



# Analysis of spectral Hamiltonian boundary value methods (SHBVMs) for the numerical solution of ODE problems

Pierluigi Amodio<sup>1</sup> · Luigi Brugnano<sup>2</sup>  · Felice Iavernaro<sup>1</sup>

Received: 16 November 2018 / Accepted: 15 May 2019 / Published online: 25 May 2019  
© Springer Science+Business Media, LLC, part of Springer Nature 2019

## Abstract

Recently, the numerical solution of stiffly/highly oscillatory Hamiltonian problems has been attacked by using Hamiltonian boundary value methods (HBVMs) as spectral methods in time. While a theoretical analysis of this spectral approach has been only partially addressed, there is enough numerical evidence that it turns out to be very effective even when applied to a wider range of problems. Here, we fill this gap by providing a thorough convergence analysis of the methods and confirm the theoretical results with the aid of a few numerical tests.

**Keywords** Spectral methods · Legendre polynomials · Hamiltonian boundary value methods · HBVMs · SHBVMs

**Mathematics Subject Classification (2010)** 65L05 · 65P10

## 1 Introduction

In recent years, the efficient numerical solution of Hamiltonian problems has been tackled via the definition of the energy-conserving Runge-Kutta methods named Hamiltonian boundary value methods (HBVMs) [10–13]. Such methods have been

---

✉ Luigi Brugnano  
luigi.brugnano@unifi.it

Pierluigi Amodio  
pierluigi.amodio@uniba.it

Felice Iavernaro  
felice.iavernaro@uniba.it

<sup>1</sup> Dipartimento di Matematica, Università di Bari, Bari, Italy

<sup>2</sup> Dipartimento di Matematica e Informatica “U. Dini”, Università di Firenze, Florence, Italy

developed along several directions (see, e.g., [4, 14, 19]), including Hamiltonian BVPs [1] and Hamiltonian PDEs (see, e.g., [6]): we also refer to the monograph [7] and to the recent review paper [8] for further details.

More recently, HBVMs have been used as spectral methods in time for solving highly oscillatory Hamiltonian problems [18], as well as stiffly oscillatory Hamiltonian problems [9] emerging from the space semi-discretization of Hamiltonian PDEs. Their spectral implementation is justified by the fact that this family of methods performs a projection of the vector field onto a finite dimensional subspace via a least squares approach based on the use of Legendre orthogonal polynomials [13]. This spectral approach, supported by a very efficient nonlinear iteration technique to handle the large nonlinear systems needed to advance the solution in time (see [12], [7, Chapter 4] and [8]), proved to be very effective. However, a thorough convergence analysis of HBVMs, used as spectral methods, was still lacking. In fact, when using large stepsizes, as is required by the spectral strategy, the notion of classical order of a method is not sufficient to explain the correct asymptotic behavior of the solutions, so that a different analysis is needed, which is the main goal of the present paper. Moreover, the theoretical achievements will be numerically confirmed by applying the methods to a number of ODE-IVPs.

It is worth mentioning that early references where numerical methods were derived by projecting the vector field onto a finite dimensional subspace are, e.g., [2, 3, 27, 28] (a related reference is [36]). A similar technique, popular for solving oscillatory problems, is that of *exponential/trigonometrically fitted* methods and, more in general, *functionally fitted* methods [22–24, 26, 29–33, 35, 37].

With these premises, the structure of the paper is as follows: in Section 2, we analyze the use of spectral methods in time; in Section 3, we discuss the efficient implementation of the fully discrete method; in Section 4, we provide numerical evidence of the effectiveness of such an approach, confirming the theoretical achievements. At last, a few conclusions are reported in Section 5.

## 2 Spectral approximation in time

This section contains the main theoretical results regarding the spectral methods that we shall use for the numerical solution of ODE-IVPs which, without loss of generality, will be assumed in the form<sup>1</sup>

$$\dot{y}(t) = f(y(t)), \quad y(0) = y_0 \in \mathbb{R}^m. \quad (1)$$

Hereafter,  $f$  is assumed to be suitably smooth (in particular, we shall assume  $f(z)$  to be analytic in a closed complex ball centered at  $y_0$ ). We consider the solution of problem (1) on the interval  $[0, h]$ , where  $h$  stands for the time-step to be used by a one-step numerical method. The same arguments will be then repeated for the subsequent integration steps. According to [13], we consider the expansion of the

<sup>1</sup>In fact, if problem (1) is non autonomous,  $t$  can be included in the state vector.

right-hand side of (1) along the shifted and scaled Legendre polynomial orthonormal basis  $\{P_j\}_{j \geq 0}$ ,

$$P_j \in \Pi_j, \quad \int_0^1 P_i(x)P_j(x)dx = \delta_{ij}, \quad i, j = 0, 1, \dots,$$

with  $\Pi_j$  the set of polynomials of degree  $j$  and  $\delta_{ij}$  the Kronecker delta. One then obtains:

$$\dot{y}(ch) = f(y(ch)) \equiv \sum_{j \geq 0} P_j(c)\gamma_j(y), \quad c \in [0, 1], \tag{2}$$

with the Fourier coefficients  $\gamma_j(y)$  given by

$$\gamma_j(y) = \int_0^1 P_j(\tau)f(y(\tau h))d\tau, \quad j = 0, 1, \dots \tag{3}$$

We recall that:

$$\|P_j\| := \max_{x \in [0,1]} |P_j(x)| = \sqrt{2j + 1}, \quad \int_0^1 P_j(x)q(x)dx = 0, \quad \forall q \in \Pi_{j-1}. \tag{4}$$

Let us now study the properties of the coefficients  $\gamma_j(y)$  defined at (3). To begin with, we report the following results.

**Lemma 1** *Let  $g : [0, h] \rightarrow V$ , with  $V$  a vector space, admit a Taylor expansion at 0. Then,*

$$\int_0^1 P_j(c)g(ch)dc = O(h^j).$$

*Proof* See [13, Lemma 1]. □

**Corollary 1** *The Fourier coefficients defined in (3) satisfy:  $\gamma_j(y) = O(h^j)$ .*

We now want to derive an estimate which generalizes the result of Corollary 1 to the case where the stepsize  $h$  is not small. For this purpose, hereafter we assume that the solution  $y(t)$  of (1) admits a complex analytic extension in a neighbourhood of 0. Moreover, we shall denote by  $\mathcal{B}(0, r)$  the closed ball of center 0 and radius  $r$  in the complex plane, and  $\mathcal{C}(0, r)$  the corresponding circumference. The following results then hold true.

**Lemma 2** *Let  $P_j$  be the  $j$ th shifted and scaled Legendre polynomial and, for  $\rho > 1$ , let us define the function*

$$Q_j(\xi) = \int_0^1 \frac{P_j(c)}{\xi - c} dc, \quad \xi \in \mathcal{C}(0, \rho). \tag{5}$$

Then,

$$\|Q_j\|_\rho := \max_{\xi \in \mathcal{C}(0, \rho)} |Q_j(\xi)| \leq \sqrt{\frac{2}{j + 1}} \frac{1}{(\rho - 1)\rho^j}. \tag{6}$$

*Proof* One has, for  $|\xi| = \rho > 1$ , and taking into account (4):

$$\begin{aligned} Q_j(\xi) &= \int_0^1 \frac{P_j(c)}{\xi - c} dc = \xi^{-1} \int_0^1 \frac{P_j(c)}{1 - \xi^{-1}c} dc = \xi^{-1} \int_0^1 P_j(c) \sum_{\ell \geq 0} \xi^{-\ell} c^\ell dc \\ &= \xi^{-1} \sum_{\ell \geq j} \xi^{-\ell} \int_0^1 P_j(c) c^\ell dc = \xi^{-j-1} \sum_{\ell \geq 0} \xi^{-\ell} \int_0^1 P_j(c) c^{\ell+j} dc. \end{aligned}$$

Passing to norms, one has:

$$\begin{aligned} \left| \int_0^1 P_j(c) c^{\ell+j} dc \right| &\leq \|P_j\| \int_0^1 c^{\ell+j} dc = \frac{\sqrt{2j+1}}{j+\ell+1} \leq \sqrt{\frac{2}{j+1}}, \\ \left| \xi^{-j-1} \sum_{\ell \geq 0} \xi^{-\ell} \right| &\leq \rho^{-j-1} \sum_{\ell \geq 0} \rho^{-\ell} = [(\rho - 1)\rho^j]^{-1}, \end{aligned}$$

from which (6) follows. □

**Lemma 3** *Let  $g(z)$  be analytic in the closed ball  $\mathcal{B}(0, r^*)$  of the complex plane, for a given  $r^* > 0$ . Then, for all  $0 < h \leq h^* < r^*$ ,*

$$g_h(\xi) := g(\xi h) \tag{7}$$

*is analytic in  $\mathcal{B}(0, \rho)$ , with*

$$\rho \equiv \rho(h) := \frac{r^*}{h} \geq \frac{r^*}{h^*} =: \rho^* > 1. \tag{8}$$

We are now in the position of stating the following result.<sup>2</sup>

**Theorem 1** *Assume that the function*

$$g(z) := f(y(z)) \tag{9}$$

*and  $h$  in (2)–(3) satisfy the hypotheses of Lemma 3. Then, there exists  $\kappa = \kappa(h^*) > 0$ , such that<sup>3</sup>*

$$|\gamma_j(y)| \leq \frac{\kappa}{\sqrt{j+1}} \rho^{-j}. \tag{10}$$

<sup>2</sup>The used arguments are mainly adapted from [21].

<sup>3</sup>Hereafter, for sake of clarity, we shall denote by  $|\cdot|$  any convenient vector or matrix norm.

*Proof* By considering the function (7) corresponding to (9), and with reference to the function  $Q_j(\xi)$  defined in (5), one has that the parameter  $\rho$ , as defined in (8), is greater than 1 and, moreover, (see (3))

$$\begin{aligned} \gamma_j(y) &= \int_0^1 P_j(c) f(y(ch)) dc \equiv \int_0^1 P_j(c) g_h(c) dc \\ &= \int_0^1 P_j(c) \left[ \frac{1}{2\pi i} \int_{C(0,\rho)} \frac{g_h(\xi)}{\xi - c} d\xi \right] dc \\ &= \frac{1}{2\pi i} \int_{C(0,\rho)} g_h(\xi) \left[ \int_0^1 \frac{P_j(c)}{\xi - c} dc \right] d\xi \equiv \frac{1}{2\pi i} \int_{C(0,\rho)} g_h(\xi) Q_j(\xi) d\xi. \end{aligned}$$

Then, passing to norms (see (6)),

$$|\gamma_j(y)| \leq \rho \|g_h\|_\rho \|Q_j\|_\rho.$$

Moreover, observing that (see (9), (7), and (8)):

$$\|g_h\|_\rho := \max_{\xi \in C(0,\rho)} |g_h(\xi)| \leq \max_{\xi \in B(0,\rho)} |g_h(\xi)| \equiv \max_{z \in B(0,r^*)} |g(z)| =: \|g\|,$$

and using (6), one has, (see (8)):

$$|\gamma_j(y)| \leq \frac{\|g\|}{(1 - \rho^{-1})} \sqrt{\frac{2}{j+1}} \rho^{-j} \leq \frac{\|g\|}{(1 - (\rho^*)^{-1})} \sqrt{\frac{2}{j+1}} \rho^{-j},$$

from which (10) eventually follows. □

*Remark 1* It is worth mentioning that, in the bound (10), the dependence on  $h$  only concerns the parameter  $\rho > 1$ , via the expression (8), from which one infers that  $\rho \sim h^{-1}$ , for all  $0 < h \leq h^* < r^*$ . This, in turn, is consistent with the result of Corollary 1, when  $h \rightarrow 0$ .

Let us now consider a polynomial approximation to (2),

$$\dot{\sigma}(ch) = \sum_{j=0}^{s-1} P_j(c) \gamma_j(\sigma), \quad c \in [0, 1], \tag{11}$$

where  $\gamma_j(\sigma)$  is defined according to (3) by formally replacing  $y$  by  $\sigma$ , i.e.,

$$\gamma_j(\sigma) = \int_0^1 P_j(\tau) f(\sigma(\tau h)) d\tau, \quad j = 0, 1, \dots, s - 1. \tag{12}$$

Integrating term by term (11), and imposing the initial condition in (1), provide us with the polynomial approximation of degree  $s$ :

$$\sigma(ch) = y_0 + h \sum_{j=0}^{s-1} \int_0^c P_j(x) dx \gamma_j(\sigma), \quad c \in [0, 1]. \tag{13}$$

We now want to assess the extent to which  $\sigma(ch)$  approximates  $y(ch)$ , for  $c \in [0, 1]$ . When  $h \rightarrow 0$ , it is known that  $y(h) - \sigma(h) = O(h^{2s+1})$  (see, e.g., [7, 8, 13]).

Nevertheless, we here discuss the approximation of  $\sigma$  to  $y$ , in the interval  $[0, h]$ , when  $h$  is finite and only assuming that the result of Theorem 1 is valid. The following result then holds true.<sup>4</sup>

**Theorem 2** *Let  $y$  be the solution of (1),  $\sigma$  be defined according to (13), and assume that  $f(\sigma(z))$  is analytic in  $\mathcal{B}(0, r^*)$ , for a given  $r^* > 0$ . Then, for all  $0 < h \leq h^* < r^*$ , there exist  $M, \bar{M} > 0$ ,  $M = M(h^*)$ ,  $\bar{M} = \bar{M}(h^*)$ , and  $\rho > 1$ ,  $\rho \sim h^{-1}$ , such that:*

- $|\sigma(ch) - y(ch)| \leq chM\rho^{-s}$ ,  $c \in [0, 1]$ ;
- $|\sigma(h) - y(h)| \leq h\bar{M}\rho^{-2s}$ .

*Proof* Let  $y(t, \xi, \eta)$  denote the solution of the problem

$$\dot{y} = f(y), \quad t \geq \xi, \quad y(\xi) = \eta,$$

and  $\Phi(t, \xi)$  be the solution of the associated variational problem,

$$\dot{\Phi}(t, \xi) = f'(y(t, \xi, \eta))\Phi(t, \xi), \quad t \geq \xi, \quad \Phi(\xi, \xi) = I,$$

having set  $f'$  the Jacobian of  $f$ . Without loss of generality, we shall assume that the function

$$\hat{g}(z) := \Phi(h, z), \tag{14}$$

is also analytic in  $\mathcal{B}(0, r^*)$ . We also recall the following well-known perturbation results:

$$\frac{\partial}{\partial \eta}y(t, \xi, \eta) = \Phi(t, \xi), \quad \frac{\partial}{\partial \xi}y(t, \xi, \eta) = -\Phi(t, \xi)f(\eta).$$

Consequently, from (12) and (13), one has:

$$\begin{aligned} \sigma(ch) - y(ch) &= y(ch, ch, \sigma(ch)) - y(ch, 0, \sigma(0)) = \int_0^{ch} \frac{d}{dt}y(ch, t, \sigma(t)) dt \\ &= \int_0^{ch} \left[ \frac{\partial}{\partial \xi}y(ch, \xi, \sigma(t)) \Big|_{\xi=t} + \frac{\partial}{\partial \eta}y(ch, t, \eta) \Big|_{\eta=\sigma(t)} \dot{\sigma}(t) \right] dt \\ &= -h \int_0^c \Phi(ch, \tau h) [f(\sigma(\tau h)) - \dot{\sigma}(\tau h)] d\tau \\ &= -h \int_0^c \Phi(ch, \tau h) \left[ \sum_{j \geq 0} P_j(\tau)\gamma_j(\sigma) - \sum_{j=0}^{s-1} P_j(\tau)\gamma_j(\sigma) \right] d\tau \\ &= -h \sum_{j \geq s} \left[ \int_0^c P_j(\tau)\Phi(ch, \tau h) d\tau \right] \gamma_j(\sigma). \end{aligned} \tag{15}$$

From the result of Theorem 1 applied to  $g(z) := f(\sigma(z))$ , we know that there exist  $\kappa = \kappa(h^*)$  and  $\rho > 1$ ,  $\rho \sim h^{-1}$ , such that, for the Fourier coefficients defined in

<sup>4</sup>The proof uses arguments similar to those of [13, Theorem 4].

(12),

$$|\gamma_j(\sigma)| \leq \frac{\kappa}{\sqrt{j+1}} \rho^{-j}. \tag{16}$$

Moreover, (see (4))  $\|P_j\| = \sqrt{2j+1}$  and, considering that, for all  $h \in (0, h^*]$ ,

$$\max_{x_1, x_2 \in [0, h]} |\Phi(x_1, x_2)| \leq \max_{x_1, x_2 \in [0, h^*]} |\Phi(x_1, x_2)| =: v \equiv v(h^*),$$

one has that

$$\left| \int_0^c P_j(\tau) \Phi(ch, \tau h) d\tau \right| \leq cv\sqrt{2j+1}.$$

Consequently, the first statement follows from (15) by setting, with reference to the parameter  $\rho^*$  defined in (8),

$$M = \frac{v\kappa\sqrt{2}}{1 - (\rho^*)^{-1}},$$

since, for all  $c \in [0, 1]$ :

$$\begin{aligned} |\sigma(ch) - y(ch)| &\leq chv\kappa \sum_{j \geq s} \sqrt{\frac{2j+1}{j+1}} \rho^{-j} \leq chv\kappa\sqrt{2} \sum_{j \geq s} \rho^{-j} \\ &= chv\kappa\sqrt{2} \frac{\rho^{-s}}{1 - \rho^{-1}} \leq ch \frac{v\kappa\sqrt{2}}{1 - (\rho^*)^{-1}} \rho^{-s} \equiv chM\rho^{-s}. \end{aligned}$$

To prove the second statement (i.e., when  $c = 1$ ), we observe that the result of Theorem 1 holds true also for the function (14) involved in the Fourier coefficients

$$\int_0^1 P_j(\tau) \Phi(h, \tau h) d\tau.$$

Consequently, there exist  $\kappa_1 = \kappa_1(h^*) > 0$ , such that

$$\left| \int_0^1 P_j(\tau) \Phi(h, \tau h) d\tau \right| \leq \frac{\kappa_1}{\sqrt{j+1}} \rho^{-j}. \tag{17}$$

The second statement then follows again from (15) by setting

$$\bar{M} = \frac{\kappa_1\kappa}{1 - (\rho^*)^{-2}},$$

so that, by using same steps as above:

$$\begin{aligned} |\sigma(h) - y(h)| &\leq h\kappa_1\kappa \sum_{j \geq s} \frac{\rho^{-2j}}{j+1} \leq h\kappa_1\kappa \sum_{j \geq s} \rho^{-2j} \\ &= h\kappa_1\kappa \frac{\rho^{-2s}}{1 - \rho^{-2}} \leq h \frac{\kappa_1\kappa}{1 - (\rho^*)^{-2}} \rho^{-2s} \equiv h\bar{M}\rho^{-2s}. \end{aligned}$$

□

Let us now introduce the use of a finite precision arithmetic, with machine precision  $u$ , for approximating (2). Then, the best we can do is to consider the polynomial

approximation (11)–(12)<sup>5</sup>

$$\dot{y}(ch) \doteq \dot{\sigma}(ch) = \sum_{j=0}^{s-1} P_j(c)\gamma_j(\sigma), \quad c \in [0, 1], \tag{18}$$

such that

$$|\gamma_s(\sigma)| < \text{tol} \cdot \max_{j < s} |\gamma_j(\sigma)|, \quad \text{tol} \sim u. \tag{19}$$

Integrating (18), and imposing that  $\sigma(0) = y_0$ , then brings back to (13). We observe that because of (16), (19) may be approximately recast as

$$\sqrt{\frac{1}{s+1}} \rho^{-s} < \text{tol} \sim u, \tag{20}$$

where  $\rho \sim h^{-1}$ . Consequently, choosing  $s$  such that (19) (or (20)) is satisfied, we obtain that:

- the polynomial  $\sigma(ch)$  defined by (18) and (13) provides a uniformly accurate approximation to  $y(ch)$ , in the whole interval  $[0, h]$ , within the possibility of the used finite precision arithmetic;
- $\sigma(h)$  is a *spectrally accurate* approximation to  $y(h)$ . Moreover, in light of the second point of the result of Theorem 2, one has that the criterion (19) can be conveniently relaxed. In fact, making the ansatz (see (16) and (17))  $\kappa = \kappa_1$ , one has that

$$|\sigma(h) - y(h)| \lesssim \frac{h\kappa^2}{1 - \rho^{-2}} \rho^{-2s} \approx \frac{h(s+1)}{1 - \rho^{-2}} |\gamma_s(\sigma)|^2. \tag{21}$$

Imposing the approximate upper bound to be smaller than the machine epsilon  $u$ , one then obtains:

$$|\gamma_s(\sigma)| \lesssim \sqrt{\frac{u(1 - \rho^{-2})}{h(s+1)}} \propto u^{1/2}, \tag{22}$$

which is generally much less restrictive than (19).<sup>6</sup> Alternatively, by considering that the use of relatively large time-steps  $h$  is sought, one can use  $\text{tol} \sim u^{1/2}$  in (19), that is,

$$|\gamma_s(\sigma)| < \text{tol} \cdot \max_{j < s} |\gamma_j(\sigma)|, \quad \text{tol} \sim u^{1/2}, \tag{23}$$

In other words, (21) means that the method maintains the property of *super-convergence*, which is known to hold when  $h \rightarrow 0$ , also in the case where the time-step  $h$  is relatively large.

<sup>5</sup>Hereafter,  $\doteq$  means “equal within the round-off error level of the used finite precision arithmetic”.

<sup>6</sup>This latter criterion was that used in [18] and [9].



*Remark 2* In particular, we observe that (19) (or (20) or (22)) can be fulfilled by varying the value of  $s$ , and/or that of the stepsize  $h$ , by considering that, by virtue of (8),

$$\rho(h_{\text{new}}) \approx \rho(h_{\text{old}}) \frac{h_{\text{old}}}{h_{\text{new}}},$$

$h_{\text{old}}$  and  $h_{\text{new}}$  being the old and new stepsizes, respectively.

It is worth mentioning that the result of Theorem 2 can be also used to define a stepsize variation, within a generic error tolerance  $\text{tol}$ , thus defining a strategy for the simultaneous order/stepsize variation.

We conclude this section mentioning that, to gain efficiency, the criterion (19) for the choice of  $s$  in (18) can be more conveniently changed to

$$|\gamma_{s-1}(\sigma)| \leq \text{tol} \cdot \max_{j \leq s-1} |\gamma_j(\sigma)|, \quad \text{tol} \sim \rho \cdot u. \tag{24}$$

Similarly, the less restrictive criterion (22) can be approximately modified as:

$$|\gamma_{s-1}(\sigma)| \lesssim \rho \sqrt{\frac{u}{hs}} \propto u^{1/2},$$

or, alternatively, one uses  $\text{tol} \sim u^{1/2}$  in (24). As is clear, computing the norms of the coefficients  $\gamma_j(\sigma)$  permits to derive estimates for the parameters  $\kappa$  and  $\rho$  in (16), as we shall see later in the numerical tests.

### 3 SHBVMs

The approximation procedure studied in the previous section does not yet provide a numerical method, in that the integrals defining  $\gamma_j(\sigma)$ ,  $j = 0, \dots, s - 1$ , in (12)–(13) need to be computed. For this purpose, one can approximate them to within machine precision through a Gauss-Legendre quadrature formula of order  $2k$  (i.e., the interpolatory quadrature rule defined at the zeros of  $P_k$ ) with  $k$  large enough. In particular, following the criterion used in [9, 18], for the double precision IEEE<sup>7</sup>, we choose

$$k = \max\{20, s + 2\}. \tag{25}$$

After that, we define the approximation to  $y(h)$  as

$$y_1 := \sigma(h) \equiv y_0 + h\gamma_0(\sigma). \tag{26}$$

In so doing, one eventually obtains a HBVM( $k, s$ ), which we sketch below. Hereafter, we shall refer to such a method as to *spectral HBVM (in short, SHBVM)*, since its parameters  $s$  and  $k$ , respectively defined in (19) (or (20) or (22)) and (25), are aimed at obtaining a numerical solution which is accurate within the round-off error level of the used finite precision arithmetic.

<sup>7</sup>In such a case, the machine precision is  $u \approx 10^{-16}$ .

For sake of completeness, let us now briefly sketch what a HBVM( $k, s$ ) is. In general, to approximate the Fourier coefficient  $\gamma_j(\sigma)$ , and assuming for sake of simplicity that full machine accuracy is gained, we use the quadrature

$$\gamma_j(\sigma) \doteq \sum_{\ell=1}^k b_\ell P_j(c_\ell) f(\sigma(c_\ell h)) =: \hat{\gamma}_j, \quad j = 0, \dots, s - 1, \quad (27)$$

where the polynomial  $\sigma$  is that defined in (13) by formally replacing  $\gamma_j(\sigma)$  with  $\hat{\gamma}_j$ , and  $(c_i, b_i)$  are the abscissae and weights of the Gauss-Legendre quadrature of order  $2k$  on the interval  $[0, 1]$ .<sup>8</sup> Setting  $Y_\ell = \sigma(c_\ell h)$ , from (27), one then obtains the *stage equations*

$$\begin{aligned} Y_i &= y_0 + h \sum_{j=0}^{s-1} \int_0^{c_i} P_j(x) dx \hat{\gamma}_j \\ &\equiv y_0 + h \sum_{j=1}^k b_j \underbrace{\left[ \sum_{\ell=0}^{s-1} \int_0^{c_i} P_\ell(x) dx P_\ell(c_j) \right]}_{=: a_{ij}} f(Y_j), \quad i = 1, \dots, k, \end{aligned} \quad (28)$$

with the new approximation given by (see (26))

$$y_1 = y_0 + h \hat{\gamma}_0 \equiv y_0 + h \sum_{i=1}^k b_i f(Y_i). \quad (29)$$

Consequently, with reference to (28), setting

$$A = (a_{ij}) \in \mathbb{R}^{k \times k}, \quad \mathbf{b} = (b_i), \mathbf{c} = (c_i) \in \mathbb{R}^k, \quad (30)$$

one easily realizes that (28) and (29) define the  $k$ -stage Runge-Kutta method with Butcher tableau:

$$\begin{array}{c|c} \mathbf{c} & A \\ \hline & \mathbf{b}^\top \end{array}.$$

From (28), one verifies that the Butcher matrix in (30) can be written as

$$A = \mathcal{I}_s \mathcal{P}_s^\top \Omega, \quad (31)$$

with

$$\mathcal{P}_s = \begin{pmatrix} P_0(c_1) & \dots & P_{s-1}(c_1) \\ \vdots & & \vdots \\ P_0(c_k) & \dots & P_{s-1}(c_k) \end{pmatrix}, \quad \mathcal{I}_s = \begin{pmatrix} \int_0^{c_1} P_0(x) dx & \dots & \int_0^{c_1} P_{s-1}(x) dx \\ \vdots & & \vdots \\ \int_0^{c_k} P_0(x) dx & \dots & \int_0^{c_k} P_{s-1}(x) dx \end{pmatrix} \in \mathbb{R}^{k \times s}, \quad (32)$$

and

$$\Omega = \begin{pmatrix} b_1 & & \\ & \ddots & \\ & & b_k \end{pmatrix} \in \mathbb{R}^{k \times k}. \quad (33)$$

<sup>8</sup>i.e.,  $0 < c_1 < \dots < c_k < 1$  are the zeros of  $P_k$ .

In fact, setting  $e_i \in \mathbb{R}^k$  the  $i$ th unit vector, and taking into account (31)–(33), one has

$$e_i^\top A e_j = e_i^\top \mathcal{I}_s \mathcal{P}_s^\top \Omega e_j = e_i^\top \mathcal{I}_s \left( e_j^\top \mathcal{P}_s \right)^\top b_j = b_j \left[ \sum_{\ell=0}^{s-1} \int_0^{c_i} P_\ell(x) dx P_\ell(c_j) \right] \equiv a_{ij},$$

as defined in (28). From well-known properties of Legendre polynomials (see, e.g., [7, Appendix A]), one has that

$$\mathcal{I}_s = \mathcal{P}_{s+1} \hat{X}_s \equiv \begin{pmatrix} P_0(c_1) & \dots & P_s(c_1) \\ \vdots & & \vdots \\ P_0(c_k) & \dots & P_s(c_k) \end{pmatrix} \begin{pmatrix} \xi_0 & -\xi_1 & & & \\ \xi_1 & 0 & \ddots & & \\ & \ddots & \ddots & & -\xi_{s-1} \\ & & & \xi_{s-1} & 0 \\ & & & & \xi_s \end{pmatrix}, \quad \xi_i = \left( 2\sqrt{|4i^2 - 1|} \right)^{-1}, \tag{34}$$

from which one easily derives the following property relating the matrices (32)–(33) (see, e.g., [7, Lemma 3.6]):

$$\mathcal{P}_s^\top \Omega \mathcal{I}_s = \begin{pmatrix} \xi_0 & -\xi_1 & & & \\ \xi_1 & 0 & \ddots & & \\ & \ddots & \ddots & & -\xi_{s-1} \\ & & & \xi_{s-1} & 0 \end{pmatrix} =: X_s \in \mathbb{R}^{s \times s}. \tag{35}$$

*Remark 3* From (32)–(34), one has that the Butcher matrix (31) can be rewritten as

$$A = \mathcal{P}_{s+1} \hat{X}_s \mathcal{P}_s^\top \Omega. \tag{36}$$

Considering that, when  $k = s$ , (see (35))  $\mathcal{P}_{s+1} \hat{X}_s = \mathcal{P}_s X_s$  and  $\mathcal{P}_s^\top \Omega = \mathcal{P}_s^{-1}$ , so that  $A$  reduces to  $\mathcal{P}_s X_s \mathcal{P}_s^{-1}$ , we observe that (36) can be also regarded as a generalization of the  $W$ -transformation in [25, Section IV.5].

At this point, we observe that the stage equation (28) can be cast in vector form, by taking into account (30)–(33), as

$$Y \equiv \begin{pmatrix} Y_1 \\ \vdots \\ Y_k \end{pmatrix} = e \otimes y_0 + h \mathcal{I}_s \mathcal{P}_s^\top \Omega \otimes I_m \cdot f(Y), \quad e = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \in \mathbb{R}^k, \tag{37}$$

with an obvious meaning of  $f(Y)$ . On the other hand, the block vector of the coefficients in (27) turns out to be given by

$$\hat{y} \equiv \begin{pmatrix} \hat{y}_0 \\ \vdots \\ \hat{y}_{s-1} \end{pmatrix} = \mathcal{P}_s^\top \Omega \otimes I_m \cdot f(Y). \tag{38}$$

Consequently, from (37), one obtains

$$Y = e \otimes y_0 + h \mathcal{I}_s \otimes I_m \cdot \hat{y},$$

and then, from (38), one eventually derives the equivalent discrete problem

$$F(\hat{y}) := \hat{y} - \mathcal{P}_s^\top \Omega \otimes I_m \cdot f(\mathbf{e} \otimes y_0 + h\mathcal{I}_s \otimes I_m \cdot \hat{y}) = \mathbf{0}, \tag{39}$$

which has (block) dimension  $s$ , independently of  $k$  (compare with (37)). Once it has been solved, the new approximation is obtained (see (29)) as  $y_1 = y_0 + h\hat{y}_0$ .

It is worth observing that the new discrete problem (39), having block dimension  $s$  independently of  $k$ , allows us to use arbitrarily high-order quadratures (see (25)), without affecting that much the computational cost.

In order to solve (39), one could in principle use a fixed-point iteration,<sup>9</sup>

$$\hat{y}^{\ell+1} := \mathcal{P}_s^\top \Omega \otimes I_m \cdot f(\mathbf{e} \otimes y_0 + h\mathcal{I}_s \otimes I_m \cdot \hat{y}^\ell), \quad \ell = 0, 1, \dots,$$

which, though straightforward, usually implies restrictions on the choice of the step-size  $h$ . For this reason, this approach is generally not useful when using the methods as spectral methods, where the use of relatively large stepsizes is sought. On the other hand, the use of the simplified Newton iteration for solving (39) reads, by virtue of (35),

$$\text{solve: } [I_s \otimes I_m - hX_s \otimes f'(y_0)] \delta^\ell = -F(\hat{y}^\ell), \quad \hat{y}^{\ell+1} := \hat{y}^\ell + \delta^\ell, \quad \ell = 0, 1, \dots \tag{40}$$

However, the coefficient matrix in (40) has a dimension  $s$  times larger than that of the continuous problem (i.e.,  $m$ ) and, therefore, this can be an issue when large values of  $s$  are to be used, as in the case of SHBVMs. Fortunately, this problem can be overcome by replacing the previous iteration (40) with a corresponding *blended iteration* [7, 8, 12] (see also [5]). In more details, once one has formally computed the  $m \times m$  matrix

$$\Sigma = (I_m - h\rho_s f'(y_0))^{-1}, \quad \rho_s = \min_{\lambda \in \sigma(X_s)} |\lambda|, \tag{41}$$

where  $\sigma(X_s)$  denotes, as is usual, the spectrum of matrix  $X_s$ , one iterates:

$$\eta^\ell := F(\hat{y}^\ell), \quad \eta_1^\ell := \rho_s X_s^{-1} \otimes I_m \eta^\ell, \quad \hat{y}^{\ell+1} := \hat{y}^\ell + I_s \otimes \Sigma [\eta_1^\ell + I_s \otimes \Sigma (\eta^\ell - \eta_1^\ell)], \quad \ell = 0, 1, \dots \tag{42}$$

Consequently, one only needs to compute, at each time-step, the matrix  $\Sigma \in \mathbb{R}^{m \times m}$  defined in (41),<sup>10</sup> having the same size as that of the continuous problem. Moreover, it is worth mentioning that for semi-linear problems with a leading linear part, the Jacobian of  $f$  can be approximated with the (constant) linear part, so that  $\Sigma$  is computed once for all [6, 9, 18].

*Remark 4* It must be stressed that it is the availability of the very efficient blended iteration (41)–(42) which makes the practical use of HBVMs as spectral methods in time possible, since relatively large values of  $s$  can be easily and effectively handled.

A thorough analysis of the blended iteration can be found in [15]. Contexts where it has been successfully implemented include stiff ODE-IVPs [16], linearly implicit DAEs up to index 3 [17] (see also the code BiMD in TestSet for

<sup>9</sup>Hereafter, the initial approximation  $\hat{y}^0 = \mathbf{0}$  is conveniently used.

<sup>10</sup>i.e., factor  $\Sigma^{-1}$ .

IVP Solvers “<https://archimede.dm.uniba.it/~testset/testsetivpsolvers/>”), and canonical Hamiltonian systems (see the Matlab code HBVM, available at “<http://web.math.unifi.it/users/brugnano/LIMbook/>”), while its implementation in the solution of RKN methods may be found in [38].

## 4 Numerical tests

The aim of this section is twofold: firstly, to assess the theoretical analysis of SHBVMs made in Section 2; secondly, to compare such methods w.r.t. some well-known ones. All numerical tests, which concern different kinds of ODE problems, have been computed on a laptop with a 2.8-GHz Intel-i7 quad-core processor and 16GB of memory, running Matlab 2017b. For the SHBVM, the criteria (23) and (25) have been respectively used to determine its parameters  $s$  and  $\kappa$ .

### 4.1 The Kepler problem

We start considering the well-known Kepler problem (see, e.g., [7, Chapter 2.5]), which is Hamiltonian, with Hamiltonian function

$$H(q, p) = \frac{1}{2} \|p\|_2^2 - \|q\|_2^{-1}, \quad q, p \in \mathbb{R}^2. \tag{43}$$

Consequently, we obtain the equations

$$\dot{q} = p, \quad \dot{p} = -\|q\|_2^{-3}q, \tag{44}$$

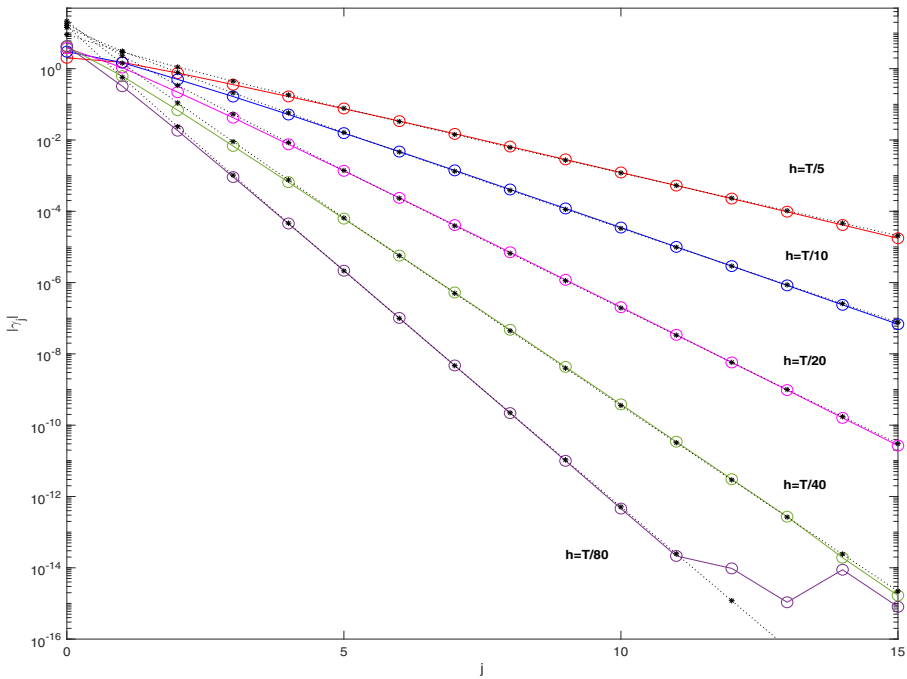
which, when coupled with the initial conditions

$$q(0) = (1 - \varepsilon, 0)^\top, \quad p(0) = \left(0, \sqrt{\frac{1+\varepsilon}{1-\varepsilon}}\right)^\top, \quad \varepsilon \in [0, 1), \tag{45}$$

provide a periodic orbit of period  $T = 2\pi$  that, in the  $q$ -plane, is given by an ellipse of eccentricity  $\varepsilon$ . In particular, we choose the value  $\varepsilon = 0.5$ . The solution of this problem has two additional (functionally independent) invariants besides the Hamiltonian (43), i.e., the *angular momentum* and one of the nonzero components of the Lenz vector [7, page 64] (in particular, we select the second one):

$$M(q, p) = q_1p_2 - p_1q_2, \quad L(q, p) = -p_1M(q, p) - q_2\|q\|_2^{-1}. \tag{46}$$

At first, we want to assess the result of Theorem 1. For this purpose, we apply the HBVM (20,16) for one step starting from the initial condition (45), and using time-steps  $h_i = 2\pi/(5 \cdot 2^{i-1})$ ,  $i = 1, \dots, 5$ . Figure 1 is the plot (see (27)) of  $|\hat{\gamma}_j|$ , for  $j = 0, 1, \dots, 15$ , (solid line with circles), which, according to (16), should behave as  $\kappa\rho^{-j}/\sqrt{j+1}$ , due to the result of Theorem 1. A least squares approximation technique has been employed to estimate the two parameters ( $\kappa$  and  $\rho$ ) appearing in the bound (16). These theoretical bounds are highlighted by the dashed line with asterisks in Fig. 1: evidently, they well fit the computed values, except those which are close to the round-off error level. Moreover, according to the arguments in the proof of Theorem 1, one also expects that the estimate of  $\kappa$  increases but is bounded,



**Fig. 1** Behavior of  $|\hat{y}_j|$  for decreasing values of the time-step  $h$  for the Kepler problem (44)–(45) solved by the HBVM(20,16) with decreasing time-steps. The line with circles are the computed norms, whereas those with the asterisks are the estimated ones. Observe that for the smallest time-steps, the computed norms stagnate near the round-off error level

as  $h \rightarrow 0$ , whereas  $\rho$  should be proportional to  $h^{-1}$ . This fact is confirmed by the results listed in Table 1.

Next, we compare the following methods for solving (44)–(45):

- the  $s$ -stage Gauss method (i.e., HBVM( $s, s$ )),  $s = 1, 2$ , which is symplectic and of order  $2s$ . Consequently, it is expected to conserve the angular momentum  $M(q, p)$  in (46), which is a quadratic invariant [34];
- the HBVM(6,  $s$ ) methods,  $s = 1, 2$ , which, for the considered stepsizes is energy-conserving and of order  $2s$ ;

**Table 1** Estimated values for the parameters  $\kappa$  and  $\rho$  for the Kepler problem (44)–(45), when using decreasing time-steps

$h$	$\kappa$	$\rho$
$2\pi/5$	9.4	2.2
$2\pi/10$	14.3	3.2
$2\pi/20$	18.2	5.6
$2\pi/40$	21.5	10.6
$2\pi/80$	16.1	19.8

**Table 2** Numerical result for the  $s$ -stage Gauss method,  $s = 1, 2$ , used for solving the Kepler problem (44)–(45),  $\varepsilon = 0.5$ , with stepsize  $h = 2\pi/n$

$n$	Time	$e_H$	Rate	$e_M$	$e_L$	Rate	$e_y$	Rate
Gauss-1								
100	2.52	6.56e−03	–	5.88e−15	4.97e−01	–	3.04e 00	–
200	4.74	1.63e−03	2.0	1.04e−14	3.54e−01	0.5	2.39e 00	0.3
400	9.78	3.82e−04	2.1	2.09e−14	9.76e−02	1.9	1.61e 00	0.6
800	17.26	3.05e−05	3.6	7.66e−15	2.45e−02	2.0	7.49e−01	1.1
1600	33.74	6.07e−07	5.6	1.93e−14	6.11e−03	2.0	2.08e−01	1.8
3200	65.46	9.65e−09	6.0	3.04e−14	1.53e−03	2.0	5.25e−02	2.0
Gauss-2								
50	2.33	2.05e−06	–	3.44e−15	3.81e−02	–	3.17e−01	–
100	4.09	5.37e−10	11.9	5.77e−15	2.43e−03	4.0	2.09e−02	3.9
200	7.52	1.44e−13	11.9	7.55e−15	1.53e−04	4.0	1.32e−03	4.0
400	14.48	9.55e−15	3.9	9.99e−14	9.55e−06	4.0	8.29e−05	4.0
800	26.75	1.53e−14	***	1.49e−14	5.97e−07	4.0	5.18e−06	4.0
1600	52.46	3.81e−14	***	1.95e−14	3.73e−08	4.0	3.24e−07	4.0
3200	101.74	3.49e−14	***	4.71e−14	2.33e−09	4.0	2.04e−08	4.0

- the SHBVM method described above, where  $s$  and  $k$  are determined according to (23) and (25), respectively, with  $\text{tol} \approx 10^{-8}$ . This tolerance, in turn, should provide us with full accuracy, according to the result of Theorem 2, because of the super-convergence of the method, which is valid for any used step-size.<sup>11</sup>

It is worth mentioning that the execution times that we shall list for the Gauss, HBVM, and SHBVM methods are perfectly comparable, since the same Matlab code has been used for all of them. This code, in turn, is a slight variant of the `hbvm` function available at the url “<http://web.math.unifi.it/users/brugnano/LIMbook/>”.

In Tables 2, 3, and 4, we list the obtained results when using a time-step  $h = 2\pi/n$  over 100 periods. In more details, we list the maximum errors, measured at each period, in the invariants (43) and (46),  $e_H, e_M, e_L$ , respectively, the solution error,  $e_y$ , and the execution times (in sec). As it is expected, the symplectic methods conserve the angular momentum (since it is a quadratic invariant), whereas the energy-conserving HBVMs conserve the Hamiltonian function.<sup>12</sup> On the other hand, the SHBVM conserves all the invariants and has a uniformly small solution error, by using very large stepsizes. Further, its execution time is the lowest one (less than 0.5 sec, when using  $h = 2\pi/5$ ), thus confirming the effectiveness of the method.

<sup>11</sup>As matter of fact, considering more stringent tolerances does not improve the accuracy of the computed numerical solution.

<sup>12</sup>In this case, the Gauss methods exhibit a super-convergence in the conservation of the Hamiltonian (3 times the usual order) and HBVMs do the same with the angular momentum. This is due to the fact that the error is measured only at the end of each period.

**Table 3** Numerical result for the HBVM(6,  $s$ ) methods,  $s = 1, 2$ , used for solving the Kepler problem (44)–(45),  $\varepsilon = 0.5$ , with stepsize  $h = 2\pi/n$

$n$	Time	$e_H$	$e_M$	Rate	$e_L$	Rate	$e_y$	Rate
HBVM(6,1)								
100	3.55	4.44e−16	9.09e−04	–	4.99e−01	–	2.94e00	–
200	7.10	4.44e−16	2.12e−05	6.0	3.52e−01	1.9	9.68e−01	1.9
400	12.47	6.66e−16	3.39e−07	6.0	9.70e−02	2.0	2.58e−01	2.0
800	22.86	4.44e−16	5.29e−09	6.0	2.44e−02	2.0	6.46e−02	2.0
1600	45.46	4.44e−16	8.26e−11	6.0	6.10e−03	2.0	1.62e−02	2.0
3200	86.34	6.66e−16	1.30e−12	6.0	1.53e−03	2.0	4.04e−03	2.0
HBVM(6,2)								
50	2.92	4.44e−16	1.09e−07	–	3.82e−02	–	4.64e−02	–
100	4.50	4.44e−16	2.72e−11	12.0	2.43e−03	4.0	2.94e−03	4.0
200	8.10	4.44e−16	5.88e−15	12.1	1.53e−04	4.0	1.84e−04	4.0
400	15.48	4.44e−16	3.89e−15	***	9.55e−06	4.0	1.15e−05	4.0
800	28.42	4.44e−16	1.40e−14	***	5.97e−07	4.0	7.20e−07	4.0
1600	52.29	6.66e−16	1.73e−14	***	3.73e−08	4.0	4.50e−08	4.0
3200	107.41	6.66e−16	1.40e−14	***	2.33e−09	4.0	2.81e−09	4.0

### 4.2 A Lotka-Volterra problem

We consider the following Poisson problem [20],

$$\dot{y} = B(y)\nabla H(y), \quad B(y)^\top = -B(y), \tag{47}$$

with  $y \in \mathbb{R}^3$  and, for arbitrary real constants  $a, b, c, v, \mu$ ,

$$B(y) = \begin{pmatrix} 0 & c y_1 y_2 & bc y_1 y_3 \\ -c y_1 y_2 & 0 & -y_2 y_3 \\ -bc y_1 y_3 & y_2 y_3 & 0 \end{pmatrix}, \quad H(y) = ab y_1 + y_2 - a y_3 + v \ln y_2 - \mu \ln y_3. \tag{48}$$

Moreover, assuming that  $abc = -1$ , there is a further invariant besides the Hamiltonian  $H$ , i.e., the *Casimir*

$$C(y) = ab \ln y_1 - b \ln y_2 + \ln y_3. \tag{49}$$

**Table 4** Numerical result for the SHBVM method used for solving the Kepler problem (44)–(45),  $\varepsilon = 0.5$ , with stepsize  $h = 2\pi/n$

$n$	$k$	$s$	Time	$e_H$	$e_M$	$e_L$	$e_y$
5	24	22	0.47	4.44e−16	2.01e−14	1.66e−14	8.00e−13
10	20	16	0.71	4.44e−16	6.22e−15	2.34e−14	6.13e−13
20	20	11	1.22	4.44e−16	6.66e−16	3.89e−15	3.87e−13
40	20	9	2.16	2.22e−16	1.89e−15	3.28e−15	5.75e−13



The solution turns out to be periodic, with period  $T \approx 2.878130103817$ , when choosing

$$a = -2, \quad b = -1, \quad c = -0.5, \quad \nu = 1, \quad \mu = 2, \quad y(0) = (1, 1.9, 0.5)^\top. \tag{50}$$

For this problem, the HBVM( $k, s$ ) is no more energy-conserving, as well as the  $s$ -stage Gauss method. As matter of fact, both exhibit a drift in the invariants and a quadratic error growth in the numerical solution. The obtained results for the SHBVM, with  $\text{tol} \approx 10^{-8}$  in (23) for choosing  $s, \kappa$  given by (25), and using a step-size  $h = T/n$ , are listed in Table 5, where it is reported the maximum Hamiltonian error,  $e_H$ , the Casimir error,  $e_C$ , and the solution error  $e_y$ , measured at each period, over 100 periods. In such a case, all the invariants turn out to be numerically conserved, and the solution error is uniformly very small. Moreover, the SHBVM using the largest time-step (i.e.,  $h = T/5 \approx 0.57$ ) turns out to be the most efficient one. For comparison, in the table, we also list the results obtained by using the Matlab solver `ode45` used with the default parameters, requiring 5600 integration steps and step-sizes approximately in the range  $[2.2 \cdot 10^{-2}, 1.1 \cdot 10^{-1}]$ , and the same solver used with parameters  $\text{AbsTol}=1e-15, \text{RelTol}=1e-10$ , now requiring 121760 integration steps, with stepsizes approximately in the range  $[10^{-3}, 4.2 \cdot 10^{-3}]$ .

### 4.3 A stiff ODE-IVP

At last, we consider a stiff ODE-IVP,

$$\dot{y}(t) = \begin{pmatrix} -9999 & 1 & 1 \\ 9900 & -100 & 1 \\ 98 & 98 & -2 \end{pmatrix} [y(t) - g(t)] + \dot{g}(t), \quad y(0) = g(0), \tag{51}$$

with  $g(t)$  a known function, having evidently solution  $y(t) = g(t)$ . We choose

$$g(t) = (\cos 2\pi t, \cos 4\pi t, \cos 6\pi t)^\top, \tag{52}$$

and consider the SHBVM with  $\text{tol} \approx 10^{-8}$  in (23) for choosing  $s$  (as before,  $\kappa$  is chosen according to (25)), so that full accuracy is expected in the numerical solution.

**Table 5** Numerical result for the SHBVM method used for solving the Lotka-Volterra problem (47)–(50) with stepsize  $h = T/n$

$n$	$k$	$s$	Time	$e_H$	$e_C$	$e_y$
5	20	16	0.88	8.26e-14	4.89e-14	4.24e-11
10	20	11	1.37	1.33e-14	1.33e-14	5.01e-11
15	20	9	1.84	3.11e-14	1.62e-14	4.92e-11
<code>ode45</code>			0.23	7.41e-01	7.27e-01	3.62e00
<code>ode45*</code>			4.12	1.14e-08	8.71e-09	8.44e-07

We also list the results obtained by using `ode45`, both with the default parameters, and with parameters  $\text{AbsTol}=1e-15, \text{RelTol}=1e-10$  (which we denote by `ode45*`)

**Table 6** Numerical result for the SHBVM method used for solving the stiff problem (51)–(52) with stepsize  $h = 100/n$ , and `ode15s` with the default parameters

$n$	$k$	$s$	Time	$e_y$
50	40	38	0.09	2.92e-11
75	32	30	0.12	1.53e-11
100	28	26	0.17	1.93e-12
125	25	23	0.21	6.28e-12
150	22	20	0.27	9.43e-12
<code>ode15s</code>			0.68	3.76e-04

The time-step used is  $h = 100/n$  for  $n$  steps. The measured errors in the last point (coinciding with the initial condition), are then reported in Table 6. For comparison, also the results obtained by the Matlab solver `ode15s`, using its default parameters, are listed in the table. This latter solver requires 6006 steps, with time-steps approximately in the range  $[1.9 \cdot 10^{-3}, 2 \cdot 10^{-2}]$ .

## 5 Conclusions

In this paper, we provide a thorough analysis of SHBVMs, namely HBVMs used as spectral methods in time, which further confirms their effectiveness. From the analysis, one obtains that the super-convergence of HBVMs is maintained also when using relatively large time-steps. SHBVMs become a practical method, due to the very efficient nonlinear *blended* iteration inherited from HBVMs. As a consequence, SHBVMs appear to be good candidates as *general ODE solvers*. This is indeed confirmed by a few numerical tests concerning a Hamiltonian problem, a Poisson (not Hamiltonian) problem, and a stiff ODE-IVP. The same tests show the numerical assessment of the theoretical achievements.

**Acknowledgments** The authors are very grateful to two unknown referees, for the careful reading of the manuscript, and for their precious comments, which allowed to formulate in a cleaner and more precise way the results presented in the paper.

## References

1. Amodio, P., Brugnano, L., Iavernaro, F.: Energy-conserving methods for Hamiltonian Boundary Value Problems and applications in astrodynamics. *Adv. Comput. Math.* **41**, 881–905 (2015)
2. Betsch, P., Steinmann, P.: Inherently energy conserving time finite elements for classical mechanics. *J. Comp. Phys.* **160**, 88–116 (2000)
3. Bottasso, C.L.: A new look at finite elements in time: a variational interpretation of Runge–Kutta methods. *Appl. Numer. Math.* **25**, 355–368 (1997)
4. Brugnano, L., Calvo, M., Montijano, J.I., Rández, L.: Energy preserving methods for Poisson systems. *J. Comput. Appl. Math.* **236**, 3890–3904 (2012)
5. Brugnano, L., Frasca Caccia, G., Iavernaro, F.: Efficient implementation of Gauss collocation and Hamiltonian Boundary Value Methods. *Numer. Algorithm.* **65**, 633–650 (2014)
6. Brugnano, L., Frasca-Caccia, G., Iavernaro, F.: Line Integral Solution of Hamiltonian PDEs. *Mathematics* 7(3), article n. 275. <https://doi.org/10.3390/math7030275> (2019)

7. Brugnano, L., Iavernaro, F.: Line Integral Methods for Conservative Problems. Chapman and Hall/CRC, Boca Raton (2016)
8. Brugnano, L., Iavernaro, F.: Line Integral Solution of Differential Problems. *Axioms* **7**(2), article n. 36. <https://doi.org/10.3390/axioms7020036> (2018)
9. Brugnano, L., Iavernaro, F., Montijano, J.I., Rández, L.: Spectrally accurate space-time solution of Hamiltonian PDEs. <https://doi.org/10.1007/s11075-018-0586-z> (2018)
10. Brugnano, L., Iavernaro, F., Trigiante, D.: Hamiltonian BVMs (HBVMs): A family of “drift-free” methods for integrating polynomial Hamiltonian systems. *AIP Conf. Proc.* **1168**, 715–718 (2009)
11. Brugnano, L., Iavernaro, F., Trigiante, D.: Hamiltonian boundary value methods (energy preserving discrete line integral methods). *JNAIAM J. Numer. Anal. Ind. Appl. Math.* **5**(1-2), 17–37 (2010)
12. Brugnano, L., Iavernaro, F., Trigiante, D.: A note on the efficient implementation of Hamiltonian BVMs. *J. Comput. Appl. Math.* **236**, 375–383 (2011)
13. Brugnano, L., Iavernaro, F., Trigiante, D.: A simple framework for the derivation and analysis of effective one-step methods for ODEs. *Appl. Math. Comput.* **218**, 8475–8485 (2012)
14. Brugnano, L., Iavernaro, F., Trigiante, D.: A two-step, fourth-order method with energy preserving properties. *Comput. Phys. Commun.* **183**, 1860–1868 (2012)
15. Brugnano, L., Magherini, C.: Blended Implementation of Block Implicit Methods for ODEs. *Appl. Numer. Math.* **42**, 29–45 (2002)
16. Brugnano, L., Magherini, C.: The BiM Code for the Numerical Solution of ODEs. *J. Comput. Appl. Math.* **164–165**, 145–158 (2004)
17. Brugnano, L., Magherini, C., Mugnai, F.: Blended implicit methods for the numerical solution of DAE problems. *J. Comput. Appl. Math.* **189**, 34–50 (2006)
18. Brugnano, L., Montijano, J.I., Rández, L.: On the effectiveness of spectral methods for the numerical solution of multi-frequency highly-oscillatory Hamiltonian problems. <https://doi.org/10.1007/s11075-018-0552-9> (2018)
19. Brugnano, L., Sun, Y.: Multiple invariants conserving Runge-Kutta type methods for Hamiltonian problems. *Numer. Algorithm.* **65**, 611–632 (2014)
20. Cohen, D., Hairer, E.: Linear energy-preserving integrators for Poisson systems. *BIT* **51**, 91–101 (2011)
21. Davis, P.J.: *Interpolation & Approximations*. Dover, New York (1975)
22. Franco, J.M.: Exponentially fitted symplectic integrators of RKN type for solving oscillatory problems. *Comput. Phys. Comm.* **177**(6), 479–492 (2007)
23. Franco, J.M.: Runge-kutta methods adapted to the numerical integration of oscillatory problems. *Appl. Numer. Math.* **50**, 427–443 (2004)
24. Gautschi, W.: Numerical integration of ordinary differential equations based on trigonometric polynomials. *Numer. Math.* **3**, 381–397 (1961)
25. Hairer, E., Wanner, G.: *Solving Ordinary Differential Equations II*, 2nd revised edition. Springer, Heidelberg (2002)
26. Hoang, N.S., Sidje, R.B., Cong, N.H.: On functionally-fitted Runge-Kutta methods. *BIT* **46**, 861–874 (2006)
27. Hulme, B.L.: One-step piecewise polynomial Galerkin methods for initial value problems. *Math. Comp.* **26**, 415–426 (1972)
28. Hulme, B.L.: Discrete Galerkin and related one-step methods for ordinary differential equations. *Math. Discret. Comp.* **26**, 881–891 (1972)
29. Martin-Vaquero, J., Vigo-Aguiar, J.: Exponential fitted Gauss, Radau and Lobatto methods of low order. *Numer. Algorithm.* **48**(4), 327–346 (2008)
30. Montijano, J.I., Van Daele, M., Calvo, M.: Functionally fitted explicit two step peer methods. *J. Sci. Comput.* **64**(3), 938–958 (2015)
31. Kalogiratou, Z., Simos, T.E.: Construction of trigonometrically and exponentially fitted Runge-Kutta-Nyström methods for the numerical solution of the Schrödinger equation and related problem - a method of 8th algebraic order. *J. Math. Chem.* **31**(2), 211–232 (2002)
32. Li, Y.-W., Wu, X.: Functionally fitted energy-preserving methods for solving oscillatory nonlinear Hamiltonian systems. *SIAM J. Numer. Anal.* **54**(4), 2036–2059 (2016)
33. Paternoster, B.: Runge-kutta (-nyström) methods for ODEs with periodic solutions based on trigonometric polynomials. *Appl. Numer. Math.* **28**, 401–412 (1998)
34. Sanz-Serna, J.M.: Runge-kutta schemes for Hamiltonian systems. *BIT* **28**(4), 877–883 (1988)

35. Simos, T.E.: A trigonometrically-fitted method for long-time integration of orbital problems. *Math. Comput. Modell.* **40**(11-12), 1263–1272 (2004)
36. Tang, W., Sun, Y.: Time finite element methods: a unified framework for numerical discretizations of ODEs. *Appl. Math. Comp.* **219**, 2158–2179 (2012)
37. Vanden Berghe, G., De Meyer, H., Van Daele, M., Van Hecke, T.: Exponentially-fitted explicit Runge-Kutta methods. *Comput. Phys. Commun.* **123**, 7–15 (1999)
38. Wang, B., Meng, F., Fang, Y.: Efficient implementation of RKN-type Fourier collocation methods for second-order differential equations. *Appl. Numer. Math.* **119**, 164–178 (2017)

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.