# Gautschi method without order reduction when integrating boundary value nonlinear wave problems

B. Cano *

IMUVA, Departamento de Matemática Aplicada,

Facultad de Ciencias, Universidad de Valladolid,

Paseo de Belén 7, 47011 Valladolid,

Spain

AND

M. J. Moreta$^{\dagger}$

IMUVA, Departamento de Fundamentos del Análisis Económico I,

Facultad de Ciencias Económicas y Empresariales, Universidad Complutense de Madrid,

Campus de Somosaguas, Pozuelo de Alarcón, 28223 Madrid,

Spain.

**Abstract**

In this paper we analyse the order reduction which turns up when integrating nonlinear wave problems with non-homogeneous and time-dependent boundary conditions with the well-known Gautschi method. Moreover, a technique is suggested to avoid that order reduction so that the classical local order 4 and global order 2 are recovered. On the other hand, the usual approximation for the derivative which is used together with this method is also analysed and a substantial improvement is suggested. Some numerical results are shown which corroborate the performed analysis.

**Keywords:** Gautschi method, initial boundary value problem, nonlinear wave equations, avoiding order reduction

**AMS codes**: 65M12 65M20

# 1   Introduction

With the recent improvement of Krylov-type methods, exponential integrators [22] have become a valuable tool to integrate partial differential equations because the linear part

---

*Corresponding author. Email: bego@mac.uva.es

$^{\dagger}$Email: mjesusmoreta@ccee.ucm.es

is integrated exactly and no stability problems turn up. In this paper we concentrate on second-order differential problems in time, for which trigonometric integrators [21] or multistep cosine methods [7] are suitable exponential integrators. More particulary, we concentrate on the simplest of these, Gautschi method [16].

On the other hand, in the recent literature, an analysis has been done on the order reduction which turns up when integrating first-order initial boundary value problems with non-homogeneous and non-periodic boundary conditions with the aid of several exponential integrators of Lawson or splitting type [2, 14]. Moreover, several techniques have been designed to avoid it [3, 4, 5, 8, 9, 12, 13] for both linear and nonlinear problems. In this paper we aim to do the same with Gautschi method, which is so useful when approximating the solution of wave problems. This type of problems are very interesting in practice and, therefore, it is important to get an approximation of them as accurate as possible. Besides, the abstract framework is necessarily different from other problems which have already been considered in the literature as well as the type of functions which turn up when modifying the method in order to avoid the order reduction.

We tackle directly the nonlinear problem for the sake of completeness. Moreover, although the standard Krylov methods may fail to converge for trigonometric functions of unbounded operators, there are efficient implementations to approximate them in the set of rational Krylov subspace methods [19].

In this paper we consider an analysis of the full discretization of the problem, both in space and time. With the classical approach, the space discretization is firstly done and, secondly, the time integration. In such a case, if periodic boundary conditions are considered, second-order in time is observed uniformly on the space grid [18]. For that, the only assumption which is required on the exact solution is that its energy is bounded. However, when the boundary conditions are not periodic but are general and time-dependent, order reduction is observed, as it is fully explained in this paper. In order to avoid that, we suggest in this paper a technique which works with the only condition that the exact solution is regular enough, which is a typical requirement if a certain accuracy can be achieved. This technique is based on integrating firstly in time determining some of the terms which define the method through some differential problems with appropriate boundary conditions which can be calculated in terms of data. Then, the space discretization is performed to approximate the solution of those new differential problems. In any case, for those who are more interested in applying the modified method than on the analysis, the final formula to be considered is written in (28) (c.f. with (11) for the classical approach). Functions $\gamma_1$ and $\gamma_2$ turn up there, which can be both calculated in terms of a cosine function according to their definition in (7) and (24).

On the other hand, although Gautschi original method just approximates the solution [16], it is usual to consider also an approximation for the derivative which is based on the differentiation of the variation-of-constants formula which defines Gautschi method [17]. In this paper we also analyse the error for the derivative with the classical approach and suggest a modification to improve it (see formulas (32) and (39)). In such a way, we obtain error bounds of the same type than those obtained in [15] for

trigonometric integrators under periodic boundary conditions. More precisely, the ones corresponding to $s = 0$ and $\alpha = 1$ in Theorem 2.1 of that paper.

For the sake of clarity, we firstly assume that the space discretization satisfies some hypotheses which simplify the proofs of theorems. In such a way, simple finite differences are for example included and, in fact, we present our numerical results in a simple problem with them. Nevertheless, the technique to avoid the order reduction is also valid with other more general space discretizations which include, for example, finite elements, that are very suitable in more complicated problems. For that case, in an appendix, we give the final formula to be implemented in (47) and also the extension of the corresponding proofs.

Therefore, we structure the paper as follows. In the preliminaries we present the problem at hand in an abstract setting of Hilbert spaces, we remind how Gautschi method is deduced in an ordinary differential system framework and describe the first assumptions which are made on the space discretization. In Section 3, the classical approach is described and results on the local and global error for the approximation of the solution are given which justify the order reduction which is observed; just order 2 for the local error and hardly 1 for the global one. Then, in Section 4, the technique to avoid order reduction is suggested and the corresponding theorems are stated, which prove that the local order 4 and global order 2 in time are achieved. In Section 5, the analysis for the approximation of the derivative is performed for both the classical and the suggested approach. In Section 6 a numerical experiment is shown which corroborates the results of previous sections. Moreover, although it is not an aim of the paper itself, a comparison between the classical and the suggested approach is also given in terms of CPU time for a particular problem and implementation. Finally, the appendix describes the suggested technique for more general abstract space discretizations.

# 2   Preliminaries

Let $H$ be a complex Hilbert space and let $A : D(A) \subset H \to H$ be a densely defined, closed linear operator. We consider the abstract initial boundary value problem of second-order in time

$$\begin{cases} \ddot{u}(t) & = & -Au(t) + f(t, u(t)), \quad t \in [0, T], \\ u(0) & = & u_0, \\ \dot{u}(0) & = & v_0, \\ \partial u(t) & = & g(t), \end{cases} \tag{1}$$

which we suppose to be well-posed for any smooth enough functions $f$, $u_0$, $v_0$ and $g$ satisfying the natural compatibility conditions (see for example Section 3 of Chapter 5 in [24]). Moreover, we assume that the linear operators $A$ and $\partial$ satisfy the following assumptions:

(A1) The boundary operator $\partial : D(A) \subset H \to Y$ is onto, where $Y$ is another Hilbert space.

(A2) $\mathrm{Ker}(\partial)$ is dense in $H$ and $-A_0 : D(A_0) = Ker(\partial) \subset H \to H$, the restriction of $-A$ to $\mathrm{Ker}(\partial)$, is such that its numerical range is contained in a parabola $\Pi_{\lambda,\alpha} = \lambda - \bar{P}_\alpha$, where $\lambda < 0$, $\alpha \geq 0$ and

$$\bar{P}_\alpha = \{x + iy : y \in \mathbb{R}, 4\alpha^2 x \geq y^2\}.$$

(A3) The steady state problem

$$\begin{aligned} Ax &= 0, \\ \partial x &= v, \end{aligned}$$

possesses a unique solution denoted by $x = K(0)v$. Moreover, there exists a constant $C$ such that the linear operator $K(0) : Y \to D(A)$ satisfies

$$\|K(0)v\|_H \leq C\|v\|_Y.$$

(A4) The function $f$ and the solution $u$ in (1) are such that $u \in C^2([0,T], D(A))$ and $f$ is globally Lipschitz.

**Remark 1.** *Hypothesis (A2) includes the case of negative definite selfadjoint operators, but also operators associated to a sesquilinear form in Hilbert spaces (see [1, 6, 11, 20]). Moreover, because of this hypothesis, using [1, 10, 11], for $\tilde{w} \in H$, the solution of*

$$\begin{aligned} \ddot{w}(\tau) &= -A_0 w(\tau), \quad 0 \leq \tau \leq T, \\ w(0) &= \tilde{w}, \\ \dot{w}(0) &= 0, \end{aligned} \qquad (2)$$

*which will be denoted by $\cos(\tau B_0)\tilde{w}$, satisfies that there exists a constant $C$ such that*

$$\|w(\tau)\| = \|\cos(\tau B_0)\tilde{w}\| \leq C\|\tilde{w}\|, \quad 0 \leq \tau \leq T.$$

*(This result comes from the fact that the cosine function is bounded in a band around the positive real axis and when $z \in \bar{\Pi}_{\lambda,\alpha}(\lambda < 0)$, $(-z)^{\frac{1}{2}}$ belongs to such a band.) Besides, this solution is understood in a generalized sense. In case $\tilde{w} \notin D(A_0)$, as $D(A_0)$ is dense in $H$, $\tilde{w}$ can be approximated by elements in $D(A_0)$ and we consider the limit of the corresponding solutions of (2).*

*On the other hand, the solution of*

$$\begin{aligned} \ddot{w}(\tau) &= -A_0 w(\tau), \quad 0 \leq \tau \leq T, \\ w(0) &= 0, \\ \dot{w}(0) &= \tilde{w}, \end{aligned}$$

*which will be denoted by $B_0^{-1}\sin(\tau B_0)\tilde{w}$, satisfies that there exists a constant $C$ such that*

$$\|\frac{1}{\tau}w(\tau)\| = \|sinc(\tau B_0)\tilde{w}\| \leq C\|\tilde{w}\|, \quad 0 \leq \tau \leq T.$$

*(Similarly, this result comes from the fact that the sinc function is bounded in a band around the positive real axis.)*

We will concentrate on Gautschi method [16] to time integrate (1). When applied to a finite-dimensional nonlinear problem like

$$\ddot{U}(t) = -B^2 U(t) + F(t, U(t)), \tag{3}$$

where $B$ is a regular matrix, this method comes from considering the variation-of-constants formula

$$U(t_n + k) = \cos(kB)U(t_n) + k\mathrm{sinc}(kB)\dot{U}(t_n)$$
$$+ \int_0^k B^{-1} \sin((k-s)B)F(t_n + s, U(t_n + s))ds, \tag{4}$$

which is summed with the one which is obtained by changing $k$ by $-k$. This leads to

$$U(t_n + k) - 2\cos(kB)U(t_n) + U(t_n - k)$$
$$= \int_0^k B^{-1} \sin((k-s)B)[F(t_n + s, U(t_n + s)) + F(t_n - s, U(t_n - s))]ds, \tag{5}$$

and substituting the term in brackets by $2F(t_n, U(t_n))$, the difference scheme which defines the method finally turns up:

$$U_{n+1} - 2\cos(kB)U_n + U_{n-1} = k^2 \gamma_1(kB)F_n, \tag{6}$$

where $F_n$ stands for $F(t_n, U_n)$ and

$$\gamma_1(\epsilon) = \frac{2(1 - \cos(\epsilon))}{\epsilon^2}. \tag{7}$$

For the sake of simplicity, in most of the paper we will consider a spatial discretization for which the nodal values on a certain grid are the key identifiers of the numerical discretization, and for which some assumptions can be made which quite simplify the analysis. This type of discretizations includes some finite differences and collocation spectral methods. Nevertheless, the analysis is also valid for other more general space discretizations which include other more sophisticated finite differences and finite elements. A thorough analysis which justifies that order reduction is also avoided in such a case with the technique suggested here will be given in the appendix. For the moment, in a similar way as in [5], we assume that, for each parameter $h$ in a sequence $\{h_j\}_{j=1}^\infty$ such that $h_j \to 0$, $X_h \subset H$ is a finite dimensional space which approximates $H$ when $h_j \to 0$ and the elements in $D(A_0)$ are approximated in a subspace $X_{h,0}$. All elements of $X_h$ (resp. $X_{h,0}$) are assumed to be determined by some nodal values in $\mathbb{C}^N$ (resp. $\mathbb{C}^{N_0}$) where $N$ (resp. $N_0$) depends on the dimension of the problem and the space grid $h$. The norm in $X_h$ is denoted by $\|\cdot\|_h$. The operator $A$ is then approximated by the matrix $A_h$, $A_0$ by $A_{h,0}$ and the solution of the elliptic problem

$$-Aw = F, \qquad \partial w = g, \tag{8}$$

5

is approximated by $(R_h w^T, Q_h g^T)^T \in \mathbb{C}^N$, where $R_h w \in \mathbb{C}^{N_0}$ is called the elliptic projection, $Q_h g \in \mathbb{C}^{N-N_0}$ discretizes the boundary values and the following is satisfied

$$-A_{h,0} R_h w - A_h Q_h g = P_h F, \qquad (9)$$

for a certain nodal projection $P_h : H \to \mathbb{C}^{N_0}$. We suppose that there exists some constant $C$ such that, for $w \in H$ and small enough $h$,

$$\|P_h w\|_h \le C \|w\|.$$

We also assume that the source function $f$ in (1) has sense as a function from $[0,T] \times \mathbb{C}^{N_0}$ on $\mathbb{C}^{N_0}$ and, for each $t \in [0,T]$ and $u \in H$,

$$P_h f(t, u) = f(t, P_h u).$$

Moreover, we will work under the following hypotheses:

(H1) There exists $B_{h,0}$ such that $B_{h,0}^2 = A_{h,0}$ and we assume that $A_{h,0}$ and $B_{h,0}$ are invertible and that their inverses can be bounded uniformly on $h$.

(H2) There exists a subspace $Z \subset D(A)$ with norm $\|\cdot\|_Z$ such that, for $u \in Z$,

$$\|A_{h,0}(R_h - P_h)u\|_h \le \varepsilon_h \|u\|_Z,$$

for $\varepsilon_h$ decreasing with $h$.

(H3) $\|A_{h,0}^{-1} A_h Q_h\|_h$ is bounded independently of $h$ for small enough $h$. Considering (9), this in fact corresponds to a discrete maximum principle, which would be simulating the continuous maximum principle which is satisfied because of (A3).

(H4) $\|\cos(\tau B_{h,0})\|_h$, $\|\sin(\tau B_{h,0})\|_h$, $\|\mathrm{sinc}(\tau B_{h,0})\|_h$, $\|e^{i\tau B_{h,0}}\|_h$ and $\|(I - e^{i\tau B_{h,0}})^{-1}\|_h$ are bounded independently of real positive and small enough $h$ and $-T \le \tau \le T$.

(H5) $f$ is globally Lipschitz as a function from $[0,T] \times \mathbb{C}^{N_0} \to \mathbb{C}^{N_0}$.

# 3 Classical approach: Discretizing firstly in space and then in time

In this section we will see that, although Gautschi method has local order 4 when applied to a non-stiff ordinary differential system where the source term is smooth enough, it does not happen the same when it is applied to the space discretization of a time-dependent boundary value problem like (1). In such a case, the following semidiscrete problem in $X_{h,0}$ arises:

$$
\begin{aligned}
\ddot{U}_h(t) &= -A_{h,0} U_h(t) - A_h Q_h g(t) + f(t, U_h(t)), \\
U_h(0) &= P_h u(0), \\
\dot{U}_h(0) &= P_h v(0),
\end{aligned}
\qquad (10)
$$

where $A_{h,0}$ and $A_h Q_h$ are not bounded with $h$. Then, Gautschi scheme reads as follows:

$$U_h^{n+2} - 2\cos(kB_{h,0})U_h^{n+1} + U_h^n = k^2\gamma_1(kB_{h,0})\big[-A_hQ_hg(t_{n+1}) + f(t_{n+1}, U_h^{n+1})\big]. \quad (11)$$

The following theorem states how the local truncation error $\rho_{h,n}$, which corresponds to integrate (10) with Gautschi scheme, behaves in such a case:

**Theorem 2.** *Under hypotheses (A1)-(A4) and (H1)-(H5), whenever $g \in C^2([0,T], Y)$, $f \in C^2([0,T] \times H, H)$ and $\dot{u} \in C([0,T], Z)$,*

$$\rho_{h,n} = O(k^2), \qquad B_{h,0}^{-1}\rho_{h,n} = O(k^3),$$

*where the constants in Landau notation are independent of $k$ and $h$.*

*Proof.* Because of the definition of local truncation error and of Gautschi method itself,

$$U_h(t_n + k) - 2\cos(kB_{h,0})U_h(t_n) + U_h(t_n - k)$$
$$= 2\int_0^k B_{h,0}^{-1}\sin((k-s)B_{h,0})[-A_hQ_hg(t_n) + f(t_n, U_h(t_n))]ds + \rho_{h,n}.$$

Comparing this with (5) substituting $U$ by $U_h$ and $F(t,U)$ by $-A_hQ_hg(t) + f(t,U)$, it turns out that

$$\rho_{h,n} = \int_0^k B_{h,0}^{-1}\sin((k-s)B_{h,0})\big[-A_hQ_h[g(t_n+s) + g(t_n-s) - 2g(t_n)]$$
$$+ f(t_n+s, U_h(t_n+s)) + f(t_n-s, U_h(t_n-s)) - 2f(t_n, U_h(t_n))\big]ds. \quad (12)$$

Now, by the regularity of $g$, the first term above can be written as

$$-\int_0^k B_{h,0}\sin((k-s)B_{h,0})\int_0^s A_{h,0}^{-1}A_hQ_h[\ddot{g}(t_n+\sigma) + \ddot{g}(t_n-\sigma)](s-\sigma)d\sigma ds$$
$$= -\cos((k-s)B_{h,0})\int_0^s A_{h,0}^{-1}A_hQ_h[\ddot{g}(t_n+\sigma) + \ddot{g}(t_n-\sigma)](s-\sigma)d\sigma\Big|_{s=0}^k$$
$$+ \int_0^k \cos((k-s)B_{h,0})\int_0^s A_{h,0}^{-1}A_hQ_h[\ddot{g}(t_n+\sigma) + \ddot{g}(t_n-\sigma)]d\sigma ds,$$
$$= O(k^2),$$

where we have used an integration-by-parts argument and hypotheses (H3) and (H4).

On the other hand, for the second term in (12), we firstly use that $U_h(t_n)$, $\dot{U}_h(t_n)$ and $\ddot{U}_h(t_n)$ are uniformly bounded with $h$, as we will justify now. Multiplying (1) by $P_h$, considering (9) and making the difference with (10),

$$P_h\ddot{u}(t) - \ddot{U}_h(t) = -A_{h,0}(P_hu(t) - U_h(t)) + f(t, P_hu(t)) - f(t, U_h(t))$$
$$+ A_{h,0}(P_hu(t) - R_hu(t)),$$
$$P_hu(0) - U_h(0) = 0,$$
$$P_h\dot{u}(0) - \dot{U}_h(0) = 0.$$

Therefore, using (4),

$$P_h u(t) - U_h(t)$$
$$= \int_0^t B_{h,0}^{-1} \sin((t-s)B_{h,0}) \big[ f(s, P_h u(s)) - f(s, U_h(s)) + A_{h,0}[P_h u(s) - R_h u(s)] \big] ds, \quad (13)$$

which implies, by using (H2),(H4),(H5) and that $u \in C([0,T], Z)$ that

$$\|P_h u(t) - U_h(t)\|_h \le L \int_0^t \|P_h u(s) - U_h(s)\| ds + T \varepsilon_h \max_{t \in [0,T]} \|u(t)\|_Z$$

for some constant $L$, and this Gronwall inequality is well-known to imply that

$$\|P_h u(t) - U_h(t)\|_h \le T \varepsilon_h \max_{t \in [0,T]} \|u(t)\|_Z e^{Lt}, \quad 0 \le t \le T, \quad (14)$$

which proves that $U_h(t)$ is bounded with $h$. Moreover, by differentiating (13),

$$\dot{U}_h(t) = P_h \dot{u}(t) - \int_0^t \cos((t-s)B_{h,0}) \big[ f(s, P_h u(s)) - f(s, U_h(s)) + A_{h,0}[P_h u(s) - R_h u(s)] \big] ds,$$

which is also uniformly bounded with $h$ and

$$
\begin{aligned}
\ddot{U}_h(t) &= P_h \ddot{u}(t) - f(t, P_h u(t)) + f(t, U_h(t)) - A_{h,0}[P_h u(t) - R_h u(t)] \\
&\quad + \int_0^t B_{h,0} \sin((t-s)B_{h,0}) \big[ f(s, P_h u(s)) - f(s, U_h(s)) + A_{h,0}[P_h u(s) - R_h u(s)] \big] ds \\
&= P_h \ddot{u}(t) - \cos(t B_{h,0}) A_{h,0}[P_h u(0) - R_h u(0)] \\
&\quad - \int_0^t \cos((t-s)B_{h,0}) \big[ f_s(s, P_h u(s)) + f_u(s, P_h u(s)) P_h \dot{u}(s) - f_s(s, U_h(s)) \\
&\qquad\qquad\qquad\qquad - f_u(s, U_h(s)) \dot{U}_h(s) + A_{h,0}[P_h \dot{u}(s) - R_h \dot{u}(s)] \big] ds,
\end{aligned}
$$

which is also uniformly bounded with $h$ considering again (H4), the regularity of $f$ and that $\dot{u} \in C([0,T], Z)$, which implies that (H2) can be applied in the last part of the previous formula. Because of all this, considering Taylor expansions in the second line of the bracket in (12), those terms are $O(s^2)$ and, together with the term outside the bracket, that part of the integral can be bounded by

$$O(k^3) \int_0^k \|\operatorname{sinc}((k-s)B_{h,0})\|_h ds = O(k^4),$$

from what the result follows and it becomes clear that, with the assumed regularity for $f$, the order reduction for the local truncation error comes just from the fact that $\ddot{g}$ does not vanish.

As for $B_{h,0}^{-1} \rho_{h,n}$, notice that the corresponding last line of (12) would also be $O(k^4)$ because of the boundedness of $\|B_{h,0}^{-1}\|$ according to (H1). As for the corresponding first line, considering (HS1), it would be directly

$$O(k^2) \int_0^k \sin((k-s)B_{h,0}) ds = O(k^3),$$

from what the result follows. $\qquad \square$

In order to study the global error, we firstly notice that (11) can be written as the explicit one-step method

$$
\begin{bmatrix} U_h^{n+1} \\ V_h^{n+1} \end{bmatrix} = \begin{bmatrix} e^{ikB_{h,0}} & 0 \\ 0 & e^{-ikB_{h,0}} \end{bmatrix} \begin{bmatrix} U_h^n \\ V_h^n \end{bmatrix}
$$
$$
+ k \begin{bmatrix} V_h^n \\ \gamma_1(kB_{h,0})[f(t_{n+1}, e^{ikB_{h,0}}U_h^n + kV_h^n) - A_hQ_hg(t_{n+1})] \end{bmatrix}, \quad (15)
$$

where $V_h^0$ is related to the starting values $U_h^0, U_h^1$ through

$$
V_h^0 = \frac{1}{k}[U_h^1 - e^{ikB_{h,0}}U_h^0]. \quad (16)
$$

Then, considering $V_h(t_n) = [U_h(t_{n+1}) - e^{ikB_{h,0}}U_h(t_n)]/k$, the first component of the local error corresponding to the one-step method vanishes and the second component, which we denote as $\bar{\rho}_{h,n}$, satisfies

$$
\begin{aligned}
\bar{\rho}_{h,n} &= V_h(t_{n+1}) - e^{-ikB_{h,0}}V_h(t_n) \\
&\quad - k\gamma_1(kB_{h,0})[P_h f(t_{n+1}, U_h(t_{n+1})) - A_hQ_hg(t_{n+1})] \\
&= \frac{1}{k}\Big[ U_h(t_{n+2}) - e^{ikB_{h,0}}U_h(t_{n+1}) - e^{-ikB_{h,0}}[U_h(t_{n+1}) - e^{ikB_{h,0}}U_h(t_n)] \\
&\quad - k^2\gamma_1(kB_{h,0})[P_h f(t_{n+1}, U_h(t_{n+1})) - A_hQ_hg(t_{n+1})]\Big] \\
&= \frac{1}{k}\rho_{h,n+1}.
\end{aligned}
$$

Then, considering the global error for the one-step method

$$
E_{h,n} = \begin{bmatrix} U_h(t_n) - U_h^n \\ V_h(t_n) - V_h^n \end{bmatrix}, \quad (17)
$$

it happens that

$$
E_{h,n+1} = R(kB_{h,0})E_{h,n} + k\phi(V_h^n, V_h(t_n), U_h^n, U_h(t_n)) + \begin{bmatrix} 0 \\ \frac{1}{k}\rho_{h,n+1} \end{bmatrix},
$$

where $R(kB_{h,0})$ is the matrix in (15) and $\phi(V_h^n, V_h(t_n), U_h^n, U_h(t_n))$ comes from making the difference of the second term in (15) with the same one which corresponds to substituting the numerical solution by the exact solution, and which satisfies

$$
\|\phi(V_h^n, V_h(t_n), U_h^n, U_h(t_n))\|_h \leq C\|E_{h,n}\|_h. \quad (18)
$$

From here, in a recursive way,

$$
E_{h,n} = R^n(kB_{h,0})E_{h,0} + \sum_{l=1}^{n} R^{n-l}(kB_{h,0}) \begin{bmatrix} 0 \\ \frac{\rho_{h,l}}{k} \end{bmatrix}
$$
$$
+ k\sum_{l=0}^{n-1} R^{n-l-1}(kB_{h,0})\phi(V_{h,l}, V_h(t_l), U_h^l, U_h(t_l)), \quad (19)
$$

9

and the classical argument does not lead to convergence although $\|R^{n-l}(kB_{h,0})\|$ is bounded because $\frac{\rho_{h,l}}{k}$ is just $O(k)$.

In order to get convergence, a summation-by-parts argument must be applied, and for that the assumptions of the following theorem must be made.

**Theorem 3.** *Under hypotheses (A1)-(A4) and (H1)-(H5), whenever the starting values are exact except for $O(k^2 + k\varepsilon_h)$, $g \in C^3([0,T], Y)$, $f \in C^2([0,T] \times X, X)$ and $\ddot{u} \in C([0,T], Z)$, if for some constant $C$, $h$ and $k$ are such that one of the following bounds are satisfied*

$$\|\sum_{r=1}^{n-1} e^{-ikrB_{h,0}} B_{h,0}\|_h \le Cn, \qquad \|\sum_{r=1}^{n-1} e^{-ikrB_{h,0}}\|_h \le C, \quad nk \le T, \qquad (20)$$

*the global error $e_{h,n} = P_h u(t_n) - U_h^n$ is bounded by $C'(k + \varepsilon_h)$ for some constant $C'$ which depends on $C$ but not on the particular values of $h$ and $k$.*

*Proof.* First of all, let us notice that

$$P_h u(t_n) - U_h^n = (P_h u(t_n) - U_h(t_n)) + (U_h(t_n) - U_h^n),$$

where the first parenthesis is $O(\varepsilon_h)$ according to (14) and the second is going to be analysed now taking into account that it is the first component of $E_{h,n}$. It suffices to notice that

$$\sum_{l=1}^{n} R^{n-l}(kB_{h,0}) \begin{bmatrix} 0 \\ \frac{\rho_{h,l}}{k} \end{bmatrix} = \left( \sum_{r=1}^{n-1} R^r(kB_{h,0}) \right) \begin{bmatrix} 0 \\ \frac{\rho_{h,1}}{k} \end{bmatrix}$$
$$+ \sum_{j=2}^{n-1} \left( \sum_{r=1}^{j-1} R^r(kB_{h,0}) \right) \begin{bmatrix} 0 \\ \frac{\rho_{h,n-j+1}}{k} - \frac{\rho_{h,n-j}}{k} \end{bmatrix} + \begin{bmatrix} 0 \\ \frac{\rho_{h,n}}{k} \end{bmatrix}.$$

Then, because of the regularity of $g$, $f$ and $u$, $\frac{\rho_{h,n-j+1}}{k} - \frac{\rho_{h,n-j}}{k}$ is $O(k^2)$ following the proof of Theorem 2. Due to this, if the second bound in (20) holds, the whole term above is bounded by $kC''$ for some constant $C''$ and the result follows by applying a discrete Gronwall lemma to the corresponding bound of (19).

On the other hand, if it is the first bound in (20) which holds, because of Theorem 2, we would have $B_{h,0}^{-1}\frac{\rho_{h,1}}{k} = O(k^2)$ and $B_{h,0}^{-1}[\frac{\rho_{h,n-j+1}}{k} - \frac{\rho_{h,n-j}}{k}] = O(k^3)$ and therefore the whole term above would again be bounded by $kC''$ for some constant $C''$ and the result would follow in the same way. $\square$

**Remark 4.** *Notice that, when $A_{h,0}$ is symmetric and the norm is the Euclidean one, the norms in (20) correspond to the spectral radius. Then, as*

$$\sum_{r=1}^{n-1} e^{-ikr\lambda} = \frac{e^{-ik\lambda} - e^{-ikn\lambda}}{1 - e^{-ik\lambda}},$$

*when $k\lambda$ is small, this behaves as $n - 1$ and, when $k\lambda$ is big enough this is bounded. Therefore, none of the assumptions (20) are satisfied uniformly in $h$ and $k$ when they both diminish but the first one is satisfied for every small enough $k$ when $h$ is moderate and the second one for every $h$ when $k$ is not too small. This explains that, in the numerical experiments of Section 6, the error behaves as the theorem states.*

# 4 Suggested approach: Discretizing firstly in time and then in space

When trying to discretize firstly (1) in time, $B^2$ in (3) must be substituted by the differential operator $A$. Therefore, when considering something similar to (4), a proposal must be given for both $\cos(kB)$ and $B^{-1}\sin((k-s)B)$ when applied to a function in $H$. More precisely, at it seems natural from the definition of the cosine function, we will propose as the corresponding $\cos(kB)u_n$ to the solution at $s = k$ of

$$
\begin{aligned}
\ddot{v}_n(s) &= -Av_n(s), \\
v_n(0) &= u_n, \\
\dot{v}_n(0) &= 0, \\
\partial v_n(s) &= \partial \hat{v}_n(s), \qquad \hat{v}_n(s) = u(t_n) - \frac{s^2}{2}Au(t_n).
\end{aligned}
\tag{21}
$$

Notice that, for the boundary, we have considered a Taylor expansion of third order for the solution of the initial value problem which corresponds to the first three lines of (21) when substituting $u_n$ by the solution $u(t_n)$ of (1). In such a way, we are near the solution we want to approximate and, moreover, that boundary can be calculated in terms of the data of the problem, either directly when the boundary conditions are Dirichlet or through the numerical approximation itself when the boundary conditions are Robin or Neumann [5].

When discretizing (21) in space, the following system arises:

$$
\begin{aligned}
\ddot{V}_{n,h}(s) &= -A_{h,0}V_{n,h}(s) - A_hQ_h\partial[u(t_n) - \frac{s^2}{2}Au(t_n)], \\
V_{n,h}(0) &= U_h^n, \\
\dot{V}_{n,h}(0) &= 0,
\end{aligned}
\tag{22}
$$

and, applying (4), the solution of this is

$$
V_{n,h}(k) = \cos(kB_{h,0})U_h^n - \int_0^k B_{h,0}^{-1}\sin((k-s)B_{h,0})A_hQ_h\partial[u(t_n) - \frac{s^2}{2}Au(t_n)]ds.
$$

Then, arguing as when deducing Gautschi method from (5) to (6),

$$
\int_0^k B_{h,0}^{-1}\sin((k-s)B_{h,0})A_hQ_h\partial u(t_n)ds = \frac{k^2}{2}\gamma_1(kB_{h,0})A_hQ_h\partial u(t_n).
$$

Besides, as integrating by parts,

$$
\int_0^k \lambda^{-1}\sin((k-s)\lambda)s^2ds = \frac{1}{\lambda^4}[(k\lambda)^2 - 2(1 - \cos(k\lambda))],
$$

it finally turns out that

$$
V_{n,h}(k) = \cos(kB_{h,0})U_h^n - \frac{k^2}{2}\gamma_1(kB_{h,0})A_hQ_h\partial u(t_n) + \frac{k^4}{2}\gamma_2(kB_{h,0})A_hQ_h\partial Au(t_n), \tag{23}
$$

where

$$\gamma_2(\epsilon) = \frac{1}{\epsilon^4}[\epsilon^2 - 2(1 - \cos(\epsilon))].\tag{24}$$

On the other hand, we propose as the corresponding $B^{-1}\sin((k-s)B)f(t_n, u_n)$ the solution at $\tau = k - s$ of

$$\begin{aligned}
\ddot{w}_n(\tau) &= -Aw_n(\tau),\\
w_n(0) &= 0,\\
\dot{w}_n(0) &= f(t_n, u_n),\\
\partial w_n(\tau) &= \hat{w}_n(\tau), \qquad \hat{w}_n(\tau) = \tau f(t_n, u(t_n)).
\end{aligned}\tag{25}$$

Notice that now, for the boundary, we have considered a Taylor expansion of second order for the solution of the initial value problem which corresponds to the first three lines of (25) when substituting $u_n$ by the solution $u(t_n)$ of (1). When discretizing in space, the following system arises:

$$\begin{aligned}
\ddot{W}_{n,h}(\tau) &= -A_{h,0}W_{n,h}(\tau) - \tau A_h Q_h \partial f(t_n, u(t_n)),\\
W_{n,h}(0) &= 0,\\
\dot{W}_{n,h}(0) &= f(t_n, U_h^n),
\end{aligned}\tag{26}$$

and again applying (4),

$$W_{h,n}(\tau) = \tau\mathrm{sinc}(\tau B_{h,0})f(t_n, U_h^n) - \int_0^\tau B_{h,0}^{-1}\sin((\tau - \sigma)B_{h,0})\sigma A_h Q_h \partial f(t_n, u(t_n))d\sigma.$$

Now, using that

$$\int_0^\tau \lambda^{-1}\sin((\tau - \sigma)\lambda)\sigma d\sigma = \frac{\tau}{\lambda^2}(1 - \mathrm{sinc}(\tau\lambda)),$$

$$W_{h,n}(\tau) = \tau\mathrm{sinc}(\tau B_{h,0})f(t_n, U_h^n) - \tau B_{h,0}^{-2}[I - \mathrm{sinc}(\tau B_{h,0})]A_h Q_h \partial f(t_n, u(t_n)).\tag{27}$$

Then, in the corresponding expression of (5) where the bracket is substituted by $2f(t_n, U_h^n)$, instead of considering $2\int_0^k B^{-1}\sin((k-s)B)f(t_n, u_n)ds$, we propose

$$\begin{aligned}
2\int_0^k &W_{h,n}(k-s)ds\\
&= 2\int_0^k \Big[B_{h,0}^{-1}\sin((k-s)B_{h,0})f(t_n, U_h^n) - (k-s)B_{h,0}^{-2}A_h Q_h \partial f(t_n, u(t_n))\\
&\qquad\qquad + B_{h,0}^{-3}\sin((k-s)B_{h,0})A_h Q_h \partial f(t_n, u(t_n))\Big]ds\\
&= 2\Big[B_{h,0}^{-2}(I - \cos(kB_{h,0}))f(t_n, U_h^n) - \frac{k^2}{2}B_{h,0}^{-2}A_h Q_h \partial f(t_n, u(t_n))\\
&\qquad\qquad + B_{h,0}^{-4}(I - \cos(kB_{h,0}))A_h Q_h \partial f(t_n, u(t_n))\Big]\\
&= k^2\gamma_1(kB_{h,0})f(t_n, U_h^n) - k^4\gamma_2(kB_{h,0})A_h Q_h \partial f(t_n, u(t_n)),
\end{aligned}$$

where $\gamma_2(\epsilon)$ is that in (24).

Gathering this with (23) and using (1), our suggestion for Gautschi method after full discretization is

$$U_h^{n+1} - 2\cos(kB_{h,0})U_h^n + U_h^{n-1}$$
$$= k^2\gamma_1(kB_{h,0})[f(t_n, U_h^n) - A_hQ_hg(t_n)] - k^4\gamma_2(kB_{h,0})A_hQ_h\ddot{g}(t_n), \qquad (28)$$

which just differs in the last term with the classical approach in (11).

## 4.1 Semidiscretization local truncation error

In this section we will analyse how the local truncation error behaves with our approach just after time discretization. For that we notice that our method in fact reads

$$u_{n+1} = 2v_n(k) - u_{n-1} + 2\int_0^k w_n(k-s)ds,$$

where $v_n$ and $w_n$ are those defined in (21) and (25) respectively. Then, the local truncation error is defined as

$$\rho_n = u(t_{n+1}) - 2\bar{v}_n(k) + u(t_{n-1}) - 2\int_0^k \bar{w}_n(k-s)ds, \qquad (29)$$

where $\bar{v}_n$ (resp. $\bar{w}_n$) corresponds to (21) (resp. (25)) when $u_n$ is replaced by $u(t_n)$. We do have the following result:

**Theorem 5.** *Under hypotheses (A1)-(A4), whenever $u \in C^4([0,T],X)$, $u \in C^2([0,T],D(A^2))$ and $f(\cdot, u(\cdot)) \in C([0,T],D(A))$, $\rho_n = O(k^4)$.*

*Proof.* The key of the proof is that $z_n(s) = \bar{v}_n(s) - \hat{v}_n(s)$ is solution of the problem

$$\ddot{z}_n(s) = -Az_n(s) + \frac{s^2}{2}A^2u(t_n),$$
$$z_n(0) = 0,$$
$$\dot{z}_n(0) = 0,$$
$$\partial z_n(s) = 0,$$

and $\nu_n(\tau) = \bar{w}_n(\tau) - \hat{w}_n(\tau)$ is solution of

$$\ddot{\nu}_n(\tau) = -A\nu_n(\tau) - \tau Af(t_n, u(t_n)),$$
$$\nu_n(0) = 0,$$
$$\dot{\nu}_n(0) = 0,$$
$$\partial\nu_n(\tau) = 0.$$

Because of this,

$$z_n(s) = \int_0^s B_0^{-1}\sin((s-\sigma)B_0)\frac{\sigma^2}{2}A^2u(t_n)d\sigma = \int_0^s (s-\sigma)\text{sinc}((s-\sigma)B_0)\frac{\sigma^2}{2}A^2u(t_n)d\sigma = O(s^4),$$

and

$$\begin{aligned}
\nu_n(\tau) &= -\int_0^\tau B_0^{-1}\sin((\tau-\sigma)B_0)\sigma Af(t_n, u(t_n))d\sigma \\
&= -\int_0^\tau (\tau-\sigma)\text{sinc}((\tau-\sigma)B_0)\sigma Af(t_n, u(t_n))d\sigma = O(\tau^3).
\end{aligned}$$

Therefore,

$$\begin{aligned}
\rho_n &= u(t_{n+1}) - 2\hat{v}_n(k) + u(t_{n-1}) - 2\int_0^k \hat{w}_n(k-s)ds + O(k^4) \\
&= u(t_n) + k\dot{u}(t_n) + \frac{k^2}{2}[-Au(t_n) + f(t_n, u(t_n))] + \frac{k^3}{6}\dddot{u}(t_n) - 2[u(t_n) - \frac{k^2}{2}Au(t_n)] \\
&\quad + u(t_n) - k\dot{u}(t_n) + \frac{k^2}{2}[-Au(t_n) + f(t_n, u(t_n))] - \frac{k^3}{6}\dddot{u}(t_n) - k^2 f(t_n, u(t_n)) + O(k^4) \\
&= O(k^4).
\end{aligned}$$

$\square$

## 4.2 Full discretization local truncation error

In this subsection we will study the local truncation error after discretizing in time and then in space. For that, we consider $\bar{V}_{h,n}$ (resp. $\bar{W}_{h,n}$) as the solution of (23) (resp. (27)) when $U_h^n$ is substituted by $P_h u(t_n)$. Then, because of the way the method is defined, we define that local truncation error after full discretization as

$$\rho_{n,h} = P_h u(t_{n+1}) - 2\bar{V}_{h,n}(k) + P_h u(t_{n-1}) - 2\int_0^k \bar{W}_{h,n}(k-s)ds. \tag{30}$$

We thus have the following result:

**Theorem 6.** *Apart from the hypotheses of Theorem 5 and (H1)-(H5), let us also assume that $u, Au \in C([0,T], Z)$ and that $f(\cdot, u(\cdot)) \in C([0,T], Z)$. Then, $\rho_{n,h} = O(k^2\varepsilon_h + k^4)$.*

*Proof.* Let us consider $\bar{V}_{h,n}(s)$ (resp. $\bar{W}_{n,h}(\tau)$) the solutions of (23) (resp. (27)) where $U_h^n$ is replaced by $P_h u(t_n)$. Then,

$$\begin{aligned}
\ddot{\bar{V}}_{h,n}(s) - P_h\ddot{\hat{v}}(s) &= -A_{h,0}\bar{V}_{h,n}(s) - A_h Q_h \partial\hat{v}_n(s) + P_h Au(t_n) \\
&= -A_{h,0}(\bar{V}_{h,n}(s) - P_h\hat{v}_n(s)) + A_{h,0}(R_h - P_h)\hat{v}_n(s) + P_h(Au(t_n) - A\hat{v}_n(s)) \\
&= -A_{h,0}(\bar{V}_{h,n}(s) - P_h\hat{v}_n(s)) + A_{h,0}(R_h - P_h)\hat{v}_n(s) + \frac{s^2}{2}P_h A^2 u(t_n),
\end{aligned}$$

$$\bar{V}_{h,n}(0) - P_h\hat{v}_n(0) = 0,$$

$$\dot{\bar{V}}_{h,n}(0) - P_h\dot{\hat{v}}_n(0) = 0,$$

and

$$\dddot{\bar{W}}_{h,n}(\tau) - P_h\dddot{\hat{w}}(\tau) = -A_{h,0}\dot{\bar{W}}_{h,n}(\tau) - A_h Q_h \partial\hat{w}_n(\tau)$$
$$= -A_{h,0}(\dot{\bar{W}}_{h,n}(\tau) - P_h\dot{\hat{w}}_n(\tau)) - A_{h,0}(P_h - R_h)\hat{w}_n(\tau) - P_h A\hat{w}_n(\tau),$$
$$\bar{W}_{h,n}(0) - P_h\hat{w}_n(0) = 0,$$
$$\dot{\bar{W}}_{h,n}(0) - P_h\dot{\hat{w}}_n(0) = 0.$$

From this,

$$\bar{V}_{h,n}(k) - P_h\hat{v}_n(k) = \int_0^k B_{h,0}^{-1}\sin((k-s)B_{h,0})[A_{h,0}(R_h - P_h)\hat{v}_n(s) + \frac{s^2}{2}P_h A^2 u(t_n)]ds$$

$$= \int_0^k (k-s)\mathrm{sinc}((k-s)B_{h,0})[O(\varepsilon_h) + O(s^2)] = O(k^2\varepsilon_h + k^4),$$

and

$$\bar{W}_{h,n}(\tau) - P_h\hat{w}_n(\tau) = \int_0^\tau B_{h,0}^{-1}\sin((\tau-\sigma)B_{h,0})[A_{h,0}(R_h - P_h)\hat{w}_n(\sigma) - \sigma P_h A f(t_n, u(t_n))]d\sigma$$

$$= \int_0^\tau (\tau-\sigma)\mathrm{sinc}((\tau-\sigma)B_{h,0})[O(\varepsilon_h) + O(\sigma)] = O(\tau^2\varepsilon_h + \tau^3).$$

Here, we have used that $\hat{v}_n$ and $\hat{w}_n$ belong to $C([0,T], Z)$ because of the hypotheses of the theorem.

Because of the above, $\rho_{n,h}$ in (30) can be written as

$$\rho_{n,h} = P_h\rho_n + O(k^2\varepsilon_h + k^4) = O(k^2\varepsilon_h + k^4),$$

where the definition of $\rho_n$ (29) and Theorem 5 has been used. □

## 4.3 Full discretization global error

In a similar way as in Section 3, the suggested implementation of Gautschi method can be written as the one-step method

$$\begin{bmatrix} U_h^{n+1} \\ V_h^{n+1} \end{bmatrix} = \begin{bmatrix} e^{ikB_{h,0}} & 0 \\ 0 & e^{-ikB_{h,0}} \end{bmatrix}\begin{bmatrix} U_h^n \\ V_h^n \end{bmatrix}$$
$$+k\begin{bmatrix} V_h^n \\ \gamma_1(kB_{h,0})[P_h f(t_{n+1}, e^{ikB_{h,0}}U_h^n + kV_h^n) - A_h Q_h g(t_{n+1})] - k^3\gamma_2(kB_{h,0})A_h Q_h\ddot{g}(t_{n+1}) \end{bmatrix},$$

where $V_h^0$ is related to the starting values $U_h^0, U_h^1$ through (16). Then, considering

$$V_h(t_{n+1}) = \frac{1}{k}[P_h u(t_{n+1}) - e^{ikB_{h,0}}P_h u(t_n)],$$

and defining

$$E_{n,h} = \begin{bmatrix} P_h u(t_n) - U_h^n \\ V_h(t_n) - V_h^n \end{bmatrix},$$

the same classical argument of convergence which is mentioned just after formula (19) but considering that now $\rho_{l,h}/k$ is $O(k^3 + k\varepsilon_h)$ leads to the following result, which proves that order reduction is avoided:

**Theorem 7.** *Under the hypotheses of Theorems 5 and 6, if the starting values are such that they satisfy $U_h^j - P_h u(t_j) = O(k^3 + k\varepsilon_h)$ $(j = 0,1)$, it happens that*

$$P_h u(t_n) - U_h^n = O(k^2 + \varepsilon_h).$$

# 5 Derivative approximation

Although Gautschi initial method just proposes an approximation for the solution $u(t)$ of (3), it is usual to consider also an approximation for the derivative which also based on a similar simple interpolation after applying the variation-of-constants formula (see [15, 17]). In such a way, from

$$\dot{U}(t_n + k) = -B\sin(kB)U(t_n) + \cos(kB)\dot{U}(t_n) + \int_0^k \cos((k-s)B)F(t_n + s, U(t_n + s)),$$

which can be obtained by differentiating (4) to respecto to $k$, it can be deduced that

$$\dot{U}(t_n+k)-\dot{U}(t_n-k) = -2B\sin(kB)U(t_n)+\int_0^k \cos((k-s)B)[F(t_n+s, U(t_n+s))+F(t_n-s, U(t_n-s))]ds,$$

and then the numerical approximation is given by

$$\dot{U}_{n+1} - \dot{U}_{n-1} = -2B\sin(kB)U(t_n) + 2k\mathrm{sinc}(kB)F_n. \tag{31}$$

In this section we will analyse the error in the derivative using also the classical approach and will suggest an approach for which order reduction is avoided in such a way that a bound for the error in the derivative is obtained as that shown in [15] with periodic boundary conditions. (More precisely, considering $s = 0$ and $\alpha = 1$ in Theorem 2.1 of that reference.)

## 5.1 Classical approach

Approximating the derivative of (10) using $U_h^n$ in (11), this additional difference system turns up:

$$\dot{U}_h^{n+1} - \dot{U}_h^{n-1} = -2B_{h,0}\sin(kB_{h,0})U_h^n + 2k\mathrm{sinc}(kB_{h,0})[-A_hQ_hg(t_n) + f(t_n, U_h^n)]. \tag{32}$$

Then, we have the following result for the local truncation error corresponding to this system:

**Theorem 8.** *Under the same hypotheses of Theorem 2,*

$$B_{h,0}^{-1}\dot{\rho}_{h,n} = O(k^2),$$

*where the constant in Landau notation is independent of $k$ and $h$.*

*Proof.* It suffices to notice that, similarly to the proof of Theorem 2,

$$
\begin{aligned}
B_{h,0}^{-1}\dot{\rho}_{h,n} &= \int_0^k B_{h,0}^{-1}\cos((k-s)B_{h,0})\Big[-A_hQ_h[g(t_n+s)+g(t_n-s)-2g(t_n)] \\
&\qquad\qquad +f(t_n+s,U_h(t_n+s))+f(t_n-s,U_h(t_n-s))-2f(t_n,U_h(t_n)))\Big]ds \\
&= -\int_0^k B_{h,0}\cos((k-s)B_{h,0})A_{h,0}^{-1}A_hQ_h[g(t_n+s)+g(t_n-s)-2g(t_n)]ds+O(k^3) \\
&= -\int_0^k B_{h,0}\cos((k-s)B_{h,0})\Big[\int_0^s A_{h,0}^{-1}A_hQ_h[\ddot{g}(t_n+\sigma)+\ddot{g}(t_n-\sigma)](s-\sigma)d\sigma\Big]ds+O(k^3) \\
&= -\int_0^k \sin((k-s)B_{h,0})\Big[\int_0^s A_{h,0}^{-1}A_hQ_h[\ddot{g}(t_n+\sigma)+\ddot{g}(t_n-\sigma)]d\sigma\Big]ds+O(k^3)=O(k^2),
\end{aligned}
$$

where (H3) and (H4) have been used. $\qquad\square$

We notice that $\dot{\rho}_{h,n}$ can be seen to behave as $O(k)$, but in such a case the Landau constant grows with $h$ and, therefore, we believe that bound is not meaningful.

As for the global error, as $U_h^n$ is present in (32), the analysis must be made considering also the approximation for the solution. However, as $\dot{U}_h^{n+1}$ does not turn up in (32), we can do it advancing two stepsizes at a time and therefore using two successive applications of (15). More precisely,

$$
\begin{aligned}
\begin{bmatrix} U_h^{n+2} \\ V_h^{n+2} \\ B_{h,0}^{-1}\dot{U}_h^{n+2} \end{bmatrix} &= \begin{bmatrix} e^{2ikB_{h,0}} & 0 & 0 \\ 0 & e^{-2ikB_{h,0}} & 0 \\ -2\sin(kB_{h,0})e^{ikB_{h,0}} & 0 & I \end{bmatrix}\begin{bmatrix} U_h^n \\ V_h^n \\ B_{h,0}^{-1}\dot{U}_h^n \end{bmatrix} \\
&\quad +k\begin{bmatrix} e^{ikA_{h,0}}V_h^n+e^{-ikA_{h,0}}V_h^n+k\cdot* \\ e^{-ikB_{h,0}}\cdot* \\ +\gamma_1(kB_{h,0})\Big[f(t_{n+2},e^{ikB_{h,0}}(e^{ikB_{h,0}}U_h^n+kV_h^n)+k[e^{-ikB_{h,0}}V_h^n+k\cdot*-A_hQ_hg(t_{n+2})]\Big] \\ -2\sin(kB_{h,0})V_h^n+2\mathrm{sinc}(kB_{h,0})[f(t_{n+1},e^{ikB_{h,0}}U_h^n+kV_h^n)-A_hQ_hg(t_{n+1})] \end{bmatrix},
\end{aligned}
$$
(33)

where $*$ stands for $\gamma_1(kB_{h,0})[f(t_{n+1},e^{ikB_{h,0}}U_h^n+kV_h^n)-A_hQ_hg(t_{n+1})]$. Then, the local truncation error corresponding to this difference system is

$$
\bar{\bar{\rho}}_{h,n}=\begin{bmatrix} \rho_{h,n+1} \\ O(\tfrac{1}{k}\rho_{h,n+1}+\tfrac{1}{k}\rho_{h,n+2}) \\ B_{h,0}^{-1}\dot{\rho}_{h,n} \end{bmatrix}=\begin{bmatrix} O(k^2) \\ O(k) \\ O(k^2) \end{bmatrix},
$$

and it happens that $\bar{\bar{E}}_{h,n}=[U_h(t_n)-U_h^n,V_h(t_n)-V_h^n,B_{h,0}^{-1}[\dot{U}_h(t_n)-\dot{U}_h^n]]^T$ satisfies

$$
\bar{\bar{E}}_{h,n+2}=\bar{\bar{R}}(kB_{h,0})\bar{\bar{E}}_{h,n}+k\bar{\bar{\phi}}(V_h^n,V_h(t_n),U_h^n,U_h(t_n))+\bar{\bar{\rho}}_{h,n},
$$

where $\bar{\bar{R}}(kB_{h,0})$ is the matrix in (33) and $\bar{\bar{\phi}}$ satisfies the same as $\phi$ in (18) with $E_{h,n}$

the first two sets of components of $\bar{\bar{E}}_{h,n}$, as described in (17). Therefore, for $j = 0, 1$,

$$\bar{\bar{E}}_{h,2n+j} = \bar{\bar{R}}^n(kB_{h,0})\bar{\bar{E}}_{h,j} + \sum_{l=1}^{n} \bar{\bar{R}}^{n-l}(kB_{h,0})\bar{\bar{p}}_{h,2l+j-2}$$

$$+k\sum_{l=0}^{n-1} \bar{\bar{R}}^{n-l-1}(kB_{h,0})\bar{\bar{\phi}}(V_{h,2l+j}, V_h(t_{2l+j}), U_h^{2l+j}, U_h(t_{2l+j})). \tag{34}$$

Now, notice that

$$\bar{\bar{R}}^n(kB_{h,0}) = \begin{bmatrix} e^{2inkB_{h,0}} & 0 & 0 \\ 0 & e^{-2inkB_{h,0}} & 0 \\ -2\sin(kB_{h,0})e^{ikB_{h,0}}(I - e^{ikB_{h,0}})^{-1}[I - e^{inkB_{h,0}}] & 0 & I \end{bmatrix}. \tag{35}$$

Because of this, under the hypotheses of Theorem 3 and assuming also that

$$B_{h,0}^{-1}[\dot{U}_h(t_j) - \dot{U}_h^j] = O(k + \varepsilon_h), \tag{36}$$

using (H4) again, the last set of components of (34) state that , for $j = 0, 1$,

$$B_{h,0}^{-1}[\dot{U}_h(t_{2n+j}) - \dot{U}_h^{2n+j}] = O(k + \varepsilon_h), \quad n > 0.$$

Then, the following result follows.

**Theorem 9.** *Under the hypotheses of Theorem 3 and assuming also (36),*

$$B_{h,0}^{-1}[P_h\dot{u}(t_n) - \dot{U}_h^n] = O(k + \varepsilon_h).$$

*Proof.* The result comes from the argument above just considering also that

$$B_{h,0}^{-1}[P_h\dot{u}(t_n) - \dot{U}_h^n] = B_{h,0}^{-1}[P_h u(t_n) - \dot{U}_h(t_n)] + B_{h,0}^{-1}[\dot{U}_h(t_n) - \dot{U}_h^n],$$

and that differentiating (13) with respect to time,

$$P_h\dot{u}(t_n) - \dot{U}_h^n = \int_0^t \cos((t-s)B_{h,0})[f(s, P_h u(s)) - f(s, U_h(s)) + A_{h,0}[P_h u(s) - R_h u(s)]]ds.$$

Using now (14) and (H5),

$$\|P_h\dot{u}(t) - \dot{U}_h(t)\|_h \leq \max(T^2 e^{LT}, T)\varepsilon_h \max_{t \in [0,T]} \|u(t)\|_Z, \quad 0 \leq t \leq T,$$

and, as $B_{h,0}^{-1}$ is bounded, the result follows. $\qquad \square$

## 5.2 Suggested approach

Discretizing firstly in time simulating (31) in some way, we obtain

$$\dot{u}_{n+1} - \dot{u}_{n-1} = -2\eta_n(k) + 2\sigma_n(k),$$

where $\eta_n$ and $\sigma_n$ satisfy

$$
\begin{aligned}
\ddot{\eta}_n(s) &= -A\eta_n(s), & \ddot{\sigma}_n(s) &= -A\eta_n(s), \\
\eta_n(0) &= 0, & \sigma_n(0) &= 0, \\
\dot{\eta}_n(0) &= Au_n, & \dot{\sigma}_n(0) &= f(t_n, u_n), \\
\partial\eta_n(s) &= \partial[sAu(t_n)], & \partial\sigma_n(s) &= \partial[sf(t_n, u(t_n))].
\end{aligned}
\tag{37}
$$

Then, after space discretization of (37), the following systems turns up

$$
\begin{aligned}
\ddot{\eta}_{n,h}(s) &= -A_{h,0}\eta_{n,h}(s) - A_hQ_h\partial[sAu(t_n)], & \ddot{\sigma}_{n,h}(s) &= -A_{h,0}\sigma_{n,h}(s) - A_hQ_h\partial[sf(t_n, u(t_n))], \\
\eta_{n,h}(0) &= 0, & \sigma_{n,h}(0) &= 0, \\
\dot{\eta}_{n,h}(0) &= A_{h,0}U_h^n + A_hQ_h\partial u(t_n), & \dot{\sigma}_{n,h}(0) &= f(t_n, U_h^n),
\end{aligned}
\tag{38}
$$

whose solutions are, using (4) and arguing as for (27),

$$\eta_{n,h}(k) = k\operatorname{sinc}(kB_{h,0})[A_{h,0}U_h^n + A_hQ_h\partial u(t_n)] - kB_{h,0}^{-2}[I - \operatorname{sinc}(kB_{h,0})]A_hQ_h\partial Au(t_n).$$

$$\sigma_{n,h}(k) = k\operatorname{sinc}(kB_{h,0})f(t_n, U_h^n) - kB_{h,0}^{-2}[I - \operatorname{sinc}(kB_{h,0})]A_hQ_h\partial f(t_n, u(t_n)),$$

and so the approximation for the derivative after full discretization is given by

$$
\begin{aligned}
\dot{U}_h^{n+1} - \dot{U}_h^{n-1} &= 2k\operatorname{sinc}(kB_{h,0})[-A_{h,0}U_h^n - A_hQ_hg(t_n) + f(t_n, U_h^n)] \\
&\quad -2kB_{h,0}^{-2}[I - \operatorname{sinc}(kB_{h,0})]A_hQ_h\ddot{g}(t_n),
\end{aligned}
\tag{39}
$$

where (1) has also been used. Notice that the last line of this formula is what is added with respect to the classical approach in (32).

As for the analysis, let us begin with the semidiscretization local truncation error

$$\dot{\rho}_n = \dot{u}(t_{n+1}) - \dot{u}(t_{n-1}) + 2\bar{\eta}_n(k) - 2\bar{\sigma}_n(k), \tag{40}$$

where $\bar{\eta}_n(s)$ and $\bar{\sigma}_n(s)$ are like those in (37) but starting from $Au(t_n)$ and $f(t_n, u(t_n))$ respectively at the derivative. Then, we have the following result.

**Theorem 10.** *Under the same hypotheses of Theorem 5, $\dot{\rho}_n = O(k^3)$.*

*Proof.* Following the same arguments of the proof of Theorem 5, it can be deduced that

$$
\begin{aligned}
\bar{\eta}_n(s) - sAu(t_n) &= -\int_0^s B_0^{-1}\sin((s-\sigma)B_0)\sigma A^2 u(t_n)d\sigma \\
&= -\int_0^s (s-\sigma)\operatorname{sinc}((s-\sigma)B_0)\sigma A^2 u(t_n)d\sigma = O(s^3),
\end{aligned}
\tag{41}
$$

and the same happens with $\bar{\sigma}_n(s) - sf(t_n, u(t_n))$. Therefore,

$$
\begin{aligned}
\dot{\rho}_n &= \dot{u}(t_{n+1}) - \dot{u}(t_{n-1}) + 2kAu(t_n) - 2kf(t_n, u(t_n)) + O(k^3) \\
&= 2k[\ddot{u}(t_n) + Au(t_n) - f(t_n, u(t_n))] + O(k^3),
\end{aligned}
\tag{42}
$$

where the bracket vanishes because of (1). $\qquad\square$

As for the full discretization local truncation error, we do have the following result.

**Theorem 11.** *Under the same hypotheses of Theorem 6, $\dot{\rho}_{n,h} = O(k^2\varepsilon_h + k^3)$.*

*Proof.* We firstly notice that

$$\dot{\rho}_{n,h} = P_h\dot{u}(t_{n+1}) - P_h\dot{u}(t_{n1}) + 2\bar{\eta}_{n,h}(k) - 2\bar{\sigma}_{n,h}(k), \qquad (43)$$

where $\bar{\eta}_{n,h}$ and $\bar{\sigma}_{n,h}$ are defined as in (38) but changing $U_h^n$ by $P_h u(t_n)$. Then, with a similar argument as that for $\bar{W}_{h,n}(\tau)$ in the proof of Theorem 6,

$$\bar{\eta}_{n,h}(k) - kP_h Au(t_n) = O(k^2\varepsilon_h + k^3),$$
$$\bar{\sigma}_{n,h}(k) - kP_h f(t_n, u(t_n)) = O(k^2\varepsilon_h + k^3).$$

From this, using the definition of $\dot{\rho}_{n,h}$ (43), $\dot{\rho}_n$ (40) and (41),

$$\dot{\rho}_{n,h} = P_h\dot{\rho}_n + O(k^2\varepsilon_h + k^3) = O(k^2\varepsilon_h + k^3),$$

where the last equality comes from Theorem 10. □

Compare this with the result in Theorem **??** where, not only less order was obtained in the timestepsize, but also it was necessary to multiply by $B_{h,0}^{-1}$ in order to get a bound which did not grow with $h$. Nevertheless, for the global error now, we will also need to consider the difference $B_{h,0}^{-1}[P_h\dot{u}(t_n) - \bar{U}_h^n]$ in order to get a bound which does not grow with $h$. This is due to the fact that $P_h\dot{u}(t_n) - \dot{U}_h^n$ would depend on $B_{h,0}\rho_{n,h}$ where $\rho_{n,h}$ is the local truncation error in the solution, and it happens that $B_{h,0}\rho_{n,h}$ can be seen to behave as $O(k^4)$ when the error in space is negligible, but with a Landau constant which grows with $h$. Therefore, in order to get a bound for the derivative for which that does not happen, we must consider an argument similar to that previous to Theorem 9. A difference system similar to (33) must be taken into account but with a slightly different right-hand side which considers formulas (28) and (39). Then, the local truncation error of this difference system is

$$\bar{\bar{\rho}}_{n,h} = \begin{bmatrix} \rho_{n+1,h} \\ O(\frac{1}{k}\rho_{n+1,h} + \frac{1}{k}\rho_{n+2,h}) \\ B_{h,0}^{-1}\dot{\rho}_{n,h} \end{bmatrix} = \begin{bmatrix} O(k^2\varepsilon_h + k^4) \\ O(k\varepsilon_h + k^3) \\ O(k^2\varepsilon_h + k^3) \end{bmatrix},$$

where, for the last equality, Theorems 6 and 11 have been applied. Then, considering again a formula similar to (34) with the same matrix $\bar{\bar{R}}^n(kB_{h,0})$ in (35), we get the following result.

**Theorem 12.** *Under the hypotheses of Theorem 7 and assuming also, for j=0,1, that $B_{h,0}^{-1}[P_h\dot{u}(t_j) - \dot{U}_h^j] = O(k^2 + \varepsilon_h)$, it happens that*

$$B_{h,0}^{-1}[P_h\dot{u}(t_n) - \dot{U}_h^n] = O(k^2 + \varepsilon_h).$$

# 6 Numerical experiments

The main aim of this section is to corroborate the results of the previous sections on the order reduction which turns up when integrating a nonlinear wave equation in time with Gautschi method in the classical way, and to check that the technique that we suggest to avoid it works. For that we have considered the following Dirichlet problem with parameter $w = 0.9$

$$
\begin{aligned}
u_{tt}(t,x) &= u_{xx}(t,x) - \sin(u(t,x)), \quad 0 \le x \le 1, \quad t \in [0,1], \\
u(0,x) &= 4\mathrm{atan}\Big(\frac{\sqrt{1-w^2}}{w}\frac{1}{\cosh(x\sqrt{1-w^2})}\Big), \\
u_t(0,x) &= 0, \\
u(t,0) &= 4\mathrm{atan}\Big(\frac{\sqrt{1-w^2}}{w}\cos(wt)\Big), \\
u(t,1) &= 4\mathrm{atan}\Big(\frac{\sqrt{1-w^2}}{w}\frac{\cos(wt)}{\cosh(\sqrt{1-w^2})}\Big),
\end{aligned}
\tag{44}
$$

which has as exact solution

$$
u(t,x) = 4\mathrm{atan}\Big(\frac{\sqrt{1-w^2}}{w}\frac{\cos(wt)}{\cosh(x\sqrt{1-w^2})}\Big).
$$

We remark that hypotheses (A1)-(A4) are satisfied for this problem when considering $H = H^2(0,1)$, $Y = \mathbb{C}^2$ and $\partial$ the Dirichlet trace operator. Moreover, in this case, $-A_0$ is negative selfadjoint.

For the space discretization of the problem we have considered the classical symmetric second-order finite difference scheme, so that

$$
A_{h,0} = \frac{1}{h^2}\mathrm{tridiag}[-1,2,-1], \quad A_h Q_h[g_0,g_1]^T = -\frac{1}{h^2}[g_0,0,\ldots,0,g_1]^T,
$$

and $P_h w$ corresponds to the interior nodal projection when $w$ is continuous. As $A_{h,0}$ is symmetric and using Gerschgorin theorem, its eigenvalues are real and positive. Moreover, a lower bound is well-known to exist for the smallest one when $h \to 0$. Besides, there exists a matrix $B_{h,0}$ which satisfies $B_{h,0}^2 = A_{h,0}$ and which has the same properties. Considering then that the matrices are real and symmetric, the Euclidean norm of $A_{h,0}^{-1}$, $B_{h,0}^{-1}$, $\cos(\tau B_{h,0})$, $\sin(\tau B_{h,0})$, $\mathrm{sinc}(\tau B_{h,0})$, $e^{i\tau B_{h,0}}$ and $(I - e^{i\tau B_{h,0}})^{-1}$ coincides with its spectral radius and therefore is bounded. Because of this, hypotheses (H1) and (H4) are satisfied. On the other hand, (H2) is satisfied with $Z = H^4(0,1)$ and $\varepsilon_h = O(h^2)$ and (H3) because of a discrete maximum principle for this discretization (look at Theorems 12.5.1 and 12.5.3 in [25] for the corroboration of both hypotheses for the five-point Laplacian. In our case, the result also applies with a similar but simpler argument). We have taken $h = 1/4000$ in the numerical experiments so that the error in space can be considered negligible and, although not shown here, we have checked that the errors remained practically the same when $h$ diminished.

| k | 0.2 | 0.1 | 0.05 | 0.025 |
|---|---|---|---|---|
| $\|P_h u(t_2) - U_h^2\|$ | 1.3962e-2 | 2.4629e-3 | 4.3389e-4 | 7.6213e-5 |
| Order | | 2.50 | 2.50 | 2.51 |
| $\|P_h u(T) - U_h^N\|$ | 4.5140e-2 | 2.0919e-2 | 1.0566e-2 | 5.5330e-3 |
| Order | | 1.11 | 0.99 | 0.93 |
| $\|B_{h,0}^{-1}[P_h u(t_2) - U_h^2]\|$ | 1.4051e-2 | 2.4809e-3 | 4.3767e-4 | 7.7272e-5 |
| Order | | 2.50 | 2.50 | 2.50 |
| $\|B_{h,0}^{-1}[P_h u(T) - U_h^N]\|$ | 5.6165e-2 | 3.1731e-2 | 1.6438e-2 | 8.1708e-3 |
| Order | | 0.83 | 0.95 | 1.01 |

Table 1: Local and global error for the solution and the derivative when integrating problem (44) with the classical approach (11)-(32) using finite differences in space

| k | 0.2 | 0.1 | 0.05 | 0.025 |
|---|---|---|---|---|
| $\|P_h u(t_2) - U_h^2\|$ | 1.8097e-5 | 1.7428e-6 | 1.1926e-7 | 7.6310e-9 |
| Order | | 3.38 | 3.87 | 3.97 |
| $\|P_h u(T) - U_h^N\|$ | 2.8633e-4 | 6.8468e-5 | 1.6935e-5 | 4.2227e-6 |
| Order | | 2.06 | 2.02 | 2.00 |
| $\|B_{h,0}^{-1}[P_h u(t_2) - U_h^2]\|$ | 1.0236e-4 | 2.0333e-5 | 2.7767e-3 | 3.5444e-7 |
| Order | | 2.33 | 2.87 | 2.97 |
| $\|B_{h,0}^{-1}[P_h u(T) - U_h^N]\|$ | 9.7941e-4 | 2.2391e-4 | 5.6404e-5 | 1.4173e-5 |
| Order | | 2.13 | 1.99 | 1.99 |

Table 2: Local and global error in the solution and the derivative when integrating problem (44) with the suggested technique (28)-(39) using finite differences in space

In Table 1 we show the results which correspond to integrate firstly in space and then in time when using exact starting values and measuring the error in the discrete $L^2$-norm. We can check that the order of the local error in the solution is at least 2, as Theorem 2 assures but it is not 4, as corresponds to the order of the method when no order reduction turns up. As for the global error in the solution, it is quite near 1, as Theorem 3 predicts. As for the derivative, local order 2 and global order 1 are observed when multiplying by $B_{h,0}^{-1}$, as it corresponds to Theorems 8 and 9. On the other hand, in Table 2, we can see the results which correspond to integrate firstly in time and then in space, and using then (28)-(39) especifically as a final formula, instead of (11)-(32). We can check that the order of the local error for the solution is very near 4 and that of the global error is very near 2, as stated by Theorems 6 and 7. Besides, local and global orders very near 2 and 3 are observed for the error in the derivative multiplied by $B_{h,0}^{-1}$, as it corresponds to Theorems 11 and 12. Moreover, the size of the errors is much smaller than with the classical approach, even for the biggest value of $k$.

Finally, we also offer a comparison in terms of computational cost because each of the formulas in (28)-(39) means to calculate one more term at each timestep than
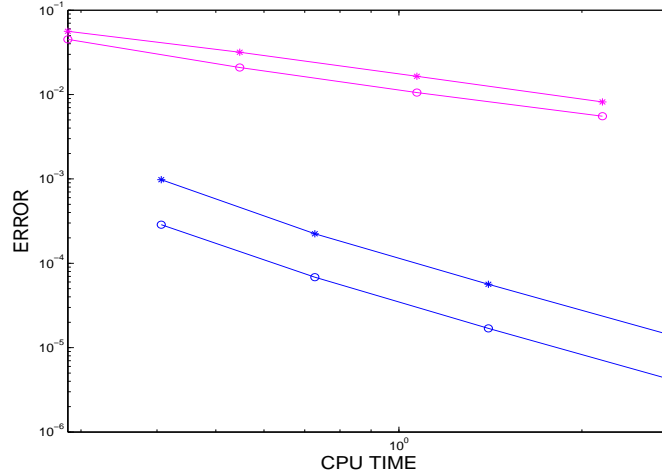
Figure 1: Error against CPU time for the solution (circles) and the derivative multiplied by $B_{h,0}^{-1}$ (asterisks) when integrating problem (44) with second-order finite differences in space and Gautschi method in time, with the classical approach (11)-(32) (pink) and the suggested technique (28)-(39) (blue)

with (11)-(32). In Figure 1 we show the global $L^2$-discrete error against CPU time for the same problem as above when calculating the terms which contain $\cos(kB_{h,0})$, $\text{sinc}(kB_{h,0})$, $\gamma_1(kB_{h,0})$ and $\gamma_2(kB_{h,0})$ through a discrete sine transform, as with fast Poisson solvers [23]. We can see that, with one second of CPU time, the errors in the solution are 300 smaller with the suggested approach than the classical approach. As for the error in the derivative multiplied by $B_{h,0}^{-1}$, the ratio is 200 smaller. Moreover, that advantage becomes bigger and bigger when the required accuracy is increased since the slope of the line corresponding to integrate firstly in time and then in space is bigger than the one which is obtained when doing things the other way round.

Doing the comparison with more general space discretizations through rational Krylov subspace methods would also be very interesting as a future research but it is out of the scope of this paper.

# 7   Acknowledgements

# References

[1] I. ALONSO-MALLO, B. CANO AND M. J. MORETA, *The stability of rational approximations of cosine functions in Hilbert spaces*, Appl. Numer. Math., Vol. 59

(2009) 21–38.

[2] I. Alonso–Mallo, B. Cano and N. Reguera, *Analysis of order reduction when integrating linear initial boundary value problems with Lawson methods*, Appl. Num. Math. 118 (2017) 64-74.

[3] I. Alonso-Mallo, B. Cano and N. Reguera, *Avoiding order reduction when integrating linear initial boundary value problems with Lawson methods*, IMA J. Num. Anal., doi: 10.1093/imanum/drw052.

[4] I. Alonso–Mallo, B. Cano and N. Reguera, *Avoiding order reduction when integrating linear initial boundary value problems with exponential splitting methods*, IMA J. Num. Anal., doi:10.1093/imanum/drx047

[5] I. Alonso-Mallo, B. Cano and N. Reguera, *Avoiding order reduction when integrating reaction-diffusion boundary value problems with exponential splitting methods*, arXiv:1705.01857, submitted for publication.

[6] W. Arendt, C.F.K. Batty, M. Hieber and F. Neubrander, *Vector-valued Laplace Transforms and Cauchy Problems*, Monographs in Mathematics, Vol. 96, Birkhäuser, Basel, 2001.

[7] B. Cano and M. J. Moreta, *Multistep cosine methods for second-order partial differential equations*, IMA J. Num. Anal. 30 (2010) 431–461.

[8] B. Cano and M. J. Moreta, *Exponential quadrature rules without order reduction for integrating linear initial boundary value problems*, to be published in SIAM J. Num. Anal.

[9] B. Cano and N. Reguera, *Avoiding order reduction when integrating nonlinear Schrödinger equation with Strang method*, J. Comp. Appl. Math., 316 (2017), 86–99.

[10] M. Crouzeix, *Operators with numerical range in a parabola*, Arch. Math. 82 (2004) 517–527.

[11] M. Crouzeix, *Numerical range and functional calculus in Hilbert spaces*, J. Funct. Anal. 244 (2007) 668–690.

[12] L. Einkemmer and A. Ostermann, *Overcoming order reduction in diffusion-reaction splitting. Part 1: Dirichlet boundary conditions*, SIAM J. Sci. Comput. **37** (3) (2015), A1577–A1592.

[13] L. Einkemmer and A. Ostermann, *Overcoming order reduction in diffusion-reaction splitting. Part 2: Oblique boundary conditions*, SIAM J. Sci. Comput. **38** (2016), A3741–A3757.

[14] E. Faou, A. Ostermann and K. Schratz, *Analysis of exponential splitting methods for inhomogeneous parabolic equations*, IMA J. Numer. Anal. **35** (1) (2015), 161–178.

[15] L. GAUCKLER, *Error analysis of trigonometric integrators for semilinear wave equations*, SIAM J. Numer. Anal. **53**, No. 2, (2015) 1082–1106.

[16] W. GAUTSHI, *Numerical integration of ordinary differential equations based on trigonometric polynomials*, Numer. Math. **3** (1961) 381–397.

[17] V. GRIMM, *On error bounds for the Gautschi-type exponential integrator applied to oscillatory second-order differential equations*, Numer. Math. **100** (2005) 71–89.

[18] V. GRIMM, *A note on the Gautschi-type method for oscillatory second-order differential equations*, Numer. Math. **102** (2005) 61–66.

[19] V. GRIMM AND M. HOCHBRUCK, *Rational approximation to trigonometric operators*, BIT Numer. Math. **48** (2008) 215–229.

[20] M. HAASE, *The Functional Calculus for Sectorial Operators*, book manuscript, available at www.mathematik.uni-ulm.de/m5/haase, to appear in: Operator Theory: Advances and Applications, Birkhäuser Verlag.

[21] E. HAIRER, CH. LUBICH AND G. WANNER, *Geometric Numerical Integration*, Springer-Verlag, 2002.

[22] M. HOCHBRUCK AND A. OSTERMANN, *Exponential integrators*, Acta Numerica (2010) 209-286.

[23] A. ISERLES, *A First Course in the Numerical Analysis of Differential Equations*, Cambridge University Press, Cambridge, 2008.

[24] J. L. LIONS AND E. MAGENES, *Non-homogeneous Boundary Value Problems and Applications*, Vol. II, Springer-Verlag, Berlin/Heidelberg/New York, 1972.

[25] J. C. STRIKWERDA, *Finite Difference Schemes and Partial Differential Equations*, (Wadsworth & Brooks, United States of America, 1989).

# 8 Appendix: More general abstract space discretizations

In this section, we will see that the technique to avoid order reduction can also be applied with more general space discretizations, which include finite elements and other more elaborate finite differences. For that, we will consider an abstract space discretization framework, which was also used in [3] when analysing how to avoid order reduction in first-order linear problems with Lawson methods. In contrast, here we are considering second-order nonlinear problems and Gautschi method.

In this section, we will assume that $Q_h : Y \to X_h$ is an interpolator operator on the boundary, that there exists a projection operator $L_h : X \to X_{h,0}$ and we will reserve the notation $P_h$ for $L_h - L_h Q_h$. Then, the projection on $X_{h,0}$ of the operator $A$ will be

approximated by $A_h : X_h \to X_{h,0}$ and that of $A_0$ by $A_{h,0} : X_{h,0} \to X_{h,0}$. In such a way, the elliptic projection $R_h : X \to X_{h,0}$ corresponding to (8) will be now the solution of

$$-A_{h,0}R_hw - A_hQ_hg = L_hF = P_hF + L_hQ_h\partial F.$$

We will assume that hypotheses (H1),(H3),(H4),(H5) are also valid substituting the word matrix by operator and $\mathbb{C}^{N_0}$ by $X_{h,0}$. Moreover, we will substitute hypothesis (H2) by

(H2') There exists a subspace $Z \subset D(A)$ with norm $\| \cdot \|_Z$ such that, for $u \in Z$,

$$\|(R_h - P_h)u\|_h \le \varepsilon_h\|u\|_Z,$$

for $\varepsilon_h$ decreasing with $h$,

and we will add the following hypothesis

(H6) There exist constants $C$ and $C'$ such that, for small enough $h$, for each $u \in H$ and $v \in Y$,
$$\|L_hu\|_h \le C\|u\|, \quad \|Q_hv\| \le C'\|v\|.$$
Besides, for each $u \in H$ such that $\partial u = g$,

$$\|u - Q_hg - L_h(u - Q_hg)\| \le \eta_h\|u\|,$$

for $\eta_h$ decreasing with $h$.

When discretizing now (21) and (25), the following systems arise instead of (22) and (26):

$$
\begin{aligned}
\ddot{V}_{h,n}(s) &= -A_{h,0}V_{h,n}(s) - A_hQ_h\partial[u(t_n) - \frac{s^2}{2}Au(t_n)] + L_hQ_h\partial Au(t_n), \\
V_{h,n}(0) &= U_h^n, \\
\dot{V}_{h,n}(0) &= 0,
\end{aligned}
\tag{45}
$$

$$
\begin{aligned}
\ddot{W}_{h,n}(\tau) &= -A_{h,0}W_{h,n}(\tau) - \tau A_hQ_h\partial f(t_n, u(t_n)), \\
W_{h,n}(0) &= 0, \\
\dot{W}_{h,n}(0) &= P_hf(t_n, U_h^n + Q_hg(t_n)),
\end{aligned}
\tag{46}
$$

and then our suggestion for Gautschi method when integrating (1) would be $U_h^n + Q_hg(t_n)$ where

$$
\begin{aligned}
U_h^{n+1} &- 2\cos(kB_{h,0})U_h^n + U_h^{n-1} \\
&= k^2\gamma_1(kB_{h,0})[P_hf(t_n, U_h^n + Q_hg(t_n)) - A_hQ_hg(t_n) + L_hQ_h\partial Au(t_n)] \\
&\quad -k^4\gamma_2(kB_{h,0})A_hQ_h\ddot{g}(t_n).
\end{aligned}
\tag{47}
$$

## 8.1 Full discretization local truncation error

Instead of (30), we define

$$\rho_{n,h} = R_h u(t_{n+1}) - 2\bar{V}_{h,n}(k) + R_h u(t_{n-1}) - 2\int_0^k \bar{W}_{h,n}(k-s)ds. \qquad (48)$$

where $\bar{V}_{h,n}$ (resp. $\bar{W}_{h,n}$) are the solutions of (45) (resp. (46)) when $U_h^n$ is substituted by $R_h u(t_n)$. Then, we have the following result:

**Theorem 13.** *Under the same hypotheses of regularity of Theorem 6, but assuming the new hypotheses on the space discretization (H2') and (H6), $\rho_{n,h} = O(k^2\varepsilon_h + k^2\eta_h + k^4)$.*

*Proof.* In a similar way to the proof of Theorem 6,

$$
\begin{aligned}
\ddot{\bar{V}}_{h,n}(s) - R_h\ddot{\hat{v}}_n(s) &= -A_{h,0}\bar{V}_{h,n}(s) - A_hQ_h\partial\hat{v}_n(s) + L_hQ_h\partial Au(t_n) + R_hAu(t_n)\\
&= -A_{h,0}(\bar{V}_{h,n}(s) - R_h\hat{v}_n(s)) - L_hA\hat{v}_n(s) + L_hQ_h\partial Au(t_n) + R_hAu(t_n)\\
&= -A_{h,0}(\bar{V}_{h,n}(s) - R_h\hat{v}_n(s)) + (R_h - P_h)Au(t_n) + \frac{s^2}{2}L_hA^2u(t_n),\\
\bar{V}_{h,n}(0) - R_h\hat{v}_n(0) &= 0,\\
\dot{\bar{V}}_{h,n}(0) - R_h\dot{\hat{v}}_n(0) &= 0,
\end{aligned}
$$

and

$$
\begin{aligned}
\ddot{\bar{W}}_{h,n}(\tau) - R_h\ddot{\hat{w}}_n(\tau) &= -A_{h,0}\bar{W}_{h,n}(\tau) - A_hQ_h\partial\hat{w}_n(\tau)\\
&= -A_{h,0}(\bar{W}_{h,n}(\tau) - R_h\hat{w}_n(\tau)) - L_hA\hat{w}_n(\tau),\\
\bar{W}_{h,n}(0) - R_h\hat{w}_n(0) &= 0,\\
\dot{\bar{W}}_{h,n}(0) - R_h\dot{\hat{w}}_n(0) &= P_hf(t_n, R_hu(t_n) + Q_hg(t_n)) - R_hf(t_n, u(t_n)).
\end{aligned}
$$

Then,

$$
\begin{aligned}
\bar{V}_{h,n}(k) - R_h\hat{v}_n(k) &= \int_0^k B_{h,0}^{-1}\sin((k-s)B_{h,0})[(R_h - P_h)Au(t_n) + \frac{s^2}{2}L_hA^2u(t_n)]ds\\
&= \int_0^k (k-s)\operatorname{sinc}((k-s)B_{h,0})[O(\varepsilon_h) + O(s^2)]ds = O(k^2\varepsilon_h + k^4),
\end{aligned}
$$

where (H2') and (H6) have been used; and

$$
\begin{aligned}
\bar{W}_{h,n}(\tau) - R_h\hat{w}_n(\tau) &= \tau\operatorname{sinc}(\tau B_{h,0})[P_hf(t_n, R_hu(t_n) + Q_hg(t_n)) - R_hf(t_n, u(t_n))]\\
&\quad - \int_0^\tau B_{h,0}^{-1}\sin((\tau-\sigma)B_{h,0})L_hA\hat{w}_n(\sigma)d\sigma\\
&= \tau\operatorname{sinc}(\tau B_{h,0})\Big[P_h[f(t_n, R_hu(t_n) + Q_hg(t_n)) - f(t_n, u(t_n))] + (P_h - R_h)f(t_n, u(t_n))\Big]\\
&\quad - \int_0^\tau (\tau-\sigma)\operatorname{sinc}((\tau-\sigma)B_{h,0})\sigma L_hAf(t_n, u(t_n))d\sigma = O(\tau\eta_h + \tau\varepsilon_h + \tau^3),
\end{aligned}
$$

where, for the last equality, (H6) and the fact that $f$ is globally Lipschitz have been considered. More precisely, we have taken into account that

$$\|R_h u(t_n) + Q_h g(t_n) - u(t_n)\|$$
$$\leq \|u(t_n) - Q_h g(t_n) - P_h u(t_n)\| + \|P_h u(t_n) - R_h u(t_n)\| = O(\eta_h + \varepsilon_h).$$

Considering now (48), in the same way than in the proof of Theorem 6,

$$\rho_{n,h} = R_h \rho_n + O(k^2 \varepsilon_h + k^2 \eta_h + k^4),$$

from what the result follows using again Theorem 5 and that $R_h$ is uniformly bounded because of (H2') and (H6).  □

## 8.2  Full discretization global error

Considering the previous result on the local error for the full discretization under the more general hypotheses for the space discretization, the following result follows for the global error using the same proof as that of Theorem 7 and the fact that

$$P_h u(t_n) - U_h^n = (P_h - R_h) u(t_n) + R_h u(t_n) - U_h^n.$$

**Theorem 14.** *Under the hypotheses of Theorems 5 and 13, if the starting values are such that they satisfy $U_h^j - P_h u(t_j) = O(k^3 + k\varepsilon_h + k\eta_h)$ $(j = 0, 1)$, it happens that*

$$P_h u(t_n) - U_h^n = O(k^2 + \varepsilon_h + \eta_h).$$

## 8.3  Approximation for the derivative

The approximation for the derivative with this more general space discretization is more similar to that in Section 5 because the second derivative with respect to $s$ of the boundaries in (37) vanish. Therefore, the only difference is in the initial condition for $\dot{\sigma}_{n,h}$ in (38), which should be now $P_h f(t_n, U_h^n + Q_h g(t_n))$. Then, the approximation for the derivative is $\dot{U}_h^n + Q_h \dot{g}(t_n)$, where

$$\dot{U}_h^{n+1} - \dot{U}_h^{n-1} = 2k\mathrm{sinc}(kB_{h,0})[-A_{h,0}U_h^n - A_h Q_h g(t_n) + P_h f(t_n, U_h^n + Q_h g(t_n))]$$
$$-2kB_{h,0}^{-2}[I - \mathrm{sinc}(kB_{h,0})]A_h Q_h \partial \ddot{g}(t_n).$$

Moreover, the result for the global error would be

**Theorem 15.** *Under the hypotheses of Theorem 14, if the starting values are such that $P_h \dot{u}(t_j) - \dot{U}_h^j = O(k^2 + \varepsilon_h + \eta_h)$, it happens that*

$$B_{h,0}^{-1}[P_h u(t_n) - U_h^n] = O(k^2 + \varepsilon_h + \eta_h).$$