonometric polynomial with coefficients  $q_j$ . We assume that the initial data  $q(0)$  and  $p(0)$  satisfy a condition corresponding to (6):

$$\left( \|q(0)\|_{s+1}^2 + \|p(0)\|_s^2 \right)^{1/2} \leq \varepsilon. \quad (15)$$

**Theorem 2** [13, Theorems 3.1 and 3.2] *Under the non-resonance condition (9) with exponent  $\sigma$  and the assumption (15) of small initial data with  $s \geq \sigma + 1$ , the near-conservation estimates*

$$\begin{aligned} \sum_{\ell=0}^M \omega_\ell^{2s+1} \frac{|I_\ell(t) - I_\ell(0)|}{\varepsilon^2} &\leq C\varepsilon && \text{for } 0 \leq t \leq \varepsilon^{-N+1} \\ \frac{|K(t) - K(0)|}{\varepsilon^2} &\leq C t \varepsilon M^{-s-1} \end{aligned}$$

for actions  $I_\ell(t) = I_\ell(q(t), p(t))$  and momentum  $K(t) = K(q(t), p(t))$  hold with a constant  $C$  that depends on  $s$ ,  $N$ , and  $C_0$ , but is independent of  $\varepsilon$ ,  $M$ , and  $t$ .

Since the expression  $\sum_{\ell=0}^M \omega_\ell^{2s+1} I_\ell(t)$  is essentially (up to the factors in the boundary terms) equal to the squared  $H^{s+1} \times H^s$  norm of the solution  $(q(t), p(t))$ , Theorem 2 implies long-time spatial regularity:

$$\left( \|q(t)\|_{s+1}^2 + \|p(t)\|_s^2 \right)^{1/2} \leq \varepsilon(1 + C\varepsilon) \quad \text{for } t \leq \varepsilon^{-N+1}. \quad (16)$$

Theorems 1 and 2 have been included as a motivation of our results. They will not be used in the following.

#### 4 Full discretization and numerical phenomena

We consider the class of time discretization methods studied in [14, Chapter XIII], which gives the exact solution for linear problems (11) with  $f(q) = 0$ , and reduces to the Störmer–Verlet/leapfrog method for (11) with  $\Omega = 0$ :

$$q^{n+1} - 2 \cos(h\Omega) q^n + q^{n-1} = h^2 \Psi f(\Phi q^n), \quad (17)$$

where  $\Psi = \psi(h\Omega)$  and  $\Phi = \phi(h\Omega)$  with filter functions  $\psi$  and  $\phi$  that are real-valued, bounded, even, and satisfy  $\psi(0) = \phi(0) = 1$ . A velocity approximation  $p^n$  is obtained from

$$2h \operatorname{sinc}(h\Omega) p^n = q^{n+1} - q^{n-1} \quad (18)$$



provided that  $\text{sinc}(h\Omega)$  is invertible. Here we use the notation  $\text{sinc } \xi = \sin \xi / \xi$ .

For an implementation it is more convenient to work with an equivalent one-step mapping  $(q^n, p^n) \mapsto (q^{n+1}, p^{n+1})$ , which is obtained from adding and subtracting the formulas (17) and (18) and which reads

$$q^{n+1} = \cos(h\Omega)q^n + h \text{sinc}(h\Omega)p^n + \frac{1}{2}h^2 \Psi f(\Phi q^n) \quad (19)$$

$$p^{n+1} = -\Omega \sin(h\Omega)q^n + \cos(h\Omega)p^n + \frac{1}{2}h \left( \Psi_0 f(\Phi q^n) + \Psi_1 f(\Phi q^{n+1}) \right).$$

Here,  $\Psi_0 = \psi_0(h\Omega)$  and  $\Psi_1 = \psi_1(h\Omega)$ , where the functions  $\psi_i(\xi)$  are defined by the relations  $\psi(\xi) = \text{sinc}(\xi)\psi_1(\xi)$  and  $\psi_0(\xi) = \cos(\xi)\psi_1(\xi)$ . These methods are symmetric for all choices of  $\psi$  and  $\phi$ ; they are symplectic if

$$\psi(\xi) = \text{sinc}(\xi) \phi(\xi) \quad \text{for all real } \xi. \quad (20)$$

The methods (19) with this property are precisely the mollified impulse methods introduced in [9]. Interpreted as a splitting method, they can be extended for fully non-linear problems.

Condition (20) will be assumed in the following. We note, however, that for non-symplectic methods, the transformation of variables

$$\widehat{q}^n = \chi(h\Omega)q^n, \quad \widehat{p}^n = \chi(h\Omega)p^n, \quad (21)$$

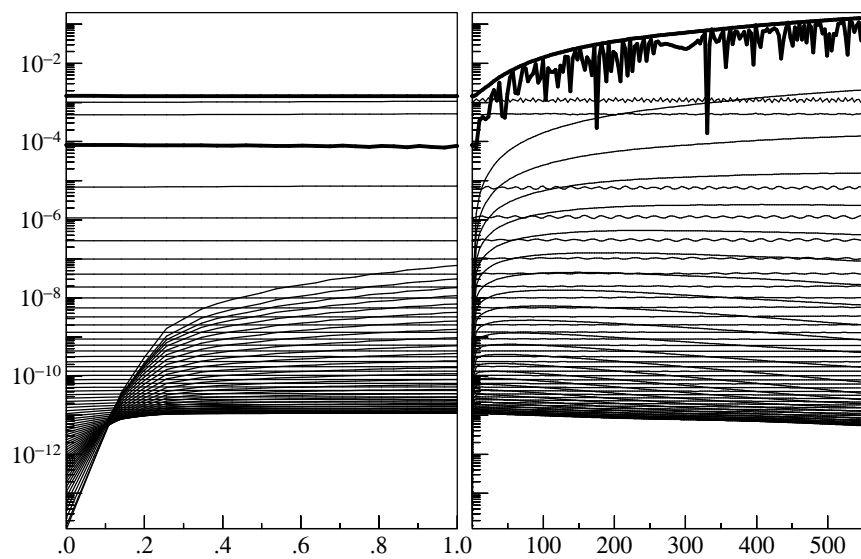
turns the method (19) into a symplectic method if  $\chi$  can be chosen as a positive solution of  $\chi(\xi)^2 = \phi(\xi) \text{sinc}(\xi) / \psi(\xi)$ .

In our numerical experiments we consider the nonlinear wave equation (1) with the following data:  $\rho = 0.5$ ,  $g(u) = -u^2$ , and initial data

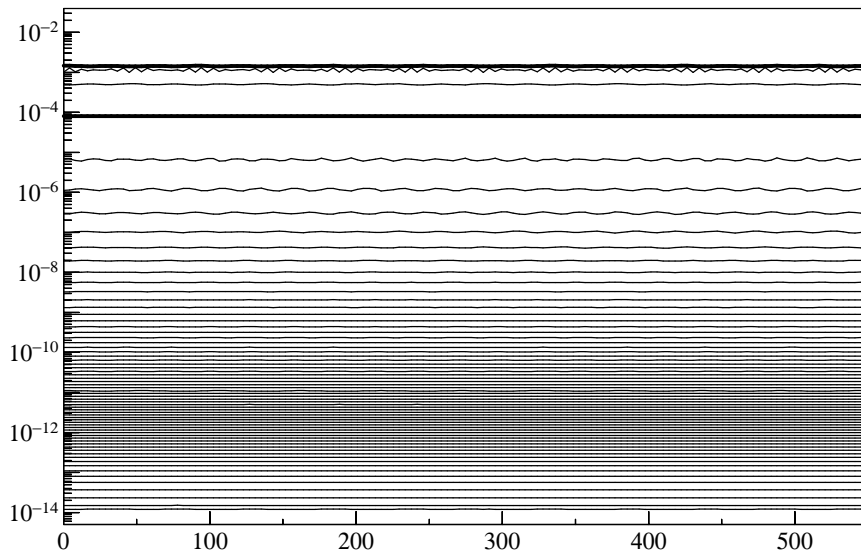
$$u(x, 0) = 0.1 \cdot \left(\frac{x}{\pi} - 1\right)^3 \left(\frac{x}{\pi} + 1\right)^2, \quad \partial_t u(x, 0) = 0.01 \cdot \frac{x}{\pi} \left(\frac{x}{\pi} - 1\right) \left(\frac{x}{\pi} + 1\right)^2$$

for  $-\pi \leq x \leq \pi$ . The spatial discretization is (11) with dimension  $2M = 2^7$ . Considered as  $2\pi$ -periodic functions, the initial data  $u(\cdot, 0)$  and  $\partial_t u(\cdot, 0)$  have a jump discontinuity in the second and first derivative, respectively. The assumption (15) is therefore satisfied for  $s < 1.5$ .

We first apply a standard explicit Runge–Kutta method in the variable stepsize implementation DOPRI5 of [15], with local error tolerances  $A_{tol} = 10^{-5}$  and  $R_{tol} = 10^{-4}$ . The program chose 32735 accepted steps for the integration over the interval  $0 \leq t \leq 550$ , which corresponds to an average stepsize  $\bar{h} = 0.0168$  and average CFL number  $\bar{h}\omega_M = 1.075$ . In both pictures of Figure 1 we plot the actions



**Fig. 1.** Actions, total energy (upper bold line), and momentum (lower bold line) along the numerical solution of DOPRI5, average CFL number 1.075.



**Fig. 2.** Actions, total energy (upper bold line), and momentum (lower bold line) along the numerical solution of the trigonometric integrator (19) with  $\psi = \text{sinc}$  and  $\phi = 1$  for the CFL number  $h\omega_M \approx 6.4$ .

$I_j$  of (5), the total energy  $H_M$  of (12), and the momentum  $K$  of (13) along the numerical solution. The left-hand picture illustrates that even on the short interval  $0 \leq t \leq 1$ , the actions with values below

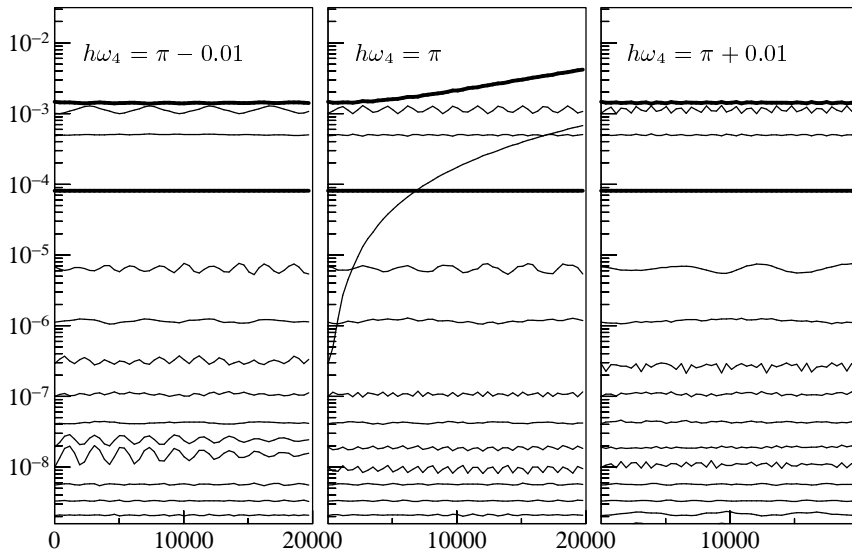


Fig. 3. Illustration of numerical resonance; method (19) with  $\psi = \text{sinc}$  and  $\phi = 1$

the tolerance are not at all conserved. The right-hand picture shows substantial drifts in all the quantities over a longer time interval.

We now consider method (19) with  $\psi = \text{sinc}$  and  $\phi = 1$ , which was originally proposed in [7]. The method can also be viewed as a special case of the impulse method used in molecular dynamics [10, 17]. We apply the method with stepsize  $h = 0.1$  to the above problem. The CFL number then is  $h\omega_M \approx 6.4$ . Figure 2 illustrates that energy, momentum and actions are very well conserved. Since the regularity of our initial data is not very high ( $s < 1.5$ ), this shows that the regularity assumption  $s \geq \sigma + 1$ , where  $\sigma = 2^9$  already for  $N = 2$  (c.f. [6]) can be relaxed in concrete examples when a fixed  $\rho$  is considered.

In a further experiment with the same problem, we choose step-sizes such that  $h\omega_4$  is close to  $\pi$  (Figure 3). In this situation of a numerical resonance, the action  $I_4$  is no longer preserved, which on longer time scales also affects the conservation of energy. The resonance behaviour depends strongly on the choice of the filter functions, cf. [14, Section XIII.2]. For example, with  $\phi = \text{sinc}$  and  $\psi = \text{sinc}^2$ , a method proposed in [9], no numerical resonance is visible.

## 5 Main results

To get rigorous statements on the good long-time behaviour illustrated in Section 4, we combine the techniques of [6], where the long-time preservation of the harmonic actions along exact solutions of

the semilinear wave equation (1) is shown, with those of [5], where the long-time behaviour of the numerical method (17) is studied for oscillatory Hamiltonian systems with a fixed number of large frequencies. As for spectral semi-discretizations (cf. [13]), we are interested in results that are valid uniformly in  $M$ , where  $2M$  is the dimension of the spatially discretized system (11).

The analytical tool for understanding the long-time behaviour of the numerical solution of (11) is given by a modulated Fourier expansion in time (see [14, Chapter XIII] and [6]),

$$\tilde{q}(t) = \sum_{\|\mathbf{k}\| \leq 2N} z^{\mathbf{k}}(\varepsilon t) e^{i(\mathbf{k} \cdot \boldsymbol{\omega})t}, \quad (22)$$

approximating the numerical solution  $q^n$  at  $t = nh$ . We use the notation introduced in (7), where now  $k_\ell = 0$  for  $\ell > M$ , since only the frequencies  $\omega_\ell$  for  $0 \leq \ell \leq M$  appear in the spatial discretization (11).

In our analysis, we must deal with small denominators (see Section 6). To control these terms, we will use non-resonance conditions. As soon as, for a given step size  $h$ , the inequality

$$\left| \sin\left(\frac{h}{2}(\omega_j - \mathbf{k} \cdot \boldsymbol{\omega})\right) \cdot \sin\left(\frac{h}{2}(\omega_j + \mathbf{k} \cdot \boldsymbol{\omega})\right) \right| \geq \varepsilon^{1/2} h^2 (\omega_j + |\mathbf{k} \cdot \boldsymbol{\omega}|) \quad (23)$$

is violated, we have to make an assumption on the pair of indices  $(j, \mathbf{k})$ . For a fixed integer  $N \geq 1$ , subsequently used in the truncation of the expansion (22), the set of near-resonant indices becomes, instead of (8),

$$\mathcal{R}_{\varepsilon, h} = \{(j, \mathbf{k}) : |j| \leq M, \|\mathbf{k}\| \leq 2N, \mathbf{k} \neq \pm(j), \text{ not satisfying (23)}\}.$$

Similar to (9), we require the following non-resonance condition: there are  $\sigma > 0$  and a constant  $C_0$  such that

$$\sup_{(j, \mathbf{k}) \in \mathcal{R}_{\varepsilon, h}} \frac{\omega_j^\sigma}{\boldsymbol{\omega}^{\sigma|\mathbf{k}|}} \varepsilon^{\|\mathbf{k}\|/2} \leq C_0 \varepsilon^N. \quad (24)$$

Notice that, in the limit  $h \rightarrow 0$ , condition (23) becomes equivalent (up to a non-zero constant factor) to  $|\omega_j^2 - (\mathbf{k} \cdot \boldsymbol{\omega})^2| \geq \varepsilon^{1/2} \cdot |\omega_j + |\mathbf{k} \cdot \boldsymbol{\omega}||$ , so that (24) corresponds precisely to the non-resonance condition (9) for the semilinear wave equation.

We assume the further numerical non-resonance condition

$$|\sin(h\omega_j)| \geq h\varepsilon^{1/2} \quad \text{for } |j| \leq M. \quad (25)$$

Yet another non-resonance condition, which leads to improved conservation estimates, reads as follows:

$$\left| \sin\left(\frac{h}{2}(\omega_j - \mathbf{k} \cdot \boldsymbol{\omega})\right) \cdot \sin\left(\frac{h}{2}(\omega_j + \mathbf{k} \cdot \boldsymbol{\omega})\right) \right| \geq c h^2 |\psi(h\omega_j)| \quad (26)$$

for  $(j, \mathbf{k})$  of the form  $j = j_1 + j_2$  and  $\mathbf{k} = \pm \langle j_1 \rangle \pm \langle j_2 \rangle$ ,

with a positive constant  $c > 0$ . In the limit  $h \rightarrow 0$ , this inequality becomes  $|\omega_j^2 - (\mathbf{k} \cdot \boldsymbol{\omega})^2| \geq 4c$  which is automatically fulfilled for the considered pairs  $(j, \mathbf{k})$ .

We are now in the position to state the main result of this paper.

**Theorem 3** *Under the symplecticity condition (20), under the non-resonance conditions (24) with exponent  $\sigma$  and (25)-(26), and under the assumption (15) of small initial data with  $s \geq \sigma + 1$  for  $(q^0, p^0) = (q(0), p(0))$ , the near-conservation estimates*

$$\begin{aligned} \frac{|H_M(q^n, p^n) - H_M(q^0, p^0)|}{\varepsilon^2} &\leq C\varepsilon \\ \frac{|K(q^n, p^n) - K(q^0, p^0)|}{\varepsilon^2} &\leq C(\varepsilon + M^{-s} + \varepsilon t M^{-s+1}) \\ \sum_{\ell=0}^M \omega_\ell^{2s+1} \frac{|I_\ell(q^n, p^n) - I_\ell(q^0, p^0)|}{\varepsilon^2} &\leq C\varepsilon \end{aligned}$$

for energy, momentum and actions hold for long times

$$0 \leq t = nh \leq \varepsilon^{-N+1}$$

with a constant  $C$  which depends on  $s$ ,  $N$ , and  $C_0$ , but is independent of the small parameter  $\varepsilon$ , the dimension  $2M$  of the spatial discretization, the time stepsize  $h$ , and the time  $t = nh$ . If condition (26) fails to be satisfied, then  $C\varepsilon$  is weakened to  $C\varepsilon^{1/2}$  in the above bounds.

In addition we obtain, by the argument of Section 6.2 in [13], that the original Hamiltonian  $H$  of (2) along the trigonometric interpolation polynomials  $(u^n(x), v^n(x))$  with Fourier coefficients  $(q_j^n, p_j^n)$  satisfies the long-time near-conservation estimate

$$\frac{|H(u^n, v^n) - H(u^0, v^0)|}{\varepsilon^2} \leq C\varepsilon \quad \text{for } 0 \leq nh \leq \varepsilon^{-N+1}.$$

For a non-symplectic symmetric method (19) the result remains valid in the transformed variables (21).

The proof of Theorem 3 is given in the subsequent Sections 6 and 7. It is based on the idea of interpolating the numerical solution

by a function where different time scales are well separated. This is done by the ansatz (22) which is a truncated series of products of  $e^{i\omega_j t}$  (oscillations with respect to the fast time  $t$ ) with coefficient functions that are smooth in the slow time  $\tau = \varepsilon t$ . The proof then proceeds as follows:

- Proving existence of smooth functions  $z^{\mathbf{k}}(\tau)$  with derivatives bounded independently of  $\varepsilon$  (on intervals of length  $\varepsilon^{-1}$ ). This is the technically difficult part and elaborated in Section 6. It requires non-resonance conditions and a careful truncation of the series.
- Establishing a Hamiltonian structure and the existence of formal invariants in the differential and algebraic equations for the functions  $z^{\mathbf{k}}(\tau)$  (Sections 7.1 – 7.3).
- Proving closeness (on intervals of length  $\varepsilon^{-1}$ ) of the formal invariants to actions  $I_\ell$ , to the total energy  $H$ , and to the momentum  $K$  (Section 7.4).
- Stretching from short to long intervals of length  $\varepsilon^{-N+1}$  by patching together previous results along an invariant.

## 6 Modulated Fourier expansion

Our principal tool for the long-time analysis of the nonlinearly perturbed wave equation is a short-time modulation expansion constructed in this section. To construct this expansion, we combine the tools and techniques developed in [5], [6], and [13].

### 6.1 Statement of the result

In this section we consider, instead of the symplecticity condition (20), the weaker condition

$$|\psi(h\omega_j)| \leq C |\operatorname{sinc}(h\omega_j)| \quad \text{for } |j| \leq M. \quad (27)$$

In the following result we use the abbreviations (7) and set

$$[[\mathbf{k}]] = \begin{cases} \frac{1}{2}(\|\mathbf{k}\| + 1), & \mathbf{k} \neq 0 \\ \frac{3}{2}, & \mathbf{k} = 0. \end{cases}$$

**Theorem 4** *Under the assumptions of Theorem 3 (with the symplecticity assumption (20) relaxed to (27)), there exist truncated asymptotic expansions (with  $N$  from (24))*

$$\begin{aligned}\tilde{q}(t) &= \sum_{\|\mathbf{k}\| \leq 2N} z^{\mathbf{k}}(\varepsilon t) e^{i(\mathbf{k} \cdot \boldsymbol{\omega})t}, \\ \tilde{p}(t) &= \operatorname{sinc}(h\Omega)^{-1} \frac{\tilde{q}(t+h) - \tilde{q}(t-h)}{2h},\end{aligned}\tag{28}$$

such that the numerical solution  $q^n, p^n$  given by method (19), satisfies

$$\|q^n - \tilde{q}(t)\|_{s+1} + \|p^n - \tilde{p}(t)\|_s \leq C\varepsilon^N \quad \text{for } 0 \leq t = nh \leq \varepsilon^{-1}.\tag{29}$$

The truncated modulated Fourier expansion is bounded by

$$\|\tilde{q}(t)\|_{s+1} + \|\tilde{p}(t)\|_s \leq C\varepsilon \quad \text{for } 0 \leq t \leq \varepsilon^{-1}.\tag{30}$$

On this time interval, we further have, for  $|j| \leq M$ ,

$$\tilde{q}_j(t) = z_j^{(j)}(\varepsilon t) e^{i\omega_j t} + z_j^{-(j)}(\varepsilon t) e^{-i\omega_j t} + r_j, \quad \text{with } \|r\|_{s+1} \leq C\varepsilon^2.\tag{31}$$

(If condition (26) fails to be satisfied, then the bound is  $\|r\|_{s+1} \leq C\varepsilon^{3/2}$ .) The modulation functions  $z^{\mathbf{k}}$  are bounded by

$$\sum_{\|\mathbf{k}\| \leq 2N} \left( \frac{\omega^{|\mathbf{k}|}}{\varepsilon^{|\mathbf{k}|}} \|z^{\mathbf{k}}(\varepsilon t)\|_s \right)^2 \leq C.\tag{32}$$

Bounds of the same type hold for any fixed number of derivatives of  $z^{\mathbf{k}}$  with respect to the slow time  $\tau = \varepsilon t$ . Moreover, the modulation functions satisfy  $z_{-j}^{-\mathbf{k}} = \overline{z_j^{\mathbf{k}}}$ . The constants  $C$  are independent of  $\varepsilon$ ,  $M$ ,  $h$ , and of  $t \leq \varepsilon^{-1}$ .

The proof of this result will cover the remainder of this section. It is organized in the same way as the proof of the analogous result for the analytical solution in [6].

## 6.2 Formal modulation equations

We are looking for a truncated series (28) such that, up to a small defect,

$$\tilde{q}(t+h) - 2\cos(h\Omega)\tilde{q}(t) + \tilde{q}(t-h) = h^2\Psi f(\Phi\tilde{q}(t))$$

with  $\tilde{q}(0) = q^0$ ,  $\tilde{p}(0) = p^0$ , see (17) and (28). We insert the ansatz (28) into this equation, expand the right-hand side into a Taylor series around zero and compare the coefficients of  $e^{i(\mathbf{k}\cdot\boldsymbol{\omega})t}$ . We then get

$$L_j^{\mathbf{k}} z_j^{\mathbf{k}} = -h^2 \psi(h\omega_j) \sum_{m \geq 2} \frac{g^{(m)}(0)}{m!} \times \sum_{\mathbf{k}^1 + \dots + \mathbf{k}^m = \mathbf{k}} \sum'_{j_1 + \dots + j_m \equiv j \pmod{2M}} \phi(h\omega_{j_1}) z_{j_1}^{\mathbf{k}^1} \cdot \dots \cdot \phi(h\omega_{j_m}) z_{j_m}^{\mathbf{k}^m}, \quad (33)$$

where the right-hand side is obtained as in [13]. The prime on the sum over  $j_1, \dots, j_m$  indicates that with every appearance of  $z_{j_i}^{\mathbf{k}^i}$  with  $j_i = \pm M$  a factor  $1/2$  is included. The operator  $L_j^{\mathbf{k}}$  is given as

$$\begin{aligned} (L_j^{\mathbf{k}} z_j^{\mathbf{k}})(\tau) &= e^{ih(\mathbf{k}\cdot\boldsymbol{\omega})} z_j^{\mathbf{k}}(\tau + \varepsilon h) - 2 \cos(h\omega_j) z_j^{\mathbf{k}}(\tau) + e^{-ih(\mathbf{k}\cdot\boldsymbol{\omega})} z_j^{\mathbf{k}}(\tau - \varepsilon h) \\ &= 4s_{\langle j \rangle + \mathbf{k}} s_{\langle j \rangle - \mathbf{k}} z_j^{\mathbf{k}}(\tau) + 2is_{2\mathbf{k}} h \varepsilon \dot{z}_j^{\mathbf{k}}(\tau) + c_{2\mathbf{k}} h^2 \varepsilon^2 \ddot{z}_j^{\mathbf{k}}(\tau) + \dots \end{aligned} \quad (34)$$

Here,  $s_{\mathbf{k}} = \sin(\frac{h}{2} \mathbf{k} \cdot \boldsymbol{\omega})$  and  $c_{\mathbf{k}} = \cos(\frac{h}{2} \mathbf{k} \cdot \boldsymbol{\omega})$ , and the dots on  $z_j^{\mathbf{k}}$  represent derivatives with respect to the slow time  $\tau = \varepsilon t$ . The higher order terms are linear combinations of the  $r$ th derivative of  $z_j^{\mathbf{k}}$  (for  $r \geq 3$ ) multiplied by  $h^r \varepsilon^r$  and containing one of the factors  $s_{2\mathbf{k}}$  or  $c_{2\mathbf{k}}$ .

The first term in (34) vanishes for  $\mathbf{k} = \pm \langle j \rangle$ , so that in this case the dominating term becomes  $\pm 2ih \sin(h\omega_j) \varepsilon \dot{z}_j^{\pm \langle j \rangle}$  due to condition (25). For  $\mathbf{k} \neq \pm \langle j \rangle$  the first term becomes dominant, if the inequality (23) holds. Else, it is not clear which term is dominant, but then the non-resonance condition (24) will ensure that the defect in simply setting  $z_j^{\mathbf{k}} \equiv 0$  is of size  $\mathcal{O}(\varepsilon^{N+1})$  in an appropriate Sobolev-type norm.

In addition, the initial conditions  $\tilde{q}(0) = q^0$  and  $\tilde{p}(0) = p^0$  need to be taken care of. The condition  $\tilde{q}(0) = q^0$  reads

$$\sum_{\|\mathbf{k}\| \leq 2N} z_j^{\mathbf{k}}(0) = q_j^0, \quad (35)$$

and for  $\tilde{p}(0) = p^0$ , we obtain from (28)

$$\frac{1}{2h \operatorname{sinc}(h\omega_j)} \sum_{\|\mathbf{k}\| \leq 2N} \left( z_j^{\mathbf{k}}(\varepsilon h) e^{i(\mathbf{k}\cdot\boldsymbol{\omega})h} - z_j^{\mathbf{k}}(-\varepsilon h) e^{-i(\mathbf{k}\cdot\boldsymbol{\omega})h} \right) = p_j^0. \quad (36)$$



### 6.3 Reverse Picard iteration

We now turn to an iterative construction of the functions  $z_j^{\mathbf{k}}$  such that after  $4N$  iteration steps, the defect in equations (33), (35), and (36) is of size  $\mathcal{O}(\varepsilon^{N+1})$  in the  $H^s$  norm. The iteration procedure we employ can be viewed as a Picard iteration on (33) to (36), where we keep only the dominant terms on the left-hand side. We call it reverse Picard iteration, because the highest appearing derivatives do not carry the new iteration number  $n + 1$ .

Indicating by  $[\cdot]^n$  the  $n$ th iterate of all appearing variables  $z_j^{\mathbf{k}}$  taken within the bracket, we set for  $\mathbf{k} = \pm\langle j \rangle$

$$\begin{aligned} \pm 2ih\varepsilon s_{2j} \left[ \dot{z}_j^{\pm\langle j \rangle} \right]^{n+1} &= \left[ -h^2 \psi(h\omega_j) \sum_{m \geq 2} \frac{g^{(m)}(0)}{m!} \times \right. \\ &\quad \sum_{\mathbf{k}^1 + \dots + \mathbf{k}^m = \mathbf{k}} \sum'_{j_1 + \dots + j_m \equiv j \pmod{2M}} \phi(h\omega_{j_1}) z_{j_1}^{\mathbf{k}^1} \cdot \dots \cdot \phi(h\omega_{j_m}) z_{j_m}^{\mathbf{k}^m} \\ &\quad \left. - \left( c_{2j} h^2 \varepsilon^2 \ddot{z}_j^{\pm\langle j \rangle} + \dots \right) \right]^n \end{aligned} \quad (37)$$

with the sines and cosines  $s_{2j}$  and  $c_{2j}$  defined after formula (34). For  $\mathbf{k} \neq \pm\langle j \rangle$  and  $j$  that are non-resonant with (23), we set

$$\begin{aligned} 4s_{\langle j \rangle + \mathbf{k}} s_{\langle j \rangle - \mathbf{k}} \left[ z_j^{\mathbf{k}} \right]^{n+1} &= \left[ -h^2 \psi(h\omega_j) \sum_{m \geq 2} \frac{g^{(m)}(0)}{m!} \times \right. \\ &\quad \sum_{\mathbf{k}^1 + \dots + \mathbf{k}^m = \mathbf{k}} \sum'_{j_1 + \dots + j_m \equiv j \pmod{2M}} \phi(h\omega_{j_1}) z_{j_1}^{\mathbf{k}^1} \cdot \dots \cdot \phi(h\omega_{j_m}) z_{j_m}^{\mathbf{k}^m} \\ &\quad \left. - \left( 2is_{2\mathbf{k}} h \varepsilon \dot{z}_j^{\mathbf{k}} + c_{2\mathbf{k}} h^2 \varepsilon^2 \ddot{z}_j^{\mathbf{k}} + \dots \right) \right]^n, \end{aligned} \quad (38)$$

whereas we let  $z_j^{\mathbf{k}} = 0$  for  $\mathbf{k} \neq \pm\langle j \rangle$  in the near-resonant set  $\mathcal{R}_{\varepsilon, h}$ . The dots indicate the remainder in (34), truncated after the  $\varepsilon^N$  term.

On the initial conditions we iterate by

$$\left[ z_j^{\langle j \rangle}(0) + z_j^{-\langle j \rangle}(0) \right]^{n+1} = \left[ q_j^0 - \sum_{\mathbf{k} \neq \pm\langle j \rangle} z_j^{\mathbf{k}}(0) \right]^n \quad (39)$$

and on (36) by

$$\begin{aligned} i\omega_j \left[ z_j^{\langle j \rangle}(0) - z_j^{-\langle j \rangle}(0) \right]^{n+1} &= p_j^0 \\ &- \frac{1}{2h \operatorname{sinc}(h\omega_j)} \left[ \sum_{\mathbf{k} \neq \pm\langle j \rangle} z_j^{\mathbf{k}}(0) \left( e^{i(\mathbf{k} \cdot \boldsymbol{\omega})h} - e^{-i(\mathbf{k} \cdot \boldsymbol{\omega})h} \right) \right. \\ &\quad \left. - \sum_{\|\mathbf{k}\| \leq K} \left( (z_j^{\mathbf{k}}(\varepsilon h) - z_j^{\mathbf{k}}(0)) e^{i(\mathbf{k} \cdot \boldsymbol{\omega})h} - (z_j^{\mathbf{k}}(-\varepsilon h) - z_j^{\mathbf{k}}(0)) e^{-i(\mathbf{k} \cdot \boldsymbol{\omega})h} \right) \right]^n. \end{aligned} \quad (40)$$

In all the above formulas, it is tacitly assumed that  $\|\mathbf{k}\| \leq K := 2N$  and  $\|\mathbf{k}^i\| \leq K$  for  $i = 1, \dots, m$ . In each iteration step, we thus have an initial value problem of first-order differential equations for  $z_j^{\pm\langle j \rangle}$  (for  $|j| \leq M$ ) and algebraic equations for  $z_j^{\mathbf{k}}$  with  $\mathbf{k} \neq \pm\langle j \rangle$ .

The starting iterates ( $n = 0$ ) are chosen as  $z_j^{\mathbf{k}}(\tau) = 0$  for  $\mathbf{k} \neq \pm\langle j \rangle$ , and  $z_j^{\pm\langle j \rangle}(\tau) = z_j^{\pm\langle j \rangle}(0)$  with  $z_j^{\pm\langle j \rangle}(0)$  determined from the above formula.

For real initial data we have  $q_{-j}^0 = \overline{q_j^0}$  and  $p_{-j}^0 = \overline{p_j^0}$ , and we observe that the above iteration yields  $[z_{-j}^{-\mathbf{k}}]^n = \overline{[z_j^{\mathbf{k}}]^n}$  for all iterates  $n$  and all  $j, \mathbf{k}$  and hence gives real approximations (28).

#### 6.4 Rescaling and estimation of the nonlinear terms

As in [6], we will work with the more convenient rescaling

$$c_j^{\mathbf{k}} = \frac{\omega^{|\mathbf{k}|}}{\varepsilon^{[\|\mathbf{k}\|]}} z_j^{\mathbf{k}}, \quad \mathbf{c}^{\mathbf{k}} = (c_j^{\mathbf{k}})_{|j| \leq M} = \frac{\omega^{|\mathbf{k}|}}{\varepsilon^{[\|\mathbf{k}\|]}} z^{\mathbf{k}}$$

considered in the space  $\mathbf{H}^s = (H^s)^{\mathcal{K}} = \{\mathbf{c} = (c^{\mathbf{k}})_{\mathbf{k} \in \mathcal{K}} : c^{\mathbf{k}} \in H^s\}$  with norm  $\|\mathbf{c}\|_s^2 = \sum_{\mathbf{k} \in \mathcal{K}} \|c^{\mathbf{k}}\|_s^2$  and where the superscripts  $\mathbf{k}$  are in the set

$$\mathcal{K} = \{\mathbf{k} = (k_\ell)_{\ell=0}^M \text{ with integers } k_\ell : \|\mathbf{k}\| \leq K\}$$

with  $K = 2N$ . The nonlinear function  $\mathbf{f} = (f_j^{\mathbf{k}})$  defined as

$$f_j^{\mathbf{k}}(\mathbf{c}) = \frac{\omega^{|\mathbf{k}|}}{\varepsilon^{[\|\mathbf{k}\|]}} \sum_{m=2}^N \frac{g^{(m)}(0)}{m!} \sum_{\mathbf{k}^1 + \dots + \mathbf{k}^m = \mathbf{k}} \frac{\varepsilon^{[\|\mathbf{k}^1\| + \dots + \|\mathbf{k}^m\|]}}{\omega^{|\mathbf{k}^1| + \dots + |\mathbf{k}^m|}} \times \\ \sum'_{j_1 + \dots + j_m \equiv j \pmod{2M}} \phi(h\omega_{j_1}) c_{j_1}^{\mathbf{k}^1} \cdot \dots \cdot \phi(h\omega_{j_m}) c_{j_m}^{\mathbf{k}^m}$$

expresses the nonlinearity in (33) in the rescaled variables. With the fact that  $H^s$  is a normed algebra, the following bounds are obtained as in [6, Section 3.5] by exploiting the connection between the  $2M$ -periodic sequence  $c_j^{\mathbf{k}}$  and the corresponding trigonometric polynomial (cf. [13]):

$$\sum_{\|\mathbf{k}\| \leq K} \|f^{\mathbf{k}}(\mathbf{c})\|_s^2 \leq \varepsilon P(\|\mathbf{c}\|_s^2) \quad (41)$$

$$\sum_{|j| \leq M} \|f^{\pm\langle j \rangle}(\mathbf{c})\|_s^2 \leq \varepsilon^3 P_1(\|\mathbf{c}\|_s^2), \quad (42)$$

where  $P$  and  $P_1$  are polynomials with coefficients bounded independently of  $\varepsilon, h$ , and  $M$ . Notice that the function  $\phi$  is bounded.

With the different rescaling

$$\hat{c}_j^{\mathbf{k}} = \frac{\omega^{s|\mathbf{k}|}}{\varepsilon^{|\mathbf{k}|}} z_j^{\mathbf{k}}, \quad \hat{\mathbf{c}}^{\mathbf{k}} = (\hat{c}_j^{\mathbf{k}})_{|j| \leq M} = \frac{\omega^{s|\mathbf{k}|}}{\varepsilon^{|\mathbf{k}|}} z^{\mathbf{k}} \quad (43)$$

considered in the space  $\mathbf{H}^1 = (H^1)^{\mathcal{K}}$  with norm  $\|\hat{\mathbf{c}}\|_1^2 = \sum_{\|\mathbf{k}\| \leq K} \|\hat{\mathbf{c}}^{\mathbf{k}}\|_1^2$ , for  $\hat{f}_j^{\mathbf{k}}$  defined as  $f_j^{\mathbf{k}}$  but with  $\omega^{|\mathbf{k}|}$  replaced by  $\omega^{s|\mathbf{k}|}$ , we have similar bounds

$$\begin{aligned} \sum_{\|\mathbf{k}\| \leq K} \|\hat{f}^{\mathbf{k}}(\hat{\mathbf{c}})\|_1^2 &\leq \varepsilon \hat{P}(\|\hat{\mathbf{c}}\|_1^2) \\ \sum_{|j| \leq M} \|\hat{f}^{\pm\langle j \rangle}(\hat{\mathbf{c}})\|_1^2 &\leq \varepsilon^3 \hat{P}_1(\|\hat{\mathbf{c}}\|_1^2) \end{aligned} \quad (44)$$

with other polynomials  $\hat{P}$  and  $\hat{P}_1$ .

### 6.5 Abstract reformulation of the iteration

For  $\mathbf{c} = (c_j^{\mathbf{k}}) \in \mathbf{H}^s$  with  $c_j^{\mathbf{k}} = 0$  for all  $\mathbf{k} \neq \pm\langle j \rangle$  with  $(j, \mathbf{k}) \in \mathcal{R}_{\varepsilon, h}$ , we split the components of  $\mathbf{c}$  corresponding to  $\mathbf{k} = \pm\langle j \rangle$  and  $\mathbf{k} \neq \pm\langle j \rangle$  and collect them in  $\mathbf{a} = (a_j^{\mathbf{k}}) \in \mathbf{H}^s$  and  $\mathbf{b} = (b_j^{\mathbf{k}}) \in \mathbf{H}^s$ , respectively:

$$\begin{aligned} a_j^{\mathbf{k}} &= c_j^{\mathbf{k}} \quad \text{if } \mathbf{k} = \pm\langle j \rangle, \quad \text{and } 0 \text{ else} \\ b_j^{\mathbf{k}} &= c_j^{\mathbf{k}} \quad \text{if (23) is satisfied,} \quad \text{and } 0 \text{ else.} \end{aligned} \quad (45)$$

We then have  $\mathbf{a} + \mathbf{b} = \mathbf{c}$  and  $\|\mathbf{a}\|_s^2 + \|\mathbf{b}\|_s^2 = \|\mathbf{c}\|_s^2$ . We now introduce differential operators  $A, B$  acting on functions  $\mathbf{a}(\tau)$  and  $\mathbf{b}(\tau)$ , respectively:

$$\begin{aligned} (A\mathbf{a})_j^{\pm\langle j \rangle}(\tau) &= \frac{1}{\pm 2i h \varepsilon s_{2j}} \left( c_{2j} h^2 \varepsilon^2 \ddot{a}_j^{\pm\langle j \rangle}(\tau) + \dots \right) \\ (B\mathbf{b})_j^{\mathbf{k}}(\tau) &= \frac{1}{4s_{\langle j \rangle + \mathbf{k} s_{\langle j \rangle - \mathbf{k}}}} \left( 2is_{2\mathbf{k}} h \varepsilon \dot{b}_j^{\mathbf{k}}(\tau) + c_{2\mathbf{k}} h^2 \varepsilon^2 \ddot{b}_j^{\mathbf{k}}(\tau) + \dots \right) \end{aligned}$$

for  $(j, \mathbf{k})$  satisfying (23). These definitions are motivated by formulas (37) and (38), and as in these formulas, the dots represent a truncation after the  $\varepsilon^N$  terms. In terms of the nonlinear function  $\mathbf{f}$  of the preceding subsection, we introduce the functions  $\mathbf{F} = (F_j^{\mathbf{k}})$  and  $\mathbf{G} = (G_j^{\mathbf{k}})$  with non-vanishing entries

$$\begin{aligned} F_j^{\pm\langle j \rangle}(\mathbf{a}, \mathbf{b}) &= -\frac{1}{\pm i \varepsilon} \frac{\psi(h\omega_j)}{\text{sinc}(h\omega_j)} f_j^{\pm\langle j \rangle}(\mathbf{a} + \mathbf{b}), \\ G_j^{\mathbf{k}}(\mathbf{a}, \mathbf{b}) &= -\frac{h^2(\omega_j + |\mathbf{k} \cdot \boldsymbol{\omega}|)}{4s_{\langle j \rangle + \mathbf{k} s_{\langle j \rangle - \mathbf{k}}}} f_j^{\mathbf{k}}(\mathbf{a} + \mathbf{b}) \end{aligned}$$

for  $(j, \mathbf{k})$  satisfying (23). Further we write

$$(\boldsymbol{\Omega}\mathbf{c})_j^{\mathbf{k}} = (\omega_j + |\mathbf{k} \cdot \boldsymbol{\omega}|) c_j^{\mathbf{k}}, \quad (\boldsymbol{\Psi}\mathbf{c})_j^{\mathbf{k}} = \psi(h\omega_j) c_j^{\mathbf{k}}.$$

In terms of  $\mathbf{a}$  and  $\mathbf{b}$ , the iterations (37) and (38) then become of the form

$$\begin{aligned} \dot{\mathbf{a}}^{(n+1)} &= \boldsymbol{\Omega}^{-1}\mathbf{F}(\mathbf{a}^{(n)}, \mathbf{b}^{(n)}) - A\mathbf{a}^{(n)} \\ \mathbf{b}^{(n+1)} &= \boldsymbol{\Omega}^{-1}\boldsymbol{\Psi}\mathbf{G}(\mathbf{a}^{(n)}, \mathbf{b}^{(n)}) - B\mathbf{b}^{(n)}. \end{aligned} \quad (46)$$

By (42), condition (27) gives the bound  $\|\mathbf{F}\|_s \leq C\varepsilon^{1/2}$ , whereas condition (25) yields  $\|\boldsymbol{\Psi}^{-1}\boldsymbol{\Omega}^{-1}\mathbf{F}\|_s \leq C$ . By (41) and (23), we have the bound  $\|\mathbf{G}\|_s \leq C$ . These bounds hold uniformly in  $\varepsilon, h, M$  on bounded subsets of  $\mathbf{H}^s$ . Analogous bounds are obtained for the derivatives of  $\mathbf{F}$  and  $\mathbf{G}$ . The operators  $A$  and  $B$  are estimated as

$$\begin{aligned} \|(A\mathbf{a})(\tau)\|_s &\leq C \sum_{l=2}^N h^{l-2} \varepsilon^{l-3/2} \left\| \frac{d^l}{d\tau^l} \mathbf{a}(\tau) \right\|_s \\ \|(B\mathbf{b})(\tau)\|_s &\leq C\varepsilon^{1/2} \|\dot{\mathbf{b}}(\tau)\|_s + C \sum_{l=2}^N h^{l-2} \varepsilon^{l-1/2} \left\| \frac{d^l}{d\tau^l} \mathbf{b}(\tau) \right\|_s. \end{aligned}$$

The bound for  $A$  is obtained with (25), that for  $B$  uses (23) and the trivial estimate  $|s_{2\mathbf{k}}| = |\sin(h\mathbf{k} \cdot \boldsymbol{\omega})| \leq h|\mathbf{k} \cdot \boldsymbol{\omega}|$ .

The initial value conditions (39) and (40) translate into an equation for  $\mathbf{a}^{(n+1)}$  of the form

$$\mathbf{a}^{(n+1)}(0) = \mathbf{v} + P\mathbf{b}^{(n)}(0) + Q(\mathbf{a} + \mathbf{b})^{(n)}(\varepsilon h) \quad (47)$$

where  $\mathbf{v}$  has the components

$$v_j^{\pm\langle j \rangle} = \frac{\omega_j}{\varepsilon} \left( \frac{1}{2} q_j^0 \mp \frac{i}{2\omega_j} p_j^0 \right).$$

By assumption (15),  $\mathbf{v}$  is bounded in  $\mathbf{H}^s$ . The operators  $P$  and  $Q$  are given by

$$\begin{aligned} (P\mathbf{b})_j^{\pm\langle j \rangle}(0) &= -\frac{\omega_j}{2\varepsilon s_{2j}} \sum_{\mathbf{k} \neq \pm\langle j \rangle} \left( \sin(\omega_j h) \pm \sin((\mathbf{k} \cdot \boldsymbol{\omega}) h) \right) \frac{\varepsilon^{[\mathbf{k}]}}{\omega|\mathbf{k}|} b_j^{\mathbf{k}}(0) \\ (Q\mathbf{c})_j^{\pm\langle j \rangle}(\tau) &= \mp \frac{\omega_j}{4i\varepsilon s_{2j}} \sum_{\|\mathbf{k}\| \leq K} \left( e^{i(\mathbf{k} \cdot \boldsymbol{\omega}) h} \frac{\varepsilon^{[\mathbf{k}]}}{\omega|\mathbf{k}|} \left( c_j^{\mathbf{k}}(\tau) - c_j^{\mathbf{k}}(0) \right) \right. \\ &\quad \left. - e^{-i(\mathbf{k} \cdot \boldsymbol{\omega}) h} \frac{\varepsilon^{[\mathbf{k}]}}{\omega|\mathbf{k}|} \left( c_j^{\mathbf{k}}(-\tau) - c_j^{\mathbf{k}}(0) \right) \right). \end{aligned}$$

For these expressions we have the bounds

$$\begin{aligned} \|(P\mathbf{b})(0)\|_s &\leq C \|\Psi^{-1}\Omega\mathbf{b}(0)\|_s \\ \|(Q\mathbf{c})(\varepsilon h)\|_s &\leq C\varepsilon \sup_{-\varepsilon h < \tau < \varepsilon h} \|\Psi^{-1}\dot{\mathbf{c}}(\tau)\|_s \end{aligned}$$

with a constant  $C$  that is independent of  $\varepsilon$ ,  $h$ , and  $M$ , but depends on  $K = 2N$ . For the first estimate we use  $|\sin(\omega_j h) \pm \sin((\mathbf{k} \cdot \boldsymbol{\omega})h)| \leq h(\omega_j + |\mathbf{k} \cdot \boldsymbol{\omega}|)$ , condition (27), and the Cauchy–Schwarz inequality together with the bound (cf. Lemma 2 of [6])

$$\sum_{\|\mathbf{k}\| \leq K} \omega^{-2|\mathbf{k}|} \leq C < \infty. \quad (48)$$

Similarly, applying the mean value theorem to  $\mathbf{c}(\tau)$  yields the second estimate.

The starting iterates are  $\mathbf{a}^{(0)}(\tau) = \mathbf{v}$  and  $\mathbf{b}^{(0)}(\tau) = 0$ .

### 6.6 Bounds of the modulation functions

In view of the non-resonance conditions (23) and (25), and using the assumption on the filter function (27), we can show by induction that the iterates  $\mathbf{a}^{(n)}$  and  $\mathbf{b}^{(n)}$  and their derivatives with respect to the slow time  $\tau = \varepsilon t$  are bounded in  $\mathbf{H}^s$  for  $0 \leq \tau \leq 1$  and  $n \leq 4N$ : more precisely, the  $(4N)$ -th iterates  $\mathbf{a} = \mathbf{a}^{(4N)}$  and  $\mathbf{b} = \mathbf{b}^{(4N)}$  satisfy

$$\begin{aligned} \|\mathbf{a}(0)\|_s &\leq C, \quad \|\Omega\dot{\mathbf{a}}(\tau)\|_s \leq C\varepsilon^{1/2}, \quad \|\Psi^{-1}\dot{\mathbf{a}}(\tau)\|_s \leq C, \\ \|\Psi^{-1}\Omega\mathbf{b}(\tau)\|_s &\leq C, \end{aligned} \quad (49)$$

with a constant  $C$  independent of  $\varepsilon, h, M$ , but dependent on  $N$ . We also obtain analogous bounds for higher derivatives of  $\mathbf{a}$  and  $\mathbf{b}$  with respect to  $\tau = \varepsilon t$ . For  $z_j^{\mathbf{k}} = \varepsilon^{|\mathbf{k}|} \omega^{-|\mathbf{k}|} c_j^{\mathbf{k}}$  with  $(c_j^{\mathbf{k}}) = \mathbf{c}^{(4N)} = \mathbf{a}^{(4N)} + \mathbf{b}^{(4N)}$ , the bounds for  $\mathbf{a}$  and  $\mathbf{b}$  together yield the bound (32).

These bounds imply  $\|\mathbf{c}(\tau) - \mathbf{a}(0)\|_{s+1} \leq C$  and as in [6, Section 3.7] give, using (48), the bound (30) for  $\tilde{q}(t)$ . For the function  $\tilde{p}(t)$ , defined in (28), we use  $z_j^{\mathbf{k}}(\varepsilon t + \varepsilon h)e^{i(\mathbf{k} \cdot \boldsymbol{\omega})h} - z_j^{\mathbf{k}}(\varepsilon t - \varepsilon h)e^{-i(\mathbf{k} \cdot \boldsymbol{\omega})h} = z_j^{\mathbf{k}}(\varepsilon t)2i \sin((\mathbf{k} \cdot \boldsymbol{\omega})h) + r_j^{\mathbf{k}}$ , where by the mean value theorem  $|r_j^{\mathbf{k}}| \leq 2\varepsilon h \max_{-\varepsilon h < \tau < \varepsilon h} |z_j^{\mathbf{k}}(\tau)|$ . Using the condition (27), the bounds (49) yield in a similar way also the statement (30) for the function  $\tilde{p}(t)$ .

Using (42) and (46) we also obtain the bound, for  $\mathbf{b} = \mathbf{b}^{(4N)}$ ,

$$\left( \sum_{\|\mathbf{k}\|=1} \|(\Psi^{-1}\Omega\mathbf{b})^{\mathbf{k}}\|_s^2 \right)^{1/2} \leq C\varepsilon.$$

Moreover, condition (26) ensures that

$$\sum_{|j| \leq M} \sum_{j_1 + j_2 = j} \sum_{\mathbf{k} = \pm \langle j_1 \rangle \pm \langle j_2 \rangle} \omega_j^{2(s+1)} |b_j^{\mathbf{k}}|^2 \leq C\varepsilon.$$

These bounds together with (49) yield (31).

With the alternative scaling (43) we obtain the same bounds (for  $\tau = \varepsilon t \leq 1$ ),

$$\|\widehat{\mathbf{a}}(0)\|_1 \leq C, \quad \|\Omega \dot{\widehat{\mathbf{a}}}(\tau)\|_1 \leq C\varepsilon^{1/2}, \quad \|\Psi^{-1} \Omega \widehat{\mathbf{b}}(\tau)\|_1 \leq C. \quad (50)$$

and again

$$\left( \sum_{\|\mathbf{k}\|=1} \|(\Psi^{-1} \Omega \widehat{\mathbf{b}})^{\mathbf{k}}\|_1^2 \right)^{1/2} \leq C\varepsilon. \quad (51)$$

For the function  $\widehat{\mathbf{a}}(\tau)$  these statements follow at once from the fact that  $\|\widehat{\mathbf{a}}^{\mathbf{k}}\|_1 = \|\mathbf{a}^{\mathbf{k}}\|_s$ . For the function  $\widehat{\mathbf{b}}(\tau)$  one has to repeat the argumentation from before, but one needs no longer take care of initial values.

In addition to these bounds, we also obtain that the map

$$B_\varepsilon \subset H^{s+1} \times H^s \rightarrow \mathbf{H}^1 : (u(0), v(0)) \mapsto \widehat{\mathbf{c}}(0)$$

(with  $B_\varepsilon$  the ball of radius  $\varepsilon$  centered at 0) is Lipschitz continuous with a Lipschitz constant proportional to  $\varepsilon^{-1}$ : at  $t = 0$ ,

$$\|\widehat{\mathbf{a}}_2 - \widehat{\mathbf{a}}_1\|_1^2 + \|\Omega(\widehat{\mathbf{b}}_2 - \widehat{\mathbf{b}}_1)\|_1^2 \leq \frac{C}{\varepsilon^2} \left( \|u_2 - u_1\|_{s+1}^2 + \|v_2 - v_1\|_s^2 \right). \quad (52)$$

### 6.7 Defects

We consider the defect  $\delta(t) = (\delta_j(t))_{|j| \leq M}$  in (17) divided by  $h^2 \psi(h\omega_j)$ :

$$\delta_j(t) = \frac{\widetilde{q}_j(t+h) - 2 \cos(h\omega_j) \widetilde{q}_j(t) + \widetilde{q}_j(t-h)}{h^2 \psi(h\omega_j)} - f_j(\Phi \widetilde{q}(t))$$

where  $f = (f_j)$  is given in (11) and the approximation  $\widetilde{q}(t) = (q_j(t))$  is given by (28) with  $z_j^{\mathbf{k}} = (z_j^{\mathbf{k}})^{(4N)}$  obtained after  $4N$  iterations of the procedure in Section 6.3. We write this defect as

$$\delta(t) = \sum_{\|\mathbf{k}\| \leq NK} d^{\mathbf{k}}(\varepsilon t) e^{i(\mathbf{k} \cdot \boldsymbol{\omega})t} + R_{N+1}(\widetilde{q})(t).$$

Here we have set

$$d_j^{\mathbf{k}} = \frac{1}{h^2 \psi(h\omega_j)} \tilde{L}_j^{\mathbf{k}} z_j^{\mathbf{k}} + \sum_{m=2}^N \frac{g^{(m)}(0)}{m!} \times \quad (53)$$

$$\sum_{\mathbf{k}^1 + \dots + \mathbf{k}^m = \mathbf{k}} \sum'_{j_1 + \dots + j_m \equiv j \pmod{2M}} \phi(h\omega_{j_1}) z_{j_1}^{\mathbf{k}^1} \cdot \dots \cdot \phi(h\omega_{j_m}) z_{j_m}^{\mathbf{k}^m},$$

which is to be considered for  $\|\mathbf{k}\| \leq NK$ , and where we set  $z_j^{\mathbf{k}} = 0$  for  $\|\mathbf{k}\| > K = 2N$ . The operator  $\tilde{L}_j^{\mathbf{k}}$  denotes the truncation of the expansion (34) after the  $\varepsilon^N$  term. The function  $R_{N+1}$  collects the remainder term of the Taylor expansion of  $f$  after  $N$  terms, and that due to the truncation of the series in (34) after the  $\varepsilon^N$  term. Using the bound (30) for the remainder in the Taylor expansion of  $f$  and the estimates (49) for the  $(N+1)$ -th derivative for  $z_j^{\mathbf{k}}(\tau)$ , we have  $\|R_{N+1}(\tilde{q})\|_{s+1} \leq C\varepsilon^{N+1}$ .

We now use the bound of [6, Section 3.8] to obtain

$$\left\| \sum_{\|\mathbf{k}\| \leq NK} d^{\mathbf{k}}(\varepsilon t) e^{i(\mathbf{k} \cdot \boldsymbol{\omega})t} \right\|_s^2 \leq C \sum_{\|\mathbf{k}\| \leq NK} \left\| \boldsymbol{\omega}^{|\mathbf{k}|} d^{\mathbf{k}}(\varepsilon t) \right\|_s^2. \quad (54)$$

In the following two subsections we estimate the right-hand side of (54) by  $C\varepsilon^{2(N+1)}$ .

### 6.8 Defect in the truncated and near-resonant modes

For  $\|\mathbf{k}\| > K = 2N$  (truncated modes) and for  $(j, \mathbf{k})$  in the set  $\mathcal{R}_{\varepsilon, h}$  of near-resonances we have by definition  $z_j^{\mathbf{k}} = 0$ . In both situations the defect reads

$$d_j^{\mathbf{k}} = \sum_{m=2}^N \frac{g^{(m)}(0)}{m!} \times$$

$$\sum_{\mathbf{k}^1 + \dots + \mathbf{k}^m = \mathbf{k}} \sum'_{j_1 + \dots + j_m \equiv j \pmod{2M}} \phi(h\omega_{j_1}) z_{j_1}^{\mathbf{k}^1} \cdot \dots \cdot \phi(h\omega_{j_m}) z_{j_m}^{\mathbf{k}^m}.$$

For truncated modes we write the defect as  $d_j^{\mathbf{k}} = \varepsilon^{[\mathbf{k}]} \boldsymbol{\omega}^{-|\mathbf{k}|} f_j^{\mathbf{k}}$ , and we notice that by (49) and (41), used with  $NK$  in place of  $K$ , the bound  $\|\mathbf{f}\|_s^2 \leq C\varepsilon$  holds. We thus have

$$\sum_{\|\mathbf{k}\| > K} \sum'_{|j| \leq M} \omega_j^{2s} |\boldsymbol{\omega}^{|\mathbf{k}|} d_j^{\mathbf{k}}|^2 = \sum_{\|\mathbf{k}\| > K} \sum'_{|j| \leq M} \omega_j^{2s} |f_j^{\mathbf{k}}|^2 \varepsilon^{2[\mathbf{k}]}$$

and hence, since  $2\llbracket \mathbf{k} \rrbracket = \|\mathbf{k}\| + 1 \geq K + 2 = 2(N + 1)$ ,

$$\sum_{\|\mathbf{k}\| > K} \sum_{|j| \leq M}' \omega_j^{2s} |\omega^{|\mathbf{k}|} d_j^{\mathbf{k}}|^2 \leq C \varepsilon^{2(N+1)}.$$

For the near-resonant modes we consider the rescaling (43), so that  $d_j^{\mathbf{k}} = \varepsilon^{\llbracket \mathbf{k} \rrbracket} \omega^{-s|\mathbf{k}|} \widehat{f}_j^{\mathbf{k}}$ . We have  $\|\widehat{\mathbf{f}}\|_1^2 \leq C\varepsilon$  by (50) and (44), so that

$$\begin{aligned} \sum_{(j, \mathbf{k}) \in \mathcal{R}_{\varepsilon, h}} \omega_j^{2s} |\omega^{|\mathbf{k}|} d_j^{\mathbf{k}}|^2 &= \sum_{(j, \mathbf{k}) \in \mathcal{R}_{\varepsilon, h}} \frac{\omega_j^{2(s-1)}}{\omega^{2(s-1)|\mathbf{k}|}} \varepsilon^{2\llbracket \mathbf{k} \rrbracket} \omega_j^2 |\widehat{f}_j^{\mathbf{k}}|^2 \\ &\leq C \sup_{(j, \mathbf{k}) \in \mathcal{R}_{\varepsilon, h}} \frac{\omega_j^{2(s-1)} \varepsilon^{2\llbracket \mathbf{k} \rrbracket + 1}}{\omega^{2(s-1)|\mathbf{k}|}}. \end{aligned}$$

The non-resonance condition (24) is formulated such that the supremum is bounded by  $C_0^2 \varepsilon^{2(N+1)}$ , and hence

$$\sum_{(j, \mathbf{k}) \in \mathcal{R}_{\varepsilon, h}} \omega_j^{2s} |\omega^{|\mathbf{k}|} d_j^{\mathbf{k}}|^2 \leq C \varepsilon^{2(N+1)}. \quad (55)$$

### 6.9 Defect in the non-resonant modes

We now assume that  $\|\mathbf{k}\| \leq K$  and that  $(j, \mathbf{k})$  satisfies the non-resonance condition (23), so that in the scaled variables  $c_j^{\mathbf{k}}$  of Section 6.4 the defect satisfies

$$\omega^{|\mathbf{k}|} d_j^{\mathbf{k}} = \varepsilon^{\llbracket \mathbf{k} \rrbracket} \left( \frac{1}{h^2 \psi(h\omega_j)} \widetilde{L}_j^{\mathbf{k}} c_j^{\mathbf{k}} + f_j^{\mathbf{k}}(\mathbf{c}) \right).$$

Written in terms of the components  $\mathbf{a}$  and  $\mathbf{b}$  of (45) we have

$$\begin{aligned} \omega_j d_j^{\pm\langle j \rangle} &= \varepsilon \left( \pm 2i\varepsilon \omega_j \frac{\text{sinc}(h\omega_j)}{\psi(h\omega_j)} (\dot{a}_j^{\pm\langle j \rangle} + (A\mathbf{a})_j^{\pm\langle j \rangle}) + f_j^{\pm\langle j \rangle}(\mathbf{a} + \mathbf{b}) \right) \\ \omega^{|\mathbf{k}|} d_j^{\mathbf{k}} &= \varepsilon^{\llbracket \mathbf{k} \rrbracket} \left( \frac{4^{s(j) + \mathbf{k}^s(j) - \mathbf{k}}}{h^2 \psi(h\omega_j)} (b_j^{\mathbf{k}} + (B\mathbf{b})_j^{\mathbf{k}}) + f_j^{\mathbf{k}}(\mathbf{a} + \mathbf{b}) \right). \end{aligned}$$

It should be noted that the functions in this defect are actually the  $4N$ -th iterates  $\mathbf{a}^{(4N)}$  and  $\mathbf{b}^{(4N)}$  of the iteration in Section 6.3. Expressing  $f_j^{\pm\langle j \rangle}(\mathbf{a} + \mathbf{b})$  and  $f_j^{\mathbf{k}}(\mathbf{a} + \mathbf{b})$  in terms of  $\mathbf{F}(\mathbf{a}, \mathbf{b})$  and  $\mathbf{G}(\mathbf{a}, \mathbf{b})$



and inserting  $\mathbf{F}$  and  $\mathbf{G}$  from (46) into this defect, relates it to the increment of the iteration in the following way:

$$\begin{aligned}\omega_j d_j^{\pm\langle j \rangle} &= 2\omega_j \alpha_j^{\pm\langle j \rangle} \left( [\dot{a}_j^{\pm\langle j \rangle}]^{(4N)} - [\dot{a}_j^{\pm\langle j \rangle}]^{(4N+1)} \right), \\ \alpha_j^{\pm\langle j \rangle} &:= \pm i \varepsilon^2 \frac{\text{sinc}(h\omega_j)}{\psi(h\omega_j)}, \\ \omega^{|\mathbf{k}|} d_j^{\mathbf{k}} &= \beta_j^{\mathbf{k}} \left( [b_j^{\mathbf{k}}]^{(4N)} - [b_j^{\mathbf{k}}]^{(4N+1)} \right), \\ \beta_j^{\mathbf{k}} &:= \varepsilon^{|\mathbf{k}|} \frac{4s_{\langle j \rangle} + \mathbf{k} s_{\langle j \rangle} - \mathbf{k}}{h^2 \psi(h\omega_j)}.\end{aligned}$$

Motivated by these relations we introduce new variables

$$\tilde{a}_j^{\pm\langle j \rangle} := \alpha_j^{\pm\langle j \rangle} a_j^{\pm\langle j \rangle}, \quad \tilde{b}_j^{\mathbf{k}} := \beta_j^{\mathbf{k}} b_j^{\mathbf{k}}. \quad (56)$$

Collecting these variables into vectors and using the transformed functions

$$\begin{aligned}\tilde{F}_j^{\pm\langle j \rangle}(\tilde{\mathbf{a}}, \tilde{\mathbf{b}}) &:= \alpha_j^{\pm\langle j \rangle} F_j^{\pm\langle j \rangle}(\boldsymbol{\alpha}^{-1}\tilde{\mathbf{a}}, \boldsymbol{\beta}^{-1}\tilde{\mathbf{b}}) \\ &= -\varepsilon f_j^{\pm\langle j \rangle}(\boldsymbol{\alpha}^{-1}\tilde{\mathbf{a}} + \boldsymbol{\beta}^{-1}\tilde{\mathbf{b}}) \\ \tilde{G}_j^{\mathbf{k}}(\tilde{\mathbf{a}}, \tilde{\mathbf{b}}) &:= \frac{\beta_j^{\mathbf{k}} \psi(h\omega_j)}{\omega_j + |\mathbf{k} \cdot \boldsymbol{\omega}|} G_j^{\mathbf{k}}(\boldsymbol{\alpha}^{-1}\tilde{\mathbf{a}}, \boldsymbol{\beta}^{-1}\tilde{\mathbf{b}}) \\ &= -\varepsilon^{|\mathbf{k}|} f_j^{\mathbf{k}}(\boldsymbol{\alpha}^{-1}\tilde{\mathbf{a}} + \boldsymbol{\beta}^{-1}\tilde{\mathbf{b}})\end{aligned}$$

the iteration (46)-(47) becomes

$$\begin{aligned}\dot{\tilde{\mathbf{a}}}^{(n+1)} &= \boldsymbol{\Omega}^{-1} \tilde{\mathbf{F}}(\tilde{\mathbf{a}}, \tilde{\mathbf{b}}) - A \tilde{\mathbf{a}}^{(n)} \\ \tilde{\mathbf{b}}^{(n+1)} &= \tilde{\mathbf{G}}(\tilde{\mathbf{a}}, \tilde{\mathbf{b}}) - B \tilde{\mathbf{b}}^{(n)} \\ \tilde{\mathbf{a}}^{(n+1)}(0) &= \boldsymbol{\alpha} \mathbf{v} + \tilde{P} \tilde{\mathbf{b}}^{(n)}(0) + Q \tilde{\mathbf{a}}^{(n)}(\varepsilon h) + \tilde{Q} \tilde{\mathbf{b}}^{(n)}(\varepsilon h).\end{aligned} \quad (57)$$

In the iteration for the initial values we abbreviate  $\tilde{P} = \boldsymbol{\alpha} P \boldsymbol{\beta}^{-1}$ ,  $\tilde{Q} = \boldsymbol{\alpha} Q \boldsymbol{\beta}^{-1}$ , which are bounded by

$$\begin{aligned}\|(\tilde{P}\tilde{\mathbf{b}})(0)\|_s &\leq C \varepsilon^{1/2} \|\tilde{\mathbf{b}}(0)\|_s \\ \|(Q\tilde{\mathbf{a}})(\varepsilon h)\|_s &\leq C \varepsilon^{1/2} \sup_{-\varepsilon h < \tau < \varepsilon h} \|\dot{\tilde{\mathbf{a}}}(\tau)\|_s \\ \|(\tilde{Q}\tilde{\mathbf{b}})(\varepsilon h)\|_s &\leq C \varepsilon^{3/2} \sup_{-\varepsilon h < \tau < \varepsilon h} \|\boldsymbol{\Omega}^{-1} \dot{\tilde{\mathbf{b}}}(\tau)\|_s.\end{aligned}$$

In an  $\mathbf{H}^s$  neighbourhood of 0 where the bounds (49) hold, the partial derivatives of  $\tilde{\mathbf{F}}$  with respect to  $\tilde{\mathbf{a}}$  and  $\tilde{\mathbf{b}}$  and those of  $\tilde{\mathbf{G}}$  with respect to

$\tilde{\mathbf{b}}$  are bounded by  $\mathcal{O}(\varepsilon^{1/2})$ , whereas the derivatives of  $\tilde{\mathbf{G}}$  with respect to  $\tilde{\mathbf{a}}$  is only  $\mathcal{O}(1)$ . This is the same situation as we had for the exact solution in [6]. As in that paper one proves

$$\begin{aligned} \|\|\| \Omega(\dot{\tilde{\mathbf{a}}}^{(4N+1)}(\tau) - \dot{\tilde{\mathbf{a}}}^{(4N)}(\tau)) \|\|\|_s &\leq C \varepsilon^{N+2} \\ \|\|\| \tilde{\mathbf{b}}^{(4N+1)}(\tau) - \tilde{\mathbf{b}}^{(4N)}(\tau) \|\|\|_s &\leq C \varepsilon^{N+2} \\ \|\|\| \tilde{\mathbf{a}}^{(4N+1)}(0) - \tilde{\mathbf{a}}^{(4N)}(0) \|\|\|_s &\leq C \varepsilon^{N+2}. \end{aligned}$$

These estimates yield the desired bound of the defect in the non-resonant modes  $(j, \mathbf{k}) \notin \mathcal{R}_{\varepsilon, h}$ . Combined with the corresponding estimates of Subsection 6.8 we obtain

$$\left( \sum_{\|\mathbf{k}\| \leq K} \|\omega^{|\mathbf{k}|} d^{\mathbf{k}}(\tau)\|_s^2 \right)^{1/2} \leq C \varepsilon^{N+1} \quad \text{for } \tau \leq 1. \quad (58)$$

Consequently, the defect  $\delta(t)$  (see Subsection 6.7) satisfies

$$\|\Omega^{-1} \delta(t)\|_{s+1} = \|\delta(t)\|_s \leq C \varepsilon^{N+1} \quad \text{for } t \leq \varepsilon^{-1}. \quad (59)$$

For the defect in the initial conditions (35) and (36) we obtain

$$\|\tilde{q}(0) - q^0\|_{s+1} + \|\tilde{p}(0) - p^0\|_s \leq C \varepsilon^{N+1}.$$

For the alternative scaling  $\hat{c}_j^{\mathbf{k}} = \omega^{s|\mathbf{k}|} z_j^{\mathbf{k}}$ , we obtain

$$\left( \sum_{\|\mathbf{k}\| \leq K} \|\omega^{s|\mathbf{k}|} d^{\mathbf{k}}(\tau)\|_1^2 \right)^{1/2} \leq C \varepsilon^{N+1} \quad \text{for } \tau \leq 1. \quad (60)$$

### 6.10 Remainder term of the modulated Fourier expansion

We write the method (19) in the form

$$\begin{aligned} \begin{pmatrix} q^{n+1} \\ \Omega^{-1} p^{n+1} \end{pmatrix} &= \begin{pmatrix} \cos(h\Omega) & \sin(h\Omega) \\ -\sin(h\Omega) & \cos(h\Omega) \end{pmatrix} \begin{pmatrix} q^n \\ \Omega^{-1} p^n \end{pmatrix} \\ &\quad + \frac{h}{2} \Psi_1 \begin{pmatrix} \sin(h\Omega) f^n \\ \cos(h\Omega) f^n + f^{n+1} \end{pmatrix} \end{aligned}$$

where  $f^n = \Omega^{-1} f(\Phi q^n)$ , and we notice that  $\Psi_1$  is a matrix, bounded independently of  $h$  and the dimension  $M$ . The differences  $\Delta q^n := \tilde{q}(t_n) - q^n$  and  $\Delta p^n := \tilde{p}(t_n) - p^n$ , where  $t_n := nh$ , satisfy the same relation with  $f^n$  replaced by  $\Omega^{-1}(f(\Phi \tilde{q}(t_n)) - f(\Phi q^n)) + \delta(t_n)$ . Using

the Lipschitz bound (cf. Section 4.2 in [13] on the relation between  $f(q)$  and  $g(u)$  of (1))

$$\begin{aligned} \|\Omega^{-1}(f(q_1) - f(q_2))\|_{s+1} &= \|f(q_1) - f(q_2)\|_s \\ &\leq C\varepsilon\|q_1 - q_2\|_s \leq C\varepsilon\|q_1 - q_2\|_{s+1} \end{aligned}$$

for  $q_1, q_2 \in H^s$  satisfying  $\|q_i\|_s \leq M\varepsilon$ , and the estimate (59) for the defect yields

$$\begin{aligned} \left\| \begin{pmatrix} \Delta q^{n+1} \\ \Omega^{-1} \Delta p^{n+1} \end{pmatrix} \right\|_{s+1} &\leq \left\| \begin{pmatrix} \Delta q^n \\ \Omega^{-1} \Delta p^n \end{pmatrix} \right\|_{s+1} \\ &\quad + \frac{h}{2} (C\varepsilon\|\Delta q^n\|_{s+1} + C\varepsilon\|\Delta q^{n+1}\|_{s+1} + C\varepsilon^{N+1}). \end{aligned}$$

Solving this inequality gives the estimate

$$\|\Delta q^n\|_{s+1} + \|\Omega^{-1} \Delta p^n\|_{s+1} \leq C(1 + t_n)\varepsilon^{N+1} \quad \text{for } t_n \leq \varepsilon^{-1}$$

and thus completes the proof of Theorem 4.

## 7 Conservation properties

We now show that the system of equations determining the modulation functions has almost-invariants close to the actions, the momentum and the total energy along numerical solutions given by the full discretization (17)–(18). The proof takes up arguments of [6] for the conservation of actions, of [13] for the conservation of momentum and aspects of the space discretization, and of [14, Ch. XIII] for the conservation of energy and for the aspects arising from the time discretization.

### 7.1 The extended potential

The defect formula (53) can be rewritten as

$$\frac{1}{h^2\psi(h\omega_j)} \tilde{L}_j^{\mathbf{k}} z_j^{\mathbf{k}} + \nabla_{-j}^{-\mathbf{k}} \mathcal{U}(\Phi_{\mathbf{z}}) = d_j^{\mathbf{k}}, \quad (61)$$

where  $\nabla_{-j}^{-\mathbf{k}} \mathcal{U}(\mathbf{y})$  is the partial derivative with respect to  $y_{-j}^{-\mathbf{k}}$  of the *extended potential* (see [13])

$$\mathcal{U}(\mathbf{y}) = \sum_{l=-N}^N \mathcal{U}_l(\mathbf{y}) \quad (62)$$

$$\mathcal{U}_l(\mathbf{y}) = \sum_{m=2}^N \frac{U^{(m+1)}(0)}{(m+1)!} \sum_{\mathbf{k}^1 + \dots + \mathbf{k}^{m+1} = \mathbf{0}} \sum_{j_1 + \dots + j_{m+1} = 2Ml} y_{j_1}^{\mathbf{k}^1} \dots y_{j_{m+1}}^{\mathbf{k}^{m+1}},$$

where again  $\|\mathbf{k}^i\| \leq 2N$  and  $|j_i| \leq M$ , and  $U(u)$  is the potential in (2).

### 7.2 Invariance under group actions

The existence of almost-invariants for the system (61) turns out to be a consequence, in the spirit of Noether's theorem, of the invariance of the extended potential under continuous group actions: for an arbitrary real sequence  $\boldsymbol{\mu} = (\mu_\ell)_{\ell \geq 0}$  and for  $\theta \in \mathbb{R}$ , let

$$\begin{aligned} S_{\boldsymbol{\mu}}(\theta)\mathbf{y} &= \left( e^{i(\mathbf{k} \cdot \boldsymbol{\mu})\theta} y_j^{\mathbf{k}} \right)_{|j| \leq M, \|\mathbf{k}\| \leq K}, \\ T(\theta)\mathbf{y} &= \left( e^{ij\theta} y_j^{\mathbf{k}} \right)_{|j| \leq M, \|\mathbf{k}\| \leq K}. \end{aligned} \quad (63)$$

Since the sum in the definition of  $\mathcal{U}$  is over  $\mathbf{k}^1 + \dots + \mathbf{k}^{m+1} = \mathbf{0}$  and that in  $\mathcal{U}_0$  over  $j_1 + \dots + j_{m+1} = 0$ , we have

$$\mathcal{U}(S_{\boldsymbol{\mu}}(\theta)\mathbf{y}) = \mathcal{U}(\mathbf{y}), \quad \mathcal{U}_0(T(\theta)\mathbf{y}) = \mathcal{U}_0(\mathbf{y}) \quad \text{for } \theta \in \mathbb{R}.$$

Differentiating these relations with respect to  $\theta$  yields

$$\begin{aligned} 0 &= \left. \frac{d}{d\theta} \right|_{\theta=0} \mathcal{U}(S_{\boldsymbol{\mu}}(\theta)\mathbf{y}) = \sum_{\|\mathbf{k}\| \leq K} \sum'_{|j| \leq M} i(\mathbf{k} \cdot \boldsymbol{\mu}) y_j^{\mathbf{k}} \nabla_j^{\mathbf{k}} \mathcal{U}(\mathbf{y}) \\ 0 &= \left. \frac{d}{d\theta} \right|_{\theta=0} \mathcal{U}_0(T(\theta)\mathbf{y}) = \sum_{\|\mathbf{k}\| \leq K} \sum'_{|j| \leq M} ij y_j^{\mathbf{k}} \nabla_j^{\mathbf{k}} \mathcal{U}_0(\mathbf{y}). \end{aligned} \quad (64)$$

### 7.3 Almost-invariants of the modulation system

We now multiply (61) once with  $i(\mathbf{k} \cdot \boldsymbol{\mu})\phi(h\omega_j)z_{-j}^{-\mathbf{k}}$  and once with  $ij\phi(h\omega_j)z_{-j}^{-\mathbf{k}}$ , and sum over  $j$  and  $\mathbf{k}$  with  $|j| \leq M$  and  $\|\mathbf{k}\| \leq K$ . Thanks to (64), we obtain

$$\begin{aligned} \sum_{\|\mathbf{k}\| \leq K} \sum'_{|j| \leq M} i(\mathbf{k} \cdot \boldsymbol{\mu}) \frac{\phi(h\omega_j)}{h^2 \psi(h\omega_j)} z_{-j}^{-\mathbf{k}} \tilde{L}_j^{\mathbf{k}} z_j^{\mathbf{k}} \\ = \sum_{\|\mathbf{k}\| \leq K} \sum'_{|j| \leq M} i(\mathbf{k} \cdot \boldsymbol{\mu}) \phi(h\omega_j) z_{-j}^{-\mathbf{k}} d_j^{\mathbf{k}}, \end{aligned} \quad (65)$$

$$\begin{aligned} \sum_{\|\mathbf{k}\| \leq K} \sum'_{|j| \leq M} ij \frac{\phi(h\omega_j)}{h^2 \psi(h\omega_j)} z_{-j}^{-\mathbf{k}} \tilde{L}_j^{\mathbf{k}} z_j^{\mathbf{k}} \\ = \sum_{\|\mathbf{k}\| \leq K} \sum'_{|j| \leq M} ij \phi(h\omega_j) z_{-j}^{-\mathbf{k}} \left( d_j^{\mathbf{k}} - \sum_{l \neq 0} \nabla_{-j}^{-\mathbf{k}} \mathcal{U}_l(\Phi \mathbf{z}) \right). \end{aligned} \quad (66)$$

By the expansion (34) of the operator  $\tilde{L}_j^{\mathbf{k}}$ , only expressions of the following type appear for  $z(\tau) = z_j^{\mathbf{k}}(\tau)$  and  $\bar{z}(\tau) = z_{-j}^{-\mathbf{k}}(\tau)$  on the left-hand side of the above equations:

$$\begin{aligned} \operatorname{Re} \bar{z} z^{(2l+1)} &= \operatorname{Re} \frac{d}{d\tau} \left( \bar{z} z^{(2l)} - \dots \pm \bar{z}^{(l-1)} z^{(l+1)} \mp \frac{1}{2} \bar{z}^{(l)} z^{(l)} \right) \\ \operatorname{Im} \bar{z} z^{(2l+2)} &= \operatorname{Im} \frac{d}{d\tau} \left( \bar{z} z^{(2l+1)} - \dot{\bar{z}} z^{(2l)} + \dots \pm \bar{z}^{(l)} z^{(l+1)} \right). \end{aligned} \quad (67)$$

Therefore, the left-hand sides can be written as total derivatives of functions  $\varepsilon \mathcal{J}_{\boldsymbol{\mu}}[\mathbf{z}](\tau)$  and  $\varepsilon \mathcal{K}[\mathbf{z}](\tau)$  which depend on  $\mathbf{z}(\tau)$  and its derivatives  $\varepsilon^\ell \mathbf{z}^{(\ell)}(\tau)$  for  $\ell = 1, \dots, N-1$ . In this way, (65) and (66) become

$$-\varepsilon \frac{d}{d\tau} \mathcal{J}_{\boldsymbol{\mu}}[\mathbf{z}] = \sum_{\|\mathbf{k}\| \leq K} \sum'_{|j| \leq M} i(\mathbf{k} \cdot \boldsymbol{\mu}) \phi(h\omega_j) z_{-j}^{-\mathbf{k}} d_j^{\mathbf{k}} \quad (68)$$

$$-\varepsilon \frac{d}{d\tau} \mathcal{K}[\mathbf{z}] = \sum_{\|\mathbf{k}\| \leq K} \sum'_{|j| \leq M} i j \phi(h\omega_j) z_{-j}^{-\mathbf{k}} \left( d_j^{\mathbf{k}} - \sum_{l \neq 0} \nabla_{-j}^{-\mathbf{k}} \mathcal{U}_l(\Phi \mathbf{z}) \right). \quad (69)$$

In the following we consider the special case of  $\boldsymbol{\mu} = \langle \ell \rangle$ . From the smallness of the right-hand sides in (68) and (69) we infer the following.

**Theorem 5** *Under the conditions of Theorem 4 we have, for  $\tau \leq 1$ ,*

$$\begin{aligned} \sum_{\ell=0}^M \omega_\ell^{2s+1} \left| \frac{d}{d\tau} \mathcal{J}_{\langle \ell \rangle}[\mathbf{z}](\tau) \right| &\leq C \varepsilon^{N+1}, \\ \left| \frac{d}{d\tau} \mathcal{K}[\mathbf{z}](\tau) \right| &\leq C (\varepsilon^{N+1} + \varepsilon^2 M^{-s+1}). \end{aligned}$$

*Proof* The result is obtained from (68) and (69) with the arguments of [6, 13] as follows. With the bounds (50) and (60), the estimate for the functions  $\mathcal{J}_{\langle \ell \rangle}[\mathbf{z}]$  follows with the proof of Theorem 3 in [6]. With the bound (32) and with the bounds  $\|\mathbf{z}\|_1 \leq C\varepsilon$  and  $\|\mathbf{d}\|_0 \leq C\varepsilon^{N+1}$ , which follow from (32) and (58), the estimate for  $\mathcal{K}[\mathbf{z}]$  is obtained as in Theorem 5.2 of [13].  $\square$

A further almost-invariant is obtained by multiplying (61) with the expression  $\phi(h\omega_j) (i(\mathbf{k} \cdot \boldsymbol{\omega}) z_{-j}^{-\mathbf{k}} + \varepsilon \dot{z}_{-j}^{-\mathbf{k}})$ , summing over  $j$  and  $\mathbf{k}$ , and using (64):

$$\begin{aligned} \sum_{\|\mathbf{k}\| \leq K} \sum'_{|j| \leq M} \frac{\phi(h\omega_j)}{h^2 \psi(h\omega_j)} (i(\mathbf{k} \cdot \boldsymbol{\omega}) z_{-j}^{-\mathbf{k}} + \varepsilon \dot{z}_{-j}^{-\mathbf{k}}) \tilde{L}_j^{\mathbf{k}} z_j^{\mathbf{k}} \\ + \varepsilon \frac{d}{d\tau} \mathcal{U}(\Phi \mathbf{z}) = \sum_{\|\mathbf{k}\| \leq K} \sum'_{|j| \leq M} \phi(h\omega_j) (i(\mathbf{k} \cdot \boldsymbol{\omega}) z_{-j}^{-\mathbf{k}} + \varepsilon \dot{z}_{-j}^{-\mathbf{k}}) d_j^{\mathbf{k}}. \end{aligned} \quad (70)$$

In addition to the identities (67) we also use

$$\begin{aligned}\operatorname{Re} \dot{\bar{z}}z^{(2l)} &= \operatorname{Re} \frac{d}{dt} \left( \dot{\bar{z}}z^{(2l-1)} - \dots \mp \bar{z}^{(l-1)}z^{(l+1)} \pm \frac{1}{2}\bar{z}^{(l)}z^{(l)} \right) \\ \operatorname{Im} \dot{\bar{z}}z^{(2l+1)} &= \operatorname{Im} \frac{d}{dt} \left( \dot{\bar{z}}z^{(2l)} - \ddot{\bar{z}}z^{(2l-1)} + \dots \mp \bar{z}^{(l)}z^{(l+1)} \right).\end{aligned}$$

Therefore, the left-hand side of (70) can be written as the total derivative of a function  $\varepsilon\mathcal{H}[\mathbf{z}](\tau)$ , so that (70) becomes

$$\varepsilon \frac{d}{d\tau} \mathcal{H}[\mathbf{z}] = \sum_{\|\mathbf{k}\| \leq K} \sum'_{|j| \leq M} \phi(h\omega_j) \left( i(\mathbf{k} \cdot \boldsymbol{\omega})z_{-j}^{-\mathbf{k}} + \varepsilon \dot{z}_{-j}^{-\mathbf{k}} \right) d_j^{\mathbf{k}}. \quad (71)$$

As in Theorem 5, the Cauchy-Schwarz inequality and the estimates for  $z_j^{\mathbf{k}}$  and  $d_j^{\mathbf{k}}$  then yield the following estimate.

**Theorem 6** *Under the conditions of Theorem 4 we have, for  $\tau \leq 1$ ,*

$$\left| \frac{d}{d\tau} \mathcal{H}[\mathbf{z}](\tau) \right| \leq C \varepsilon^{N+1}.$$

#### 7.4 Relationship with actions, momentum, and energy

We now show that the almost-invariant  $\mathcal{J}_{(\ell)}$  of the modulated Fourier expansion is close to the corresponding harmonic action (13) of the numerical solution,

$$J_\ell = I_\ell + I_{-\ell} = 2I_\ell \quad \text{for } 0 < \ell < M, \quad J_0 = I_0, \quad J_M = I_M,$$

and that  $\mathcal{H}$  and  $\mathcal{K}$  are close to the Hamiltonian  $H_M$  and the momentum  $K$  of (12) and (13), respectively.

**Theorem 7** *Under the conditions of Theorem 3, along the numerical solution  $(q^n, p^n)$  of (19) and the associated modulation sequence  $\mathbf{z}(\varepsilon t)$ , it holds that*

$$\begin{aligned}\mathcal{H}[\mathbf{z}](\varepsilon t_n) &= H_M(q^n, p^n) + \mathcal{O}(\varepsilon^3) \\ \mathcal{K}[\mathbf{z}](\varepsilon t_n) &= K(q^n, p^n) + \mathcal{O}(\varepsilon^3) + \mathcal{O}(\varepsilon^2 M^{-s}) \\ \mathcal{J}_{(\ell)}[\mathbf{z}](\varepsilon t_n) &= J_\ell(q^n, p^n) + \gamma_\ell(t_n) \varepsilon^3\end{aligned}$$

with  $\sum_{\ell=0}^M \omega_\ell^{2s+1} \gamma_\ell(t_n) \leq C$  for  $t_n \leq \varepsilon^{-1}$ . All appearing constants are independent of  $\varepsilon$ ,  $M$ ,  $h$ , and  $n$ .

*Proof* With the identities (67) we obtain from (66) that

$$\mathcal{K}[\mathbf{z}] = \sum_{\|\mathbf{k}\| \leq K} \sum'_{|j| \leq M} j \frac{\phi(h\omega_j)}{\psi(h\omega_j)} \left( (\mathbf{k} \cdot \boldsymbol{\omega}) \operatorname{sinc}(h \mathbf{k} \cdot \boldsymbol{\omega}) |z_j^{\mathbf{k}}|^2 + 2\varepsilon c_{2\mathbf{k}} \operatorname{Im}(z_{-j}^{-\mathbf{k}} z_j^{\mathbf{k}}) + \dots \right).$$

Separating the terms with  $\mathbf{k} = \pm \langle j \rangle$  and using the symplecticity condition (20), and applying the bounds (50) and (51) to the remaining terms, we find

$$\mathcal{K}[\mathbf{z}] = \sum'_{|j| \leq M} j \omega_j \left( |z_j^{\langle j \rangle}|^2 - |z_j^{-\langle j \rangle}|^2 \right) + \mathcal{O}(\varepsilon^3).$$

In terms of the Fourier coefficients of the modulated Fourier expansion  $\tilde{q}_j(t) = \sum_{\|\mathbf{k}\| \leq K} z_j^{\mathbf{k}}(\varepsilon t) e^{i(\mathbf{k} \cdot \boldsymbol{\omega})t}$ , we have at  $t = t_n$

$$\begin{aligned} \mathcal{K}[\mathbf{z}] &= \sum'_{|j| \leq M} j \frac{\omega_j}{4} \left( |\tilde{q}_j + (i\omega_j)^{-1} \tilde{p}_j|^2 - |\tilde{q}_j - (i\omega_j)^{-1} \tilde{p}_j|^2 \right) + \mathcal{O}(\varepsilon^3) \\ &= K(\tilde{q}, \tilde{p}) + \mathcal{O}(\varepsilon^3) + \mathcal{O}(\varepsilon^2 M^{-s}) \\ &= K(q^n, p^n) + \mathcal{O}(\varepsilon^3) + \mathcal{O}(\varepsilon^2 M^{-s}), \end{aligned}$$

where we have used (31). The  $\mathcal{O}(\varepsilon^2 M^{-s})$  terms come from the boundary terms in the sum. The last equality is a consequence of the remainder bound of Theorem 4.

Similarly, we obtain from (70) that

$$\mathcal{H}[\mathbf{z}] = \sum_{\|\mathbf{k}\| \leq K} \sum'_{|j| \leq M} (\mathbf{k} \cdot \boldsymbol{\omega}) \frac{\phi(h\omega_j)}{\psi(h\omega_j)} \left( (\mathbf{k} \cdot \boldsymbol{\omega}) \operatorname{sinc}(h \mathbf{k} \cdot \boldsymbol{\omega}) |z_j^{\mathbf{k}}|^2 + \dots \right) + \mathcal{U}(\Phi \mathbf{z}),$$

which yields, using in addition  $\mathcal{U}(\Phi \mathbf{z}) = \mathcal{O}(\varepsilon^3)$ ,

$$\mathcal{H}[\mathbf{z}] = \sum'_{|j| \leq M} \omega_j^2 \left( |z_j^{\langle j \rangle}|^2 + |z_j^{-\langle j \rangle}|^2 \right) + \mathcal{O}(\varepsilon^3),$$

and shows that  $\mathcal{H}[\mathbf{z}] = H_M(q^n, p^n) + \mathcal{O}(\varepsilon^3)$ .

The result for  $J_\ell$  is obtained in the same way, using in addition Lemma 3 of [6] to estimate the remainder terms.  $\square$

With an identical argument to that of [6, Section 4.5], Theorems 5 and 6 together with the estimates of Theorem 4 and the Lipschitz continuity (52) yield the statement of Theorem 3 by patching together many intervals of length  $\varepsilon^{-1}$ .

## 8 The Störmer–Verlet/leapfrog discretization

The leapfrog discretization of (11) reads, in the two-step formulation,

$$q^{n+1} - 2q^n + q^{n-1} = h^2(-\Omega^2 q^n + f(q^n)), \quad (72)$$

with the velocity approximation  $p^n$  given by

$$2hp^n = q^{n+1} - q^{n-1}. \quad (73)$$

The starting value is chosen as  $q^1 = q^0 + hp^0 + \frac{h^2}{2} f(q^0)$ . Conservation properties of this method will be obtained by reinterpreting it as a trigonometric method (17) with modified frequencies  $\widehat{\omega}_j$  satisfying  $1 - \frac{1}{2}h^2\omega_j^2 = \cos(h\widehat{\omega}_j)$ , that is,

$$\sin\left(\frac{1}{2}h\widehat{\omega}_j\right) = \frac{1}{2}h\omega_j. \quad (74)$$

This is possible as long as  $h\omega_j \leq 2$ .

**Theorem 8** *Under the stepsize restriction  $h\omega_M \leq c < 2$ , under the non-resonance conditions (24) and (26) for the modified frequencies  $\widehat{\omega}_j$  of (74), and under the assumption (15) of small initial data with  $s \geq \sigma + 1$  for  $(q^0, p^0) = (q(0), p(0))$ , the near-conservation estimates*

$$\begin{aligned} \frac{|H_M(q^n, p^n) - H_M(q^0, p^0)|}{\varepsilon^2} &\leq C(\varepsilon + h^2) \\ \frac{|K(q^n, p^n) - K(q^0, p^0)|}{\varepsilon^2} &\leq C(\varepsilon + h^2 + M^{-s} + \varepsilon t M^{-s+1}) \\ \sum_{\ell=0}^M \omega_\ell^{2s-1} \frac{|I_\ell(q^n, p^n) - I_\ell(q^0, p^0)|}{\varepsilon^2} &\leq C(\varepsilon + h^2) \end{aligned}$$

for energy, momentum and actions hold for long times

$$0 \leq t = nh \leq \varepsilon^{-N+1}$$

with a constant  $C$  which depends on  $s$ ,  $N$ ,  $C_0$ , and  $c$ , but is independent of  $\varepsilon$ ,  $M$ ,  $h$ , and  $t$ .

*Proof* Denoting by  $\widehat{\Omega}$  the diagonal matrix with entries  $\widehat{\omega}_j$ , we introduce the transformed variables

$$\widehat{q}^n = \text{sinc}(h\widehat{\Omega})q^n, \quad \widehat{p}^n = p^n,$$

which are solutions to the symplectic trigonometric method (17)-(18) with  $\psi = \text{sinc}$  and  $\phi = 1$ . Under the stepsize restriction  $h\omega_M \leq c < 2$



the non-resonance condition (25) is trivially satisfied for  $\widehat{\omega}_j$ , and we have

$$\omega_j \leq \widehat{\omega}_j \leq C\omega_j,$$

where  $C$  depends only on  $c$ . Hence, the assumption (15) of small initial data is satisfied with the same exponent  $s$  for the weighted norms defined with  $\widehat{\omega}_j$  or  $\omega_j$ . We can therefore apply Theorem 3 in the transformed variables  $(\widehat{q}^n, \widehat{p}^n)$ . With the estimate  $|\text{sinc}(h\widehat{\omega}_j) - 1| \leq \frac{1}{6}h^2\widehat{\omega}_j^2$ , the result stated for the original variables  $(q^n, p^n)$  then follows.  $\square$

We apply the leapfrog method to the problem of Section 4 with stepsize  $h = 0.3$ , so that the CFL number  $h\omega_M \approx 1.92$  is close to the linear stability limit. In Figure 4 we observe oscillations with large relative amplitude proportional to  $h^2\omega_j^2$  for the actions  $I_j$  corresponding to high frequencies, but no drift in actions, energy, and momentum.

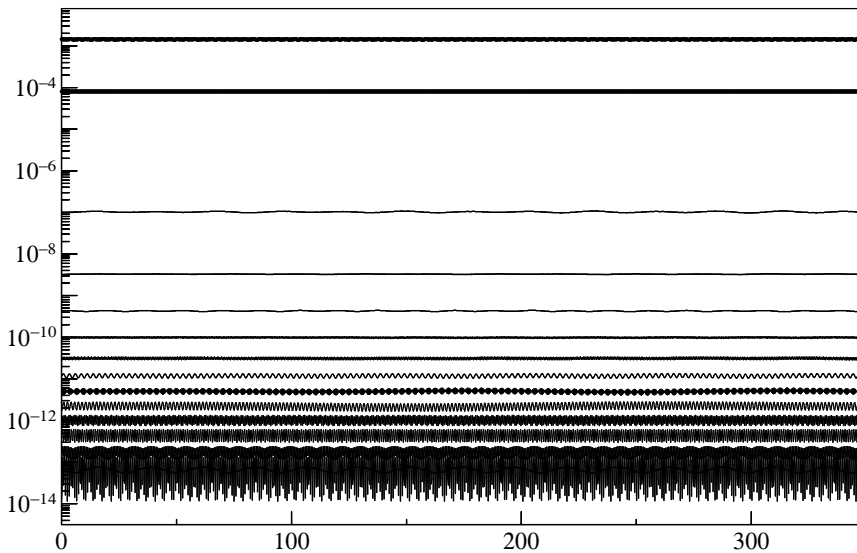


Fig. 4. Actions, energy, and momentum along the numerical solution of the leapfrog method, every 5th action is plotted.

### Acknowledgement

This work was partially supported by the Fonds National Suisse, project No. 200020-113249/1, and by DFG, Project LU 532/4-1. A large part of this work was carried out when the authors visited the Isaac Newton Institute in Cambridge.

## References

1. D. Bambusi, *Birkhoff normal form for some nonlinear PDEs*, Comm. Math. Phys. **234** (2003), 253–285.
2. G. Benettin and A. Giorgilli, *On the Hamiltonian interpolation of near to the identity symplectic mappings with application to symplectic integration algorithms*, J. Statist. Phys. **74** (1994), 1117–1143.
3. J. Bourgain, *Construction of approximative and almost periodic solutions of perturbed linear Schrödinger and wave equations*, Geom. Funct. Anal. **6** (1996), 201–230.
4. B. Cano, *Conserved quantities of some Hamiltonian wave equations after full discretization*, Numer. Math. **103** (2006), 197–223.
5. D. Cohen, E. Hairer, and C. Lubich, *Numerical energy conservation for multi-frequency oscillatory differential equations*, BIT **45** (2005), 287–305.
6. ———, *Long-time analysis of nonlinearly perturbed wave equations via modulated Fourier expansions*, Arch. Ration. Mech. Anal. **187** (2008), 341–368.
7. P. Deuffhard, *A study of extrapolation methods based on multistep schemes without parasitic solutions*, Z. Angew. Math. Phys. **30** (1979), 177–189.
8. G. Dujardin and E. Faou, *Normal form and long time analysis of splitting schemes for the linear Schrödinger equation with small potential*, Numer. Math. **108** (2007), 223–262.
9. B. García-Archilla, J. M. Sanz-Serna, and R. D. Skeel, *Long-time-step methods for oscillatory differential equations*, SIAM J. Sci. Comput. **20** (1999), 930–963.
10. H. Grubmüller, H. Heller, A. Windemuth, and P. Tavan, *Generalized Verlet algorithm for efficient molecular dynamics simulations with long-range interactions*, Mol. Sim. **6** (1991), 121–142.
11. E. Hairer and C. Lubich, *The life-span of backward error analysis for numerical integrators*, Numer. Math. **76** (1997), 441–462, Erratum: <http://www.unige.ch/math/folks/hairer/>.
12. ———, *Long-time energy conservation of numerical methods for oscillatory differential equations*, SIAM J. Numer. Anal. **38** (2001), 414–441.
13. ———, *Spectral semi-discretisations of nonlinear wave equations over long times*, Foundations of Comput. Math. (2008).
14. E. Hairer, C. Lubich, and G. Wanner, *Geometric numerical integration. Structure-preserving algorithms for ordinary differential equations*, 2nd ed., Springer Series in Computational Mathematics 31, Springer-Verlag, Berlin, 2006.
15. E. Hairer, S. P. Nørsett, and G. Wanner, *Solving ordinary differential equations I. Nonstiff problems*, 2nd ed., Springer Series in Computational Mathematics 8, Springer, Berlin, 1993.
16. S. Reich, *Backward error analysis for numerical integrators*, SIAM J. Numer. Anal. **36** (1999), 1549–1570.
17. M. Tuckerman, B. J. Berne, and G. J. Martyna, *Reversible multiple time scale molecular dynamics*, J. Chem. Phys. **97** (1992), 1990–2001.