

## A Stroboscopic Numerical Method for Highly Oscillatory Problems

Maripaz Calvo, Philippe Chartier, Ander Murua, Jesus Maria Sanz-Serna

► **To cite this version:**

Maripaz Calvo, Philippe Chartier, Ander Murua, Jesus Maria Sanz-Serna. A Stroboscopic Numerical Method for Highly Oscillatory Problems. Engquist Björn, Runborg Olof, Tsai Yen-Hsi R. Numerical analysis of multiscale computations: proceedings of a winter workshop at the Banff International Research Station 2009, 82, Springer, pp.71-85, 2012, 978-3-642-21942-9. 10.1007/978-3-642-21943-6 . hal-00777182

**HAL Id: hal-00777182**

**<https://hal.inria.fr/hal-00777182>**

Submitted on 17 Jan 2013

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# A Stroboscopic Numerical Method for Highly Oscillatory Problems

Mari Paz Calvo, Philippe Chartier, Ander Murua, and Jesús María Sanz-Serna

**Abstract** We suggest a method for the integration of highly oscillatory systems with a single high frequency. The new method may be seen as a purely numerical way of implementing the analytical technique of stroboscopic averaging. The technique may be easily implemented in combination with standard software and may be applied with variable step sizes. Numerical experiments show that the suggested algorithms may be substantially more efficient than standard numerical integrators.

## 1 Introduction

We suggest a numerical method for the integration of highly oscillatory differential equations  $dy/dt = f(y,t)$  with a single high frequency  $2\pi/\varepsilon$ ,  $\varepsilon \ll 1$ . The new method may be seen as a purely numerical way of implementing the analytical technique of *stroboscopic averaging* [13] which constructs an averaged differential system  $dY/dt = F(Y)$  whose solutions  $Y$  (approximately) interpolate the sought highly oscillatory solution  $y$  at times  $t = t_0 + 2\pi\varepsilon n$ , ( $n$  integer). In the spirit of the heterogeneous multiscale methods (see [6, 5, 8, 16, 7, 1], cf. [14, 3]), we integrate numeri-

---

M. P. Calvo

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

P. Chartier

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

A. Murua

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

J. M. Sanz-Serna

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

cally the averaged system without using the analytic expression of  $F$ ; all information on  $F$  required by the algorithm is gathered on the fly by numerically integrating the original system in small time windows. The technique may be easily implemented in combination with standard software and may be applied with *variable step sizes*.

Section 2, based on [4], presents the theoretical foundation of the algorithm. Sect. 3 contains a description of the new method along with a brief discussion of related literature. Examples of oscillatory systems that may be treated with our approach are provided in Sect. 4 and the final section presents numerical examples. It is found that the suggested algorithms may be substantially more efficient than standard numerical integrators.

## 2 A Modified Equation Approach to Averaging

We wish to integrate numerically initial value problems for differential systems of the form

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

where  $y$  is a  $D$ -dimensional real vector,  $\varepsilon$  is a small parameter and the smooth function  $f$  is assumed to depend  $2\pi$ -periodically on the variable  $t/\varepsilon$ . Our interest is in situations where, as  $\varepsilon \rightarrow 0$ , the solutions or some of their derivatives with respect to  $t$  become *unbounded*; relevant examples will be presented in Sect. 4.

If we denote by  $\varphi_{t_0,t;\varepsilon} : \mathcal{R}^D \rightarrow \mathcal{R}^D$  the solution operator of (1), so that

$$y(t) = \varphi_{t_0,t;\varepsilon}(y_0)$$

is the solution that satisfies the initial condition  $y(t_0) = y_0$ , then the *one-period* map  $\Psi_{t_0;\varepsilon} = \varphi_{t_0,t_0+2\pi\varepsilon;\varepsilon}$  depends on  $t_0$  in a  $2\pi\varepsilon$ -periodic manner; this is proved by noting that both  $\varphi_{t_0,t;\varepsilon}(y_0)$  and  $\varphi_{t_0+2\pi\varepsilon,t+2\pi\varepsilon;\varepsilon}(y_0)$  satisfy the same initial value problem

$$\frac{d}{dt}y(t) = f\left(y(t), \frac{t}{\varepsilon}; \varepsilon\right) = f\left(y(t), \frac{t+2\pi\varepsilon}{\varepsilon}; \varepsilon\right), \quad y(t_0) = y_0.$$

It follows that, at the *stroboscopic times*  $t_n = t_0 + 2\pi\varepsilon n$ ,  $n = 0, \pm 1, \pm 2, \dots$ ,

$$y(t_n) = \varphi_{t_0,t_n;\varepsilon}(y_0) = \varphi_{t_{n-1},t_n;\varepsilon}(\varphi_{t_0,t_{n-1};\varepsilon}(y_0)) = \varphi_{t_0,t_0+2\pi\varepsilon;\varepsilon}(\varphi_{t_0,t_{n-1};\varepsilon}(y_0))$$

and, hence, we arrive at the fundamental formula:

$$y(t_n) = (\Psi_{t_0;\varepsilon})^n(y_0), \quad n = 0, \pm 1, \pm 2, \dots \quad (2)$$

For the problems we are interested in (see Sect. 4) there is an expansion

$$\Psi_{t_0;\varepsilon}(y_0) = y_0 + \sum_{j=1}^{\infty} \varepsilon^j M_j(y_0), \quad (3)$$

with suitable smooth maps  $M_j : \mathcal{R}^D \rightarrow \mathcal{R}^D$  independent of  $\varepsilon$ , and thus  $\Psi_{t_0;\varepsilon}$  is a smooth *near-to-identity map*. Standard results from the backward error analysis of numerical integrators [15, 9] show then the existence of an *autonomous system* (the modified system of  $\Psi_{t_0;\varepsilon}$ )

$$\frac{d}{dt}Y = F(Y; \varepsilon) = F_1(Y) + \varepsilon F_2(Y) + \varepsilon^2 F_3(Y) + \dots \quad (4)$$

whose (formal) solutions satisfy that  $Y(t_n) = \Psi_{t_0;\varepsilon}(Y(t_{n-1}))$  for  $n = 0, \pm 1, \pm 2, \dots$  so that

$$Y(t_n) = (\Psi_{t_0;\varepsilon})^n(Y_0), \quad n = 0, \pm 1, \pm 2, \dots \quad (5)$$

( $F$  and the  $F_j$  depend on  $t_0$ —because  $\Psi_{t_0;\varepsilon}$  does—, but this dependence has not been incorporated to the notation.) We conclude from (2) and (5) that, if one chooses  $Y(t_0) = y(t_0)$ , then  $Y(t)$  *exactly coincides with*  $y(t)$  *at the stroboscopic times*  $t_n = t_0 + 2\pi\varepsilon n$ . In this way it is possible in principle to find  $y(t_n)$  by solving the system (4), *where all  $t$ -derivatives of  $Y$  remain bounded as  $\varepsilon \rightarrow 0$* . Furthermore  $y$  may be recovered from  $Y$  even at values of  $t$  that do not coincide with one of the stroboscopic times. In fact,

$$y(t) = (\varphi_{t_n,t;\varepsilon} \circ \Phi_{t_n-t;\varepsilon})(Y(t)), \quad (6)$$

where  $t_n$  is the largest stroboscopic time  $\leq t$  and  $\Phi_{\cdot;\varepsilon}$  denotes the flow of (4). In this way,  $y$  is ‘enslaved’ to  $Y$  through the mapping  $\varphi_{t_n,t;\varepsilon} \circ \Phi_{t_n-t;\varepsilon}$  whose dependence on  $t$  is easily seen to be  $2\pi\varepsilon$ -periodic.

For future reference we note that an alternative way of writing (5) is

$$\Psi_{t_0;\varepsilon}^n \equiv \Phi_{2\pi\varepsilon n;\varepsilon}; \quad (7)$$

after a whole number  $n$  of periods the solution operator  $\Psi_{t_0;\varepsilon}^n = \varphi_{t_0,t_0+2\pi\varepsilon n}$  of the non-autonomous system (1) coincides with the flow of the autonomous (4).

It is well known that the series (4) does not converge in general, and in order to get rigorous results one has to consider a truncated version ( $J \geq 1$  is an arbitrarily large integer)

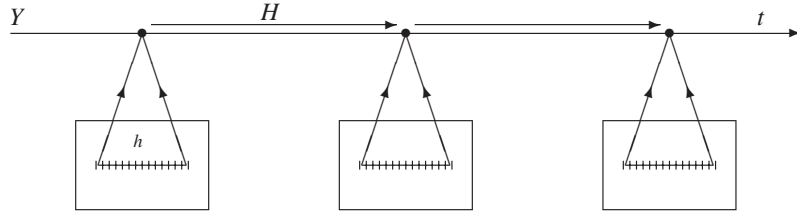
$$\frac{d}{dt}Y = F^{(J)}(Y; \varepsilon) = F_1(Y) + \varepsilon F_2(Y) + \varepsilon^2 F_3(Y) + \dots + \varepsilon^{J-1} F_J(Y), \quad (8)$$

whose solutions satisfy that  $Y(t_n) - \Psi_{t_0;\varepsilon}(Y(t_{n-1})) = \mathcal{O}(\varepsilon^{J+1})$ . If  $Y$  solves (8) with  $Y(t_0) = y(t_0)$  then  $Y(t_n)$  and  $y(t_n)$  differ by an  $\mathcal{O}(\varepsilon^J)$  amount, where the constant implied in the  $\mathcal{O}$  notation is uniform as the stroboscopic time  $t_n$  ranges in a time interval  $t_0 \leq t_n \leq t_0 + T$  with  $T = \mathcal{O}(1)$  as  $\varepsilon \rightarrow 0$ .

The process of obtaining the autonomous system (4) (or (8)) from the original system (1) is referred to in the averaging literature [13] as high-order stroboscopic averaging. As a rule, the amount of work required to find analytically the functions  $F_j$  is formidable, even when the interest is limited to lowest values of  $j$ .

### 3 A Numerical Method

In this section we propose a purely numerical method that bypasses the need for finding analytically the functions  $F_j$ . To simplify the exposition, we will ignore hereafter the  $\mathcal{O}(\varepsilon^J)$  remainder that arises from truncating (4), i.e. we will proceed as if the series (4) were convergent. Since  $J$  may be chosen arbitrarily large, the disregarded truncation errors are, as  $\varepsilon \rightarrow 0$ , negligible when compared with other errors present in the method to be described.



**Fig. 1** Schematic view of the numerical integration. The  $t$ -axis above represents the macro-integration of the averaged system with (large) macro-steps  $H$ . Whenever the macro-solver requires information on the averaged system, the algorithm carries out a micro-integration of the original problem in a small time-window. The micro-step size  $h$  is small with respect to  $\varepsilon$

In order to integrate the highly oscillatory system (1) with initial condition  $y(t_0) = y_0$ , we (approximately) compute the corresponding smooth interpolant  $Y(t)$ , i.e. the solution of the initial value problem specified by the *averaged system* (4) along with the initial condition  $Y(t_0) = y_0$ . We integrate (4) by a standard numerical method, the so-called *macro-solver*, with a macro-step  $H$  that ideally should be substantially larger than the small period  $2\pi\varepsilon$ . In the spirit of heterogeneous multi-scale methods, the information on  $F$  required by the macro-solver is gathered on the fly by integrating, with a micro-step  $h$ , the original system (1) in time-windows of length  $\mathcal{O}(\varepsilon)$ . These auxiliary integrations are also performed by means of a standard numerical method, the *micro-solver*, see Fig. 1. (It is not necessary that the choices of macro and micro-solver coincide.)

If the macro-solver is a linear multistep or Runge-Kutta method, then the only information on the system (4) required by the solver are function values  $F(Y^*; \varepsilon)$  at given values of the argument  $Y^*$ . Since, by definition,  $\Phi_{t;\varepsilon}$  is the flow of (4) we may write

$$F(Y^*; \varepsilon) = \left. \frac{d}{dt} \Phi_{t;\varepsilon}(Y^*) \right|_{t=0},$$

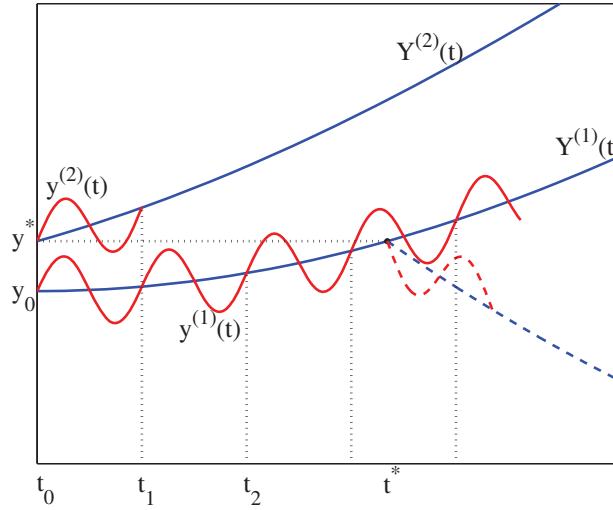
or, after approximating the time-derivative by central differences,

$$F(Y^*; \varepsilon) = \frac{1}{2\delta} [\Phi_{\delta;\varepsilon}(Y^*) - \Phi_{-\delta;\varepsilon}(Y^*)] + \mathcal{O}(\delta^2).$$

We now set  $\delta = 2\pi\varepsilon$  and use the identity (7) to get

$$F(Y^*; \varepsilon) = \frac{1}{4\pi\varepsilon} [\Psi_{t_0; \varepsilon}(Y^*) - \Psi_{t_0; \varepsilon}^{-1}(Y^*)] + \mathcal{O}(\varepsilon^2), \quad (9)$$

a formula that may be used to compute approximately  $F(Y^*; \varepsilon)$  since  $\Psi_{t_0; \varepsilon}(Y^*)$  and  $\Psi_{t_0; \varepsilon}^{-1}(Y^*)$  may be found numerically through micro-integrations. In fact one has to integrate (1) with initial condition  $y(t_0) = Y^*$ , first from  $t = t_0$  to  $t = t_0 + 2\pi\varepsilon$  and then from  $t = t_0$  to  $t = t_0 - 2\pi\varepsilon$ .



**Fig. 2** The wiggly solid lines represent the solutions  $y^{(1)}(t)$  and  $y^{(2)}(t)$  of the oscillatory problem with initial conditions  $y^{(1)}(t_0) = y_0$  and  $y^{(2)}(t_0) = y^*$ . We have also represented the solutions of the averaged system with  $Y^{(1)}(t_0) = y_0$  and  $Y^{(2)}(t_0) = y^*$ ; the graphs of  $Y^{(1)}(t)$  and  $Y^{(2)}(t)$  are translates along the time-axis of one another because the averaged system is autonomous. At stroboscopic times each oscillatory solution  $y^{(i)}(t)$  coincides with the corresponding averaged solution  $Y^{(i)}(t)$ . Now assume that we are computing numerically  $Y^{(1)}$ , that the macro-solver has reached the point  $(t^*, y^*)$  ( $t^*$  is not a stroboscopic time) and that it requires the value of the slope  $F(y^*; \varepsilon)$ . The correct procedure is based on the fact that the slope of  $Y^{(1)}(t)$  at  $(t^*, y^*)$  coincides with the slope of  $Y^{(2)}(t)$  at  $(t_0, y^*)$ ; micro-integrations on the intervals  $t_0 \leq t \leq t_0 + 2\pi\varepsilon$  and  $t_0 \geq t \geq t_0 - 2\pi\varepsilon$  (this is not shown in the figure) are performed to find  $y^{(2)}(t_0 \pm 2\pi\varepsilon) = Y^{(2)}(t_0 \pm 2\pi\varepsilon)$  and the values  $Y^{(2)}(t_0 \pm 2\pi\varepsilon)$  are then used to find the slope by means of finite differences. Micro-integrating in the intervals  $t^* \leq t \leq t^* + 2\pi\varepsilon$  and  $t^* \geq t \geq t^* - 2\pi\varepsilon$  will not do: the averaged system depends on  $t_0$ —see Sect. 2—and such micro-integrations (discontinuous wiggly lines) would provide information on a solution (discontinuous line without wiggles) of the wrong averaged system.

Some important remarks are in order. The initial condition for each micro-integration is *always* prescribed at  $t = t_0$ , regardless of the point of the time axis the macro-solver may have reached when the micro-integration is performed. We have

tried to make this fact apparent in Fig. 1 by enclosing different micro-integrations in boxes that are not connected by a common time-axis (cf. Fig. 1.1 in [8] or Fig. 2 in [16]). All micro-integrations find solutions of (1) in the interval  $[t_0 - 2\pi\varepsilon, t_0 + 2\pi\varepsilon]$ . With the terminology of [3] we may say that the algorithm suggested here is *asynchronous*. Fig. 2 may be of assistance in understanding the situation. This figure should also make it clear that it is not at all necessary that the step-points used by the macro-integrator be stroboscopic times; this is a particularly valuable feature if the macro-solver employs variable steps. We also emphasize that if the macro-solver outputs (an approximation to) the averaged solution  $Y$  at a stroboscopic time  $t_n$ , then the output is an approximation to  $y(t_n)$ ; if output occurs at a non-stroboscopic value of  $t$  it is still possible to recover an approximation to  $y(t)$  by using (6).

Of course, other difference formulae may also be used instead of (9). For instance, we may approximate  $F(Y^*; \varepsilon)$  with an  $\mathcal{O}(\varepsilon^4)$  error by means of

$$\begin{aligned} & \frac{1}{24\pi\varepsilon} \left( -\Phi_{4\pi\varepsilon;\varepsilon}(Y^*) + 8\Phi_{2\pi\varepsilon;\varepsilon}(Y^*) - 8\Phi_{-\pi\varepsilon;\varepsilon}(Y^*) + \Phi_{-4\pi\varepsilon;\varepsilon}(Y^*) \right) \quad (10) \\ & = \frac{1}{24\pi\varepsilon} \left( -\Psi_{t_0;\varepsilon}^2(Y^*) + 8\Psi_{t_0;\varepsilon}(Y^*) - 8\Psi_{t_0;\varepsilon}^{-1}(Y^*) + \Psi_{t_0;\varepsilon}^{-2}(Y^*) \right). \end{aligned}$$

Now the integrations to be carried out to find  $\Psi_{t_0;\varepsilon}^2(Y^*) = \varphi_{t_0, t_0+4\pi\varepsilon;\varepsilon}(Y^*)$  and  $\Psi_{t_0;\varepsilon}^{-2}(Y^*) = \varphi_{t_0, t_0-4\pi\varepsilon;\varepsilon}(Y^*)$  work in the intervals  $t_0 \leq t \leq t_0 + 4\pi\varepsilon$  and  $t_0 \geq t \geq t_0 - 4\pi\varepsilon$  respectively. Difference formulae of arbitrarily high orders may also be employed, but higher order implies a wider stencil and costlier micro-integrations.

The approach suggested here is related to methods called envelop-following or multi-revolution (see [12, 2] and their references) that go back to the 1960's and have been successfully used in a number of application areas, including celestial mechanics and circuit theory. Note that, while in this paper both the macro- and micro-integrators are standard ODE solvers, the multi-revolution technique requires the construction of new special formulae. The closest relative of the algorithm described above is perhaps the LIPS method of Kirchgraber [10] that, in lieu of the finite difference formulae employed here, retrieves values of  $F(Y^*; \varepsilon)$  through Runge-Kutta like formulae. Again those formulae have to be build on purpose and reference [10] provides coefficients for the orders  $\mathcal{O}(\varepsilon^2)$ ,  $\mathcal{O}(\varepsilon^3)$ ,  $\mathcal{O}(\varepsilon^4)$ .<sup>1</sup>

## 4 Examples

In order that a highly-oscillatory problem (1) may be integrated by the procedure outlined above, it is necessary that the corresponding one-period map  $\Psi_{t_0;\varepsilon}$  be a smooth near-to-identity transformation as in (3). In this section we present families of systems that satisfy this condition.

<sup>1</sup> The possibility of using finite-difference formulae to approximate modified equations —this is essentially the problem solved by Kirchgraber's formulae— was already pointed out in reference [11], page 228.

(i) If  $f$  in (1) is of the form

$$f(y, \tau; \varepsilon) = \sum_{j=1}^{\infty} \varepsilon^{j-1} f_j(y, \tau). \quad (11)$$

where the  $f_j(y, \tau)$  are smooth  $2\pi$ -periodic functions of  $\tau$ , then  $f = \mathcal{O}(1)$  as  $\varepsilon \rightarrow 0$  and therefore  $y(t) - y(t_0)$  undergoes  $\mathcal{O}(\varepsilon)$  changes in the interval  $t_0 \leq t \leq t_0 + 2\pi\varepsilon$  and (3) holds. Presented in [4] is a way of systematically constructing with the help of rooted trees the functions  $M_j$  that feature in (3).

The format (11) is the standard starting point to perform analytically averaging so that any system to be averaged has first to be brought to that format via suitable changes of variables. We show next that those preliminary changes of variables are not needed to implement the numerical method of Sect. 3.

(ii) Consider second order systems of the form

$$\frac{d^2}{dt^2} q = G\left(q, \frac{t}{\varepsilon}; \varepsilon\right), \quad (12)$$

where  $q \in \mathcal{R}^d$  and the force  $G$  has an expansion

$$G(q, \tau; \varepsilon) = \sum_{j=0}^{\infty} \varepsilon^{j-1} G_j(q, \tau)$$

(the  $G_j$  are  $2\pi$ -periodic in  $\tau$ ).

To treat this case, we begin by rewriting (12) as a first order system

$$\frac{d}{dt} q = p, \quad \frac{d}{dt} p = G\left(q, \frac{t}{\varepsilon}; \varepsilon\right) \quad (13)$$

for the vector  $y = (q, p)$  in  $\mathcal{R}^D$ ,  $D = 2d$ . Note that here  $G = \mathcal{O}(1/\varepsilon)$  and the solution  $y$  will undergo  $\mathcal{O}(1)$  changes in the interval  $t_0 \leq t \leq t_0 + 2\pi\varepsilon$ . However if the leading term  $(1/\varepsilon)G_0$  of  $G$  averages to zero over one period, i.e.

$$\int_0^{2\pi} G_0(q, \tau) d\tau = 0, \quad (14)$$

then (3) holds as proved in [4], a reference that presents a technique for explicitly constructing the functions  $M_j$ . An alternative proof will be given here. Consider the system

$$\frac{d}{dt} q = 0, \quad \frac{d}{dt} p = \frac{1}{\varepsilon} G_0\left(q, \frac{t}{\varepsilon}\right), \quad (15)$$

denote by  $\widehat{\varphi}_{t_0, t; \varepsilon}(q_0, p_0)$  the corresponding solution operator and introduce the time-dependent change of variables

$$(q(t), p(t)) = \widehat{\varphi}_{t_0, t; \varepsilon}(\widehat{q}(t), \widehat{p}(t)).$$



Of course, this change reduces the system (15) to the trivial form  $(d/dt)\hat{q} = 0$  and  $(d/dt)\hat{p} = 0$ . When applied to the full (13), the change reduces the system to the format (11) (i.e. the new right-hand side contains no  $\mathcal{O}(1/\varepsilon)$  term). From case (i) above we conclude that (3) holds *after changing variables*. However the solution operator is explicitly given by

$$\widehat{\Phi}_{t_0,t;\varepsilon}(q_0, p_0) = \left( q_0, p_0 + \int_{t_0}^t \frac{1}{\varepsilon} G_0 \left( q_0, \frac{t'}{\varepsilon} \right) dt' \right)$$

an expression that, in tandem with (14), shows that the associated one-period map  $\widehat{\Phi}_{t_0,t_0+2\pi\varepsilon;\varepsilon}$  is the identity. Therefore at stroboscopic times  $t_n$  the values of the new  $(\hat{q}, \hat{p})$  variables coincide with the values of the old variables  $(q, p)$  and (3) also holds *without changing variables*. As a consequence the numerical method works for the given system (13) without any need to previously perform any analytic manipulations.

Note that the expression of the change of variables reveals that in the interval  $t_0 \leq t \leq t_0 + 2\pi\varepsilon$ , the variations of the variable  $p(t)$  are  $\mathcal{O}(1)$  and those in  $q(t)$  are  $\mathcal{O}(\varepsilon)$ . At the end of the interval, both  $q(t_0 + 2\pi\varepsilon)$  and  $p(t_0 + 2\pi\varepsilon)$  are  $\mathcal{O}(\varepsilon)$  away from their initial values  $q(t_0)$  and  $p(t_0)$  in view of (3).

A well known example of (12) is given by the vibrated inverted pendulum equation

$$\frac{d^2}{dt^2} q = G \left( q, \frac{t}{\varepsilon}; \varepsilon \right) = \left( \frac{1}{\varepsilon} \frac{v_{max}}{\ell} \cos \left( \frac{t}{\varepsilon} + \theta_0 \right) + \frac{g}{\ell} \right) \sin q. \quad (16)$$

(iii) The reader is referred to [10] and [4] for further examples (including perturbed Kepler problems, perturbed harmonic oscillators, Fermi-Pasta-Ulam like problems) of systems for which (3) holds because they may be brought to the format (11) through a change of variables that coincides with the identity map at stroboscopic times.

## 5 Numerical Experiments

Our aim in this section is to illustrate by means of simple examples the use of the stroboscopic technique described in Sect. 3. For this reason we only report on experiments performed when the macro-integrator is either the ‘classical’ fourth-order, four stages Runge-Kutta (RK) method with constant step-sizes or the variable-step code ode45 from MATLAB. Extensive numerical experiments, including detailed comparisons with alternative techniques and wider choices of macro- and micro-solvers, will be presented elsewhere.

As a test problem, we integrate in the interval  $t_0 = 0 \leq t \leq \pi$  the inverted (Kapitsa) pendulum equation (16) with parameter values  $v_{max} = 4$ ,  $\ell = 0.2$ ,  $\theta_0 = 2$ ,  $g = 9.8$ , and initial conditions  $q(0) = 0.25$ ,  $p(0) = 0$ . This equation has been used as a test example in [16] to illustrate the power of the heterogeneous multiscale approach (see also [14, 3, 4]). Unlike the algorithms described in this paper, those

analyzed in [16] require some preliminary analytical work to derive formulae that relate macro- and micro-states.

### 5.1 Constant Step-Sizes

We first take the classical RK method with constant step-sizes as macro- and micro-integrator. This is run, for different values of  $\varepsilon$ , for combinations of macro- and micro-steps  $(H, h)$  of the form  $(2\pi 2^{-\nu}/50, 2\pi\varepsilon 2^{-\nu}/4)$ ,  $\nu = 0, 1, 2, \dots$  and with either second- or fourth-order differences (see (9) or (10) respectively).<sup>2</sup> The results are summarized in Tables 1 and 2 respectively. In the former, the symbol \*\*\* means that the corresponding run was not carried out: when  $H$  is smaller than  $2\pi\varepsilon$  the stroboscopic algorithm does not make any sense.

**Table 1** Errors in stroboscopic algorithm: 2nd-order finite differences

$H$	Micro evaluations	$1/\varepsilon$			
		3,200	6,400	12,800	25,600
$2\pi/50$	3,200	3.12(-1)	3.12(-1)	3.12(-1)	3.12(-1)
$2\pi/100$	12,800	2.14(-2)	2.16(-2)	2.17(-2)	2.17(-2)
$2\pi/200$	51,200	3.22(-3)	2.17(-3)	1.94(-3)	1.88(-3)
$2\pi/400$	204,800	1.59(-3)	5.31(-4)	2.67(-4)	2.02(-4)
$2\pi/800$	819,200	1.42(-3)	3.65(-4)	1.01(-4)	3.54(-5)
$2\pi/1,600$	3,276,800	1.41(-3)	3.53(-4)	8.88(-5)	2.29(-5)
$2\pi/3,200$	13,107,200	1.41(-3)	3.52(-4)	8.80(-5)	2.20(-5)
$2\pi/6,400$	52,428,800	***	3.52(-4)	8.79(-5)	2.20(-5)
$2\pi/12,800$	209,715,200	***	***	8.79(-5)	2.20(-5)
$2\pi/25,600$	838,860,800	***	***	***	2.20(-5)

Let us first discuss the computational cost. Since each micro-integration takes place in an interval of width  $4\pi\varepsilon$  (or  $8\pi\varepsilon$ ) and, for given  $H$ , the value of  $h$  is chosen to be proportional to  $\varepsilon$ , the cost of the algorithm is *independent of*  $\varepsilon$ . Furthermore when  $H$  is halved so is  $h$  and therefore the total number of micro-steps in a run is multiplied by four (see the second column of the tables that display the total number of function evaluations required by the micro-integrations).

We report errors measured as the maximum, over all macro-step-points, of the (absolute value of the) difference between the  $q$  component of a very accurate nu-

<sup>2</sup> Our experience indicates that standard central differences of order 6 are not competitive in terms of efficiency with those of orders 2 or 4.

**Table 2** Errors in stroboscopic algorithm: 4th-order finite differences

$H$	Micro evaluations	$1/\varepsilon$			
		3,200	6,400	12,800	25,600
$2\pi/50$	6,400	3.12(-1)	3.12(-1)	3.12(-1)	3.12(-1)
$2\pi/100$	25,600	2.18(-2)	2.17(-2)	2.17(-2)	2.17(-2)
$2\pi/200$	102,400	1.87(-3)	1.86(-3)	1.86(-3)	1.86(-3)
$2\pi/400$	409,600	1.81(-4)	1.81(-4)	1.80(-4)	1.80(-4)
$2\pi/800$	1,638,400	1.36(-5)	1.35(-5)	1.34(-5)	1.34(-5)
$2\pi/1,600$	6,553,600	1.05(-6)	9.18(-7)	9.09(-7)	9.04(-7)
$2\pi/3,200$	26,214,400	2.01(-7)	6.74(-8)	5.89(-8)	5.45(-8)

merical approximation to the true solution of the oscillatory problem and the solution  $Q$  provided by the stroboscopic algorithm; errors in  $p$  behave in exactly the same way as those in  $q$ . There are three sources of error (cf. [14]): (i) the recovery of the right-hand side  $F$  of the averaged system by the finite-difference formula (9) (or (10)), (ii) the replacement in (9) (or (10)) of the exact values of  $\Psi_{t_0;\varepsilon}^k(Y^*)$  by numerical approximations based on micro-integrations, (iii) the discretization error introduced by the macro-integrator. We consider these sources in turn.

As  $H$  and  $h$  tend to 0, the errors arising from (ii) and (iii) vanish and only the source (i) remains. At each evaluation of  $F$  the error from this source is  $\mathcal{O}(\varepsilon^2)$  (or  $\mathcal{O}(\varepsilon^4)$ ) and, due to the stability of the macro-solver, these evaluation errors introduce  $\mathcal{O}(\varepsilon^2)$  (or  $\mathcal{O}(\varepsilon^4)$ ) errors in the values of  $Q$ . This is apparent in Table 1, where the errors at the bottom of the different columns, clearly behave as  $\mathcal{O}(\varepsilon^2)$ . For fourth-order differences Table 2 does not report results for very small  $H$  and  $h$  due to the cost of obtaining a sufficiently accurate reference solution to measure errors.

To analyze the micro-integration errors, it is best to rewrite (16) in terms of the fast, non-dimensional time  $\tau = t/\varepsilon$ , i.e.

$$\frac{d}{d\tau}q = \varepsilon p, \quad \frac{d}{d\tau}p = \varepsilon G\left(q, \frac{t}{\varepsilon}; \varepsilon\right) = \left(\frac{v_{\max}}{\ell} \cos(\tau + \theta_0) + \varepsilon \frac{g}{\ell}\right) \sin q. \quad (17)$$

Now the force  $\varepsilon G$  is bounded as  $\varepsilon \rightarrow 0$ , the micro-integrations span intervals of fixed length  $4\pi$  (or  $8\pi$ ) and (because the micro-step  $h$  in the variable  $t$  is chosen proportional to  $\varepsilon$ ) the step-length  $h/\varepsilon$  in  $\tau$  is also independent of  $\varepsilon$ . Therefore, standard results show that the error in finding each value  $\Psi_{t_0;\varepsilon}^k(Y^*)$  is  $\mathcal{O}((h/\varepsilon)^4)$ . Furthermore it can be shown that the constant  $C$  implied in the  $\mathcal{O}$  notation is itself  $\mathcal{O}(\varepsilon)$ ;<sup>3</sup> the extra factor in  $C$  makes up for the factor  $\varepsilon$  that features in the denom-

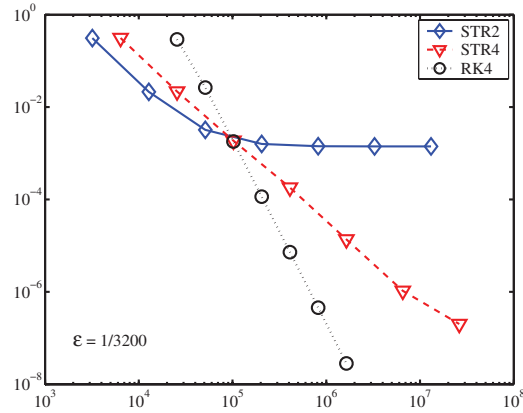
<sup>3</sup> The proof of the estimate  $C = \mathcal{O}(\varepsilon)$  is easy after noting that for  $\varepsilon = 0$  the RK micro-integrator finds the solution of (17) at  $\tau = 2\pi$  without any error. (In fact finding the solution at  $\tau = 2\pi$  of

inator of (9) (or (10)) and therefore, in each evaluation of  $F$ , the error due to the micro-integrator is  $\mathcal{O}((h/\varepsilon)^4)$ , where the implied constant is  $\varepsilon$ -independent. Again the stability of the macro-solver entails that the corresponding effect in the macro-solution  $Q$  is itself  $\mathcal{O}((h/\varepsilon)^4)$ , or, with our choice of  $H$  and  $h$ ,  $\mathcal{O}(H^4)$ . Since the error due to discretizing the averaged equation is itself  $\mathcal{O}(H^4)$ , we conclude that the combined effect of sources (ii) and (iii) is  $\mathcal{O}(H^4)$ , *uniformly in  $\varepsilon$* . In this way, the overall algorithm yields approximations to the true  $q$  and  $p$  of sizes  $\mathcal{O}(\varepsilon^\mu + H^4)$ , where the implied constant is independent of  $\varepsilon$  and  $\mu = 2$  or  $\mu = 4$  for second and fourth-order differences respectively. Thus, unless  $H$  is chosen to be so small that the contribution of size  $\varepsilon^\mu$  manifests itself, the algorithm yields errors that behave as  $\mathcal{O}(H^4)$  *uniformly in  $\varepsilon$  at a cost that is also independent of  $\varepsilon$* . Once more this is borne out by the tables, where the errors in the top rows are independent of  $\varepsilon$  and of the finite-difference formula and show a reduction by a factor of  $\approx 16$  when  $H$  is halved.

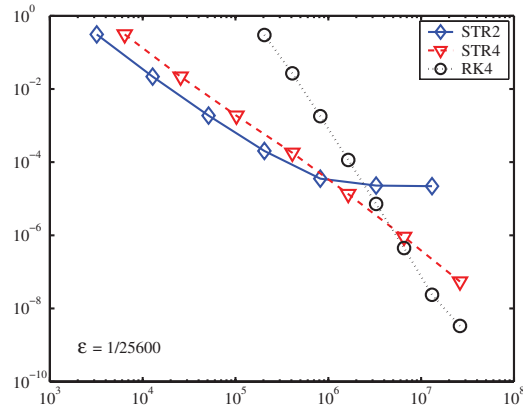
Figs. 3 and 4 are based on Tables 1 and 2 and compare the efficiency of the stroboscopic algorithm with second or fourth-order differences with that of a straight-forward integration of the oscillatory problem with the classical RK method. For errors of size  $\approx 10^{-2}$ , Fig. 3 reveals that for  $\varepsilon = 1/3, 200$  the second-difference algorithm needs an amount of work that is less than  $1/5$  of that required by the classical method. For  $\varepsilon = 1/25, 600$ , we see in Fig. 4 that the same ratio is less than  $1/30$ . Also note that for the algorithm based on fourth-order differences, the lines in Figs. 3 and 4 are virtually identical, indicating an  $\varepsilon$ -independent behavior. The line corresponding to the classical RK method undergoes a marked translation to the right when  $\varepsilon$  is decreased, indicating an efficiency loss. For the algorithm with second-order differences, the lines in both figures coincide for larger values of the errors (larger values of  $H$ ); however in Fig. 3 errors saturate at a larger value than that in Fig. 4 in agreement with earlier discussions. Finally we point out that the lines of the stroboscopic algorithms possess a smaller slope than those of the RK method: while to divide the error by a factor of 16 the classical method has to work twice as hard, the new algorithms must toil four times as hard, as they require both more macro-steps and more accurate micro-integrations.

---

(17) with  $\varepsilon = 0$  essentially requires the computation of the integral in (14); the RK numerical solution may be written down in closed form as a trigonometric sum whose value vanishes.) The key point here is that the micro-integrator is such that when applied to the system (15) it generates a one-period map that *exactly* coincides with the identity, thus mimicking a key property of the system being integrated. For micro-integrators that do not possess this property the error behavior is not so favorable as for those considered here because estimates suffer from the factor  $\varepsilon$  in the denominator of the finite-difference formulae (cf. our analysis with that in [10]). Similarly, when integrating perturbed Kepler problems, perturbed harmonic oscillators, etc. as in [10] or [4], it is important that the micro-integration be performed in such a way that for the unperturbed problem ( $\varepsilon = 0$ ) it results in a one-period map that coincides *exactly* with the identity. This may be achieved by using splitting methods.



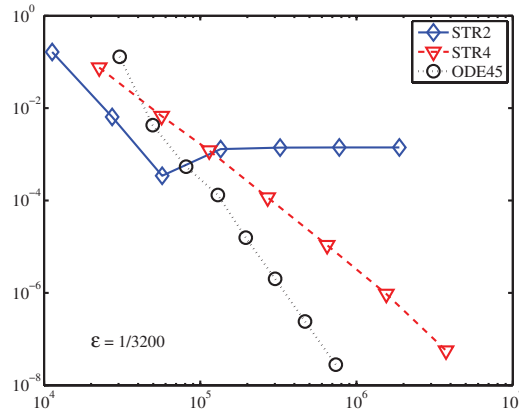
**Fig. 3** Efficiency comparison: errors vs. number of evaluations of the micro-force. Constant step-sizes, ‘larger’  $\epsilon$



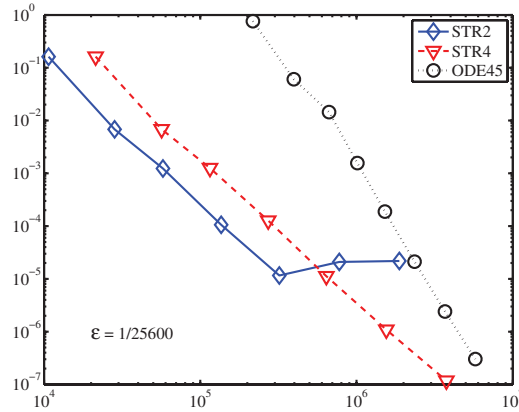
**Fig. 4** Efficiency comparison: errors vs. number of evaluations of the micro-force. Constant step-sizes, smaller  $\epsilon$

## 5.2 Variable Step-Sizes

To illustrate the use of the stroboscopic algorithm with variable macro-step sizes we ran the ode45 MATLAB as macro-integrator with absolute error tolerances  $Tol$  from the sequence  $10^{-2}, 10^{-3}, \dots, 10^{-8}$  (the relative error tolerance was taken to be equal to the absolute tolerance). For reasons discussed in the preceding subsection is important that the micro-integration is performed by a method that solves (17) exactly at  $\tau = 2\pi$  for  $\epsilon = 0$ ; we decided to micro-integrate, *with constant step-*



**Fig. 5** Efficiency comparison: errors vs. number of evaluations of the micro-force. Variable-step macro-solver, 'larger'  $\epsilon$



**Fig. 6** Efficiency comparison: errors vs. number of evaluations of the micro-force. Variable-step macro-solver, smaller  $\epsilon$

sizes, by means of the fifth-order RK formula of the pair used by ode45.<sup>4</sup> We took  $h = (2\pi\epsilon)/\nu$  where  $\nu$  is the smallest integer for which  $(2\pi/\nu)^5 \leq 1000 \times Tol$ ; this equilibrates the accuracy of the macro- and micro-integrations in a way similar to that analyzed in the preceding subsection. (The values of  $\nu$  for the seven values of  $Tol$  turn out to be 4, 7, 10, 16, 26, 40, 63.) The variable-step macro-integrator chooses step-points that of course do not coincide with stroboscopic times but, as discussed in Sect. 3, this causes no problem to the stroboscopic algorithm. To mea-

<sup>4</sup> The use of the variable-step code ode45 as micro-integrator for (17) with  $\epsilon = 0$  yields errors that, after one period, are small but not exactly zero.

sure errors we took advantage of the dense output capabilities of ode45 and generated output of the macro-integration at each stroboscopic time. Errors were then measured as the maximum, over all stroboscopic times, of the (absolute value of the) difference between the  $q$  component of the reference solution and the output  $Q$  provided by the algorithms.

Figures 5 and 6 compare the efficiency of the stroboscopic algorithms with that of a straightforward integration of the oscillatory problem with ode45. Again the stroboscopic algorithm exhibits a behavior that, unless  $Tol$  is so small that errors saturate, is  $\varepsilon$ -independent. Clearly, for small values of  $\varepsilon$ , this uniformity renders them more efficient than the conventional integrator, whose performance is degraded as  $\varepsilon \downarrow 0$ .

## Acknowledgement

This research has been supported by ‘Acción Integrada entre España y Francia’ HF2008-0105. M.P. Calvo and J.M. Sanz-Serna are also supported by project MTM2007-63257 (Ministerio de Educación, España). A. Murua is also supported by projects MTM2007-61572 (Ministerio de Educación, España) and EHU08/43 (Universidad del País Vasco/Euskal Herriko Unibertsitatea).

## References

1. Ariel, G., Engquist, B., Tsai, R.: A multiscale method for highly oscillatory ordinary differential equations with resonance. *Math. Comput.* **78**, 929–956 (2009)
2. Calvo, M., Jay, L.O., Montijano, J.I., Rández, L.: Approximate compositions of a near identity map by multi-revolution Runge-Kutta methods. *Numer. Math.* **97**, 635–666 (2004)
3. Calvo, M.P., Sanz-Serna, J.M.: Heterogeneous Multiscale Methods for mechanical systems with vibrations. *SIAM J. Sci. Comput.* **32**, 2029–2046, (2010)
4. Chartier, Ph., Murua, A., Sanz-Serna, J.M.: Higher-order averaging, formal series and numerical integration I: B-series. *Found. Comput. Math* **10**, 695–727 (2010)
5. E., W.: Analysis of the heterogeneous multiscale method for ordinary differential equations. *Comm. Math. Sci.* **1**, 423–436 (2003)
6. E., W., Engquist, B.: The heterogeneous multiscale methods. *Comm. Math. Sci.* **1**, 87–132 (2003)
7. E., W., Engquist, B., Li, X., Ren, W., Vanden-Eijnden, E.: Heterogeneous multiscale methods: A review. *Commun. Comput. Phys.* **2**, 367–450 (2007)
8. Engquist, B., Tsai, R.: Heterogeneous multiscale methods for stiff ordinary differential equations. *Math. Comput.* **74**, 1707–1742 (2005)
9. Hairer, E., Lubich, Ch., Wanner, G.: *Geometric Numerical Integration*, 2nd ed. Springer, Berlin (2006)
10. Kirchgraber, U.: An Ode-solver based on the method of averaging. *Numer. Math.* **53**, 621–652 (1988)
11. Murua, A.: Formal series and numerical integrators, Part I: Systems of ODEs and symplectic integrators. *Appl. Numer. Math.* **29**, 221–251 (1999)

12. Petzold, L.R., Jay, L.O., Yen, J.: Numerical solution of highly oscillatory ordinary differential equations. *Acta Numerica* **6**, 437–484 (1997)
13. Sanders, J.A., Verhulst, F., Murdock, J.: *Averaging Methods in Nonlinear Dynamical Systems*, 2nd ed. Springer, New York (2007)
14. Sanz-Serna, J.M.: Modulated Fourier expansions and heterogeneous multiscale methods. *IMA J. Numer. Anal.* **29**, 595–605 (2009)
15. Sanz-Serna, J.M., Calvo, M.P.: *Numerical Hamiltonian Problems*. Chapman and Hall, London (1994)
16. Sharp, R., Tsai, Y.-H., Engquist, B.: Multiple time scale numerical methods for the inverted pendulum problem. In: Engquist, B., Lötsdtedt, P., Runborg, O. (eds) *Multiscale Methods in Science and Engineering*, *Lect. Notes Comput. Sci. Eng.* **44**, pp. 241–261. Springer, Berlin (2005)