



High-order linearly implicit exponential integrators conserving quadratic invariants with application to scalar auxiliary variable approach

Shun Sato¹

Received: 21 August 2023 / Accepted: 11 February 2024 / Published online: 29 February 2024
© The Author(s) 2024, corrected publication 2024

Abstract

This paper proposes a framework for constructing high-order linearly implicit exponential integrators that conserve a quadratic invariant. This is then applied to the scalar auxiliary variable (SAV) approach. Quadratic invariants are significant objects that are present in various physical equations and also in computationally efficient conservative schemes for general invariants. For instance, the SAV approach converts the invariant into a quadratic form by introducing scalar auxiliary variables, which have been intensively studied in recent years. In this vein, Sato et al. (Appl. Numer. Math. **187**, 71–88 2023) proposed high-order linearly implicit schemes that conserve a quadratic invariant. In this study, it is shown that their method can be effectively merged with the Lawson transformation, a technique commonly utilized in the construction of exponential integrators. It is also demonstrated that combining the constructed exponential integrators and the SAV approach yields schemes that are computationally less expensive. Specifically, the main part of the computational cost is the product of several matrix exponentials and vectors, which are parallelizable. Moreover, we conduct some mathematical analyses on the proposed schemes.

Keywords Ordinary differential equations · Quadratic invariants · Geometric numerical integration · Exponential integrators · Scalar auxiliary variable approach

Mathematics Subject Classification (2010) 65L05 · 65M06 · 65P10

✉ Shun Sato
shun@mist.i.u-tokyo.ac.jp

¹ Department of Mathematical Informatics, Graduate School of Information Science and Technology, The University of Tokyo, Bunkyo-ku, Tokyo 1138656, Japan

1 Introduction

In this paper, we consider ordinary differential equations (ODEs)

$$\dot{y} = S(y)\nabla V(y), \quad (1)$$

on a finite-dimensional real Hilbert space $(\mathcal{X}, \langle \cdot, \cdot \rangle_{\mathcal{X}})$, where $y: [0, T] \rightarrow \mathcal{X}$ is a dependent variable, $S: \mathcal{X} \rightarrow \mathcal{L}(\mathcal{X})$ is a “skew-symmetric matrix function” ($\mathcal{L}(\mathcal{X})$ denotes the space of linear operators on \mathcal{X}), i.e., $\langle x, S(z)y \rangle_{\mathcal{X}} = -\langle S(z)x, y \rangle_{\mathcal{X}}$ holds for all $x, y, z \in \mathcal{X}$, $V: \mathcal{X} \rightarrow \mathbb{R}$ is a differentiable function, and $\nabla V: \mathcal{X} \rightarrow \mathcal{X}$ denotes the gradient of V with respect to the inner product of \mathcal{X} . Although the methods in this paper can be applied even when \mathcal{X} is infinite-dimensional, we limit ourselves to finite-dimensional cases for the sake of mathematical analysis (see Remark 4).

The class of ODEs in the form (1) includes many examples: Hamiltonian systems, Poisson systems, and spatial discretization of variational partial differential equations (PDEs) (see, e.g., [1]). These ODEs satisfy the conservation law with respect to V :

$$\frac{d}{dt} V(y(t)) = \langle \nabla V(y(t)), \dot{y}(t) \rangle_{\mathcal{X}} = \langle \nabla V(y(t)), S(y(t))\nabla V(y(t)) \rangle_{\mathcal{X}} = 0,$$

where the last equality holds due to the skew-symmetry of $S(y(t))$.

Since this conservation law is an important property of the differential equation (1), numerical methods that preserve it have been studied in the literature, for example, the discrete gradient method [2–4] for the gradient ODEs and the discrete variational derivative method [5, 6] (see also [7]) for variational PDEs. In addition, Cohen and Hairer [8] proposed a high-order extension of the discrete gradient method.

Since these conservative numerical schemes are fully implicit, several techniques have been devised to reduce computational cost. For example, Besse [9] and Zhang, Pérez-García, and Vázquez [10] proposed linearly implicit conservative schemes for the nonlinear Schrödinger equation. For polynomial invariants, Matsuo and Furihata [11] proposed multistep linearly implicit DVDM (see also Dahlby and Owren [12]). Recently, Yang and Han [13] proposed the invariant energy quadratization (IEQ) approach, and Shen et al. [14] proposed the scalar auxiliary variable (SAV) approach (see also [15] and the references therein). The SAV approach employs scalar auxiliary variables to convert the invariant into a quadratic form, which is then preserved by a linearly implicit scheme (see Sect. 2.4 for details).

Although the above computationally inexpensive methods differ, they have in common that they attribute the invariant to a quadratic function in some way. Given this, Sato et al. [16] recently proposed a framework for constructing high-order and linearly implicit schemes conserving a quadratic invariant.

Exponential integrators (cf. [17]) are efficient numerical methods for solving semi-linear differential equations, and combining it with geometric numerical integration has also been studied. Celledoni et al. [18] deal with L^2 norm conservation for the Schrödinger equation and derive the condition to preserve it in exponential Runge–Kutta methods. Mei et al. [19] deal with (1) with $V(y) = \frac{1}{2}\langle y, y \rangle_{\mathcal{X}}$ and construct conservative exponential integrators. Some researchers deal with the equation in the

form $\dot{y} = J(Ly + \nabla E(y))$, where J is skew-symmetric, L is symmetric, and E is a function: Li and Wu [20] propose second-order exponential discrete gradient schemes; Mei et al. [21] discuss how to design high-order conservative schemes based on Li and Wu [20], modified differential equation, and order condition in terms of B-series; and Li [22] proposes a multistep linearly implicit conservative exponential scheme based on polarization.

Combinations of exponential integrators and the SAV approach have also been investigated [23–26]. In particular, Jiang et al. [27] propose high-order linearly implicit structure-preserving exponential integrators for the nonlinear Schrödinger equation based on the SAV approach and the Lawson transformation [28]. They also mention that the same method can be applied to some general differential equations, and numerically confirm the superiority of their method. However, no theoretical guarantee of high accuracy is given. In this paper, we reveal that linearly implicit high-order conservative schemes proposed by Sato et al. [16] can be combined with exponential integrators (Sect. 3). The resulting scheme, although linearly implicit, is computationally not very cheap, since it includes matrix exponentials in the coefficient matrix. However, when combined with the SAV approach, its structure can be exploited to provide a computationally efficient implementation (Sect. 4). Specifically, the main part of the computational cost is a few products of matrix exponential functions and vectors with the size of the dimension of the given differential equation. Furthermore, these products can be computed in parallel. Theoretical guarantees such as accuracy are also provided, albeit being limited to finite dimensions. (The relationship with [27] is discussed in Remark 6.)

Here, we note a limitation of the proposed method: it cannot be applied to dissipative systems. In many cases, if a conservative numerical method can be constructed, a dissipative numerical method, one that replicates the dissipation law, can be constructed correspondingly: [16] can also be applied to dissipative cases. However, this is not the case for the proposed method in the present paper. The proposed method is essentially driven by the fact that under a mild assumption, the quadratic invariants are invariant by the Lawson transformation (see Lemma 4), and the extension to dissipative cases is nontrivial (see Remark 2).

The remainder of the paper is organized as follows. In Sect. 2, we briefly review the canonical Runge–Kutta methods, linearly implicit high-order conservative schemes proposed by Sato et al. [16], the Lawson transformation, and the SAV approach. The contents in Sects. 3 and 4 are already described above. The proposed schemes are numerically examined in Sect. 5.

2 Preliminaries

2.1 Runge–Kutta methods and quadratic invariants

For the autonomous system

$$\dot{y} = f(y),$$

general Runge–Kutta methods that compute an approximation $y_1 \approx y(h)$ from $y_0 = y(0)$ can be written as

$$\begin{cases} Y_i = y_0 + h \sum_{j \in [s]} a_{ij} f(Y_j) & (i \in [s]), \\ y_1 = y_0 + h \sum_{i \in [s]} b_i f(Y_i), \end{cases}$$

where $[s] := \{1, 2, \dots, s\}$. Throughout this paper, the fixed step size is considered, but as is common in the context of conservative methods, variable step sizes can also be employed without difficulty.

A subclass of Runge–Kutta methods preserves all quadratic invariants.

Proposition 1 (Cooper [29]) *Runge–Kutta methods satisfying*

$$b_i a_{ij} + b_j a_{ji} = b_i b_j \quad i, j \in [s] \tag{2}$$

automatically preserve all quadratic invariants.

The Runge–Kutta methods satisfying (2) are said to be *canonical* (see [30] for details on canonical Runge–Kutta methods). Note that a canonical Runge–Kutta method must be implicit, yet diagonally implicit Runge–Kutta methods can be canonical.

For example, the second-order Gauss, fourth-order Gauss, and third-order diagonally implicit canonical Runge–Kutta methods are as follows:

$$\begin{array}{c|c} \frac{1}{2} & \frac{1}{2} \\ \hline & 1 \end{array}, \quad \begin{array}{c|cc} \frac{1}{2} - \frac{\sqrt{3}}{6} & \frac{1}{4} & \frac{1}{4} - \frac{\sqrt{3}}{6} \\ \hline \frac{1}{2} + \frac{\sqrt{3}}{6} & \frac{1}{4} + \frac{\sqrt{3}}{6} & \frac{1}{4} \end{array}, \quad \begin{array}{c|cc} \frac{\alpha}{2} & \frac{\alpha}{2} & 0 \\ \frac{3}{2}\alpha & \alpha & \frac{\alpha}{2} \\ \hline \frac{1}{2} + \alpha & \alpha & \alpha \frac{1}{2} - \alpha \end{array}, \quad \begin{array}{c|cc} & \alpha & 0 \\ \hline \alpha & \alpha & 1 - 2\alpha \end{array},$$

where $\alpha = \frac{1}{3} \left(2 + \frac{1}{2^{1/3}} + 2^{1/3} \right)$.

2.2 High-order linearly implicit schemes that conserve quadratic invariants

In this section, we briefly review the high-order linearly implicit schemes for the gradient system (1) with a quadratic invariant V proposed by [16].

Definition 1

Step 0 Prepare $Y_i^{(0)} \approx y(c_i h)$ and set $k = 1$.

Step 1 Solve the linear equation system

$$Y_i^{(k)} = y_0 + h \sum_{j \in [s]} a_{ij} S \left(Y_j^{(k-1)} \right) \nabla V \left(Y_j^{(k)} \right) \quad (i \in [s]) \tag{3}$$

to obtain $Y_i^{(k)}$'s. If some criteria hold, go to Step 2. Otherwise, set $k = k + 1$ and repeat Step 1.

Step 2 Output

$$y_1^{(k)} = y_0 + h \sum_{j \in [s]} b_j S \left(Y_j^{(k-1)} \right) \nabla V \left(Y_j^{(k)} \right).$$

Theorem 2 Suppose that A and b satisfy (2), and V is quadratic. Then, the solution $y_1^{(k)}$ of Definition 1 satisfies $V \left(y_1^{(k)} \right) = V(y_0)$ for any $h > 0$.

Theorem 3 Assume the following conditions:

- (A1) $Y_i^{(0)}$ satisfies $\|Y_i^{(0)} - y(c_i h)\|_{\mathcal{X}} \leq Ch^q$ for each $i \in [s]$.
- (A2) $S: \mathcal{X} \rightarrow \mathcal{L}(\mathcal{X})$ is Lipschitz continuous.
- (A3) The base Runge–Kutta method is of order p .

Then, the numerical solution $y_1^{(k)}$ of the scheme defined by Definition 1 satisfies

$$\|y_1^{(k)} - y(h)\|_{\mathcal{X}} \leq C'h^{\min\{p, q+k-1\}+1}$$

for sufficiently small $h > 0$, i.e., the scheme with k iteration is of order $\min\{p, q+k-1\}$. Here, C' is a constant depending only on the exact solution y , k , S , V , A , b , and the constant C in (A1).

Remark 1 In view of Theorem 3, $k = p - q + 1$ is a natural choice as the criterion in Step 1 of Definition 1. In addition, there are several other possibilities. As shown in [16, Theorem 4.1], $Y_i^{(k)}$ linearly coversges to the corresponding inner stage of the base Runge–Kutta method. Based on this fact, one can choose k so that the update $\|Y_i^{(k)} - Y_i^{(k-1)}\|_{\mathcal{X}}$ is sufficiently small. Another possibility is to choose k so that $\|y_1^{(k)} - y_1^{(k-1)}\|_{\mathcal{X}}$, an approximation of the local error, is sufficiently small (note that, when $k \leq p - q + 1$, $y_1^{(k)}$ and $y_1^{(k-1)}$ are approximations of $y(h)$ with different order of accuracy, which are known to be useful to approximate the local error). However, in the numerical experiments in Sect. 5, we simply choose some fixed k to verify the theoretical order of accuracy.

2.3 Lawson transformation

We briefly review the Lawson transformation [28] which is useful to construct exponential integrators. For the ODE in the form

$$\dot{y} = My + f(y), \tag{4}$$

we consider the transformation $w(t) = \exp(-tM)y(t)$ to obtain

$$\dot{w} = \exp(-tM)f(\exp(tM)w(t)), \tag{5}$$

where \exp denotes the matrix exponential function. It is known that, when M has large eigenvalues, the transformed system (5) is easier to solve numerically than the original ODE. Therefore, a good approach to solve the system (4) is to apply a numerical

method to (5) and use the inverse transformation to return to the original variable. If a Runge–Kutta method is adopted as the numerical method, the resulting numerical method is written as

$$\begin{cases} Y_i = \exp(c_i h M) y_0 + h \sum_{j \in [s]} a_{ij} \exp((c_i - c_j) h M) f(Y_j) & (i \in [s]), \\ y_1 = \exp(h M) y_0 + h \sum_{i \in [s]} b_i \exp((1 - c_i) h M) f(Y_i). \end{cases} \quad (6)$$

Hereafter, the scheme is designated the *Lawson method*.

2.4 Scalar auxiliary variable approach

In this section, we review the scalar auxiliary variable (SAV) approach for the Hamiltonian system

$$\dot{u} = J \nabla H(u) \quad (7)$$

on a finite-dimensional real Hilbert space $(\mathcal{V}, \langle \cdot, \cdot \rangle)$, where $u: [0, T) \rightarrow \mathcal{V}$ is a dependent variable, $J \in \mathcal{L}(\mathcal{V})$ is a (constant) skew-symmetric linear operator, and $H: \mathcal{V} \rightarrow \mathbb{R}$ is written in the form

$$H(u) = \frac{1}{2} \langle Lu, u \rangle + E(u),$$

where $L \in \mathcal{L}(\mathcal{V})$ is symmetric.

Here, for brevity, we assume that the function $E: \mathcal{V} \rightarrow \mathbb{R}$ is bounded from below, i.e., $\alpha := -\inf E(u) + \epsilon$ is finite (see [15] for unbounded cases). Using this assumption, we introduce a scalar auxiliary variable $r := \sqrt{E(u) + \alpha}$. Then, (7) is rewritten as

$$\begin{cases} \dot{u} = J(Lu + 2r\phi(u)), \\ \dot{r} = \langle \phi(u), \dot{u} \rangle, \end{cases} \quad (8)$$

where $\phi(u) := \nabla E(u) / (2\sqrt{E(u) + \alpha})$.

The system (8) has a quadratic modified invariant $V(u, r) = \frac{1}{2} \langle Lu, u \rangle + r^2 - \alpha$. Since quadratic invariants are much easier to preserve in numerical schemes than general invariants, this property enables us to construct computationally efficient conservative schemes (see, e.g., [14]).

A comprehensive way to construct conservative schemes based on the SAV approach is given in [15]; the reformulated system (8) can be rewritten for it to be regarded as a special case of (1). To this end, the inner product space \mathcal{X} is defined by $\mathcal{X} = \mathcal{V} \times \mathbb{R}$, and the associated inner product is defined as $\langle (x_1, r_1), (x_2, r_2) \rangle_{\mathcal{X}} = \langle x_1, x_2 \rangle + r_1 r_2$. Then, the system (8) can be written as

$$\frac{d}{dt} \begin{bmatrix} u \\ r \end{bmatrix} = \begin{bmatrix} I \phi(u) \\ 1 \end{bmatrix}^* \begin{bmatrix} J \\ 0 \end{bmatrix} \begin{bmatrix} I \phi(u) \\ 1 \end{bmatrix} \nabla V(u, r), \quad (9)$$

where the superscript $*$ denotes the adjoint.

Since the scheme in Sect. 2.2 can be applied to this system, such reformulation enables the construction of high-order linearly implicit schemes that conserve the modified invariant V . However, high-order schemes require solving large linear equation systems (see Appendix 4 for details).

3 Exponential Runge–Kutta methods conserving quadratic invariants

Let us consider the system

$$\dot{y} = My + S(y)\nabla V(y), \tag{10}$$

where $S: \mathcal{X} \rightarrow \mathcal{L}(\mathcal{X})$ is a skew-symmetric matrix function, and V is a quadratic function ($V(y) = \frac{1}{2}\langle y, Qy \rangle_{\mathcal{X}}$). We further assume that V is also an invariant of the linear part $\dot{y} = My$. In Sect. 3.1, we show that the Lawson method based on the canonical Runge–Kutta method conserves the quadratic invariant V . Then, Sect. 3.2 shows that the combination of the Lawson transformation and the scheme defined by Definition 1 yields a linearly implicit exponential integrators that conserves the quadratic invariant V . Although the linearly implicit exponential integrators are not always computationally efficient, but as we will see in Sect. 4, it works very well with the SAV approach.

3.1 Lawson transformation and quadratic invariants

The following lemma, which is a slight extension of [20, Lemma 2.2], is crucial in constructing conservative schemes (see Remark 2).

Lemma 4 *Let $V: \mathcal{X} \rightarrow \mathbb{R}$ be a quadratic function, i.e., $V(y) = \frac{1}{2}\langle y, Qy \rangle_{\mathcal{X}}$, where $Q \in \mathcal{L}(\mathcal{X})$ is symmetric. Then, V is an invariant of the linear ODE $\dot{y} = My$ if and only if $(\exp(tM))^* Q \exp(tM) = Q$ holds for any $t \in \mathbb{R}$.*

Proof Since $y(t) = \exp(tM)y(0)$ holds for a solution of the linear ODE $\dot{y} = My$, $V(y(t))$ can be written as

$$\begin{aligned} V(y(t)) &= \frac{1}{2}\langle \exp(tM)y(0), Q \exp(tM)y(0) \rangle_{\mathcal{X}} \\ &= \frac{1}{2}\langle y(0), (\exp(tM))^* Q \exp(tM)y(0) \rangle_{\mathcal{X}}. \end{aligned}$$

Therefore, V is an invariant, i.e., $V(y(t)) = V(y(0))$ holds for any $t \in \mathbb{R}$ and initial value $y(0)$, if and only if $(\exp(tM))^* Q \exp(tM) = Q$ holds for any $t \in \mathbb{R}$. \square

Assumption that V is an invariant of the linear ODE $\dot{y} = My$ may seem restrictive, but there are many examples satisfying this assumption. For example, since the difference operators often commute with each other, ODEs obtained by a finite difference discretization of PDEs often satisfy this assumption. Moreover, as shown in Sect. 4, ODEs obtained by the SAV approach satisfy this assumption.

Remark 2 In [20], it is assumed that “ $M = SQ$ holds with a skew-symmetric matrix S ” instead of “ V is an invariant.” The latter condition is weaker in the sense that, even when V is an invariant of the linear ODE $\dot{y} = My$, we cannot conclude the existence of a skew-symmetric matrix S that satisfies $M = SQ$; the pair

$$Q = \begin{bmatrix} 1 & & \\ & 1 & \\ & & 0 \end{bmatrix}, \quad M = \begin{bmatrix} & 1 & \\ -1 & & \\ & & 1 \end{bmatrix}$$

is a counterexample.

The lemma in [20] also deals with dissipative cases. Lemma 4 can also be extended to the dissipative cases. Under the setting of Lemma 4, if V is a weak Lyapunov function, then $(\exp(tM))^* Q \exp(tM) \leq Q$ holds for any $t > 0$. However, extending the following theorem to dissipative systems is nontrivial and will be the subject of future research.

By using the lemma above, we show a sufficient condition for conserving quadratic invariants.

Theorem 5 Suppose that V is a quadratic invariant of the semilinear ODE (4). Assume that V is also an invariant of the linear part $\dot{y} = My$. In addition, we assume (A, b) satisfies (2). Then, the solution y_1 of the Lawson method (6) satisfies $V(y_1) = V(y_0)$.

Proof Using the Lawson transformation $w(t) = \exp(-tM)y(t)$ and Lemma 4, we see

$$\begin{aligned} V(y(t)) &= \frac{1}{2} \langle y(t), Qy(t) \rangle_{\mathcal{X}} \\ &= \frac{1}{2} \langle \exp(tM)y(t), Q \exp(tM)y(t) \rangle_{\mathcal{X}} \\ &= \frac{1}{2} \langle w(t), Qw(t) \rangle_{\mathcal{X}}. \end{aligned}$$

Since this implies the transformed ODE (5) also has a quadratic invariant in the form $\frac{1}{2} \langle w, Qw \rangle_{\mathcal{X}}$, a canonical RK method applied to the system (5) preserves the invariant (Proposition 1). \square

3.2 Linearly implicit exponential integrators conserving quadratic invariants

As shown in the previous section, a class of Lawson methods automatically preserves quadratic invariants under the additional assumption that the linear part also preserves the invariants. However, since all canonical RK methods are implicit, the resulting scheme must be fully implicit. Therefore, in this section, we construct linearly implicit schemes conserving quadratic invariants. To this end, we apply the schemes introduced in Sect. 2.2 instead of the Lawson methods.

For (10), the Lawson transformation $w(t) = \exp(-tM)y(t)$ gives

$$\dot{w} = \exp(-tM)S(\exp(tM)w(t)) Q \exp(tM)w(t).$$

Using Lemma 4 and $(\exp(tM))^{-1} = \exp(-tM)$, we can further rewrite the equation as

$$\dot{w} = \exp(-tM)S(\exp(tM)w(t))(\exp(-tM))^* \nabla V(w(t)).$$

Since the map $w \mapsto \exp(-tM)S(\exp(tM)w)(\exp(-tM))^*$ is a skew-symmetric matrix function, we can apply the scheme in Sect. 2.2. The inverse transformation reads the following scheme.

Definition 2

Step 0 Prepare $Y_i^{(0)} \approx y(c_i h)$ and set $k = 1$.

Step 1 Solve the linear equation system

$$Y_i^{(k)} = \exp(c_i h M) y_0 + h \sum_{j \in [s]} a_{ij} \exp((c_i - c_j)hM) S(Y_j^{(k-1)}) \nabla V(Y_j^{(k)}) \quad (i \in [s]) \quad (11)$$

to obtain $Y_i^{(k)}$'s. If some criteria hold, go to Step 2. Otherwise, set $k = k + 1$ and repeat Step 1.

Step 2 Output

$$y_1^{(k)} = \exp(hM)y_0 + h \sum_{j \in [s]} b_j \exp((1 - c_j)hM) S(Y_j^{(k-1)}) \nabla V(Y_j^{(k)}).$$

Similar to Theorem 5, the discrete conservation law for the scheme above can be proved.

Theorem 6 Suppose that V is a quadratic invariant of the semilinear ODE (4). Assume that V is also an invariant of the linear part $\dot{y} = My$. In addition, we assume (A, b) satisfies (2). Then, the solution $y_1^{(k)}$ of the scheme defined by Definition 2 satisfies $V(y_1^{(k)}) = V(y_0)$.

Moreover, Theorem 3 implies the following theorem showing the accuracy of the scheme defined by Definition 2. Although the proof of Theorem 7 is similar to that for Theorem 3, we present the complete proof here for the reader's convenience.

Theorem 7 Assume the following conditions:

(A1) $Y_i^{(0)}$ satisfies $\|Y_i^{(0)} - y(c_i h)\|_{\mathcal{X}} \leq Ch^q$ for each $i \in [s]$.

(A2) $S: \mathcal{X} \rightarrow \mathcal{L}(\mathcal{X})$ is L_S -Lipschitz continuous.

(A3) The base Runge–Kutta method is of order p .

Then, the numerical solution $y_1^{(k)}$ of the scheme defined by Definition 2 satisfies

$$\|y_1^{(k)} - y(h)\|_{\mathcal{X}} \leq C'h^{\min\{p, q+k-1\}+1}$$

for sufficiently small $h > 0$, i.e., the scheme with k iteration is of order $\min\{p, q+k-1\}$. Here, C' is a constant depending only on the exact solution y , k , M , S , V , A , b , and the constant C in (A1).

Proof Due to the assumption (A1), there exists a function $y^{(0)}: [0, h] \rightarrow \mathcal{X}$ satisfying $y^{(0)}(c_i h) = Y_i^{(0)}$ and $\sup_{t \in [0, h]} \|y^{(0)}(t) - y(t)\|_{\mathcal{X}} \leq Ch^q$. Then, the scheme defined by Definition 2 can be regarded as the usual Runge–Kutta method corresponding to A, b for the system

$$\begin{cases} \dot{w}^{(1)}(t) = \exp(-tM)S(\exp(tM)w^{(0)}(t))(\exp(-tM))^* \nabla V(w^{(1)}(t)), \\ \dot{w}^{(2)}(t) = \exp(-tM)S(\exp(tM)w^{(1)}(t))(\exp(-tM))^* \nabla V(w^{(2)}(t)), \\ \vdots \\ \dot{w}^{(k)}(t) = \exp(-tM)S(\exp(tM)w^{(k-1)}(t))(\exp(-tM))^* \nabla V(w^{(k)}(t)), \end{cases}$$

where $w^{(0)}(t) = \exp(-tM)y^{(0)}(t)$. Otherwise expressed, the scheme can be regarded as the usual Lawson method for the system

$$\begin{cases} \dot{y}^{(1)}(t) = My^{(1)} + S(y^{(0)}(t)) \nabla V(y^{(1)}(t)), \\ \dot{y}^{(2)}(t) = My^{(2)} + S(y^{(1)}(t)) \nabla V(y^{(2)}(t)), \\ \vdots \\ \dot{y}^{(k)}(t) = My^{(k)} + S(y^{(k-1)}(t)) \nabla V(y^{(k)}(t)), \end{cases}$$

where $y^{(i)}(t) = \exp(tM)w^{(i)}(t)$ for $i = 1, \dots, k$. Therefore, it is sufficient to prove $\|y^{(k)}(h) - y(h)\|_{\mathcal{X}} \leq C^{(k)}h^{q+k}$.

To this end, we prove

$$\sup_{t \in [0, h]} \|y^{(j)}(t) - y(t)\|_{\mathcal{X}} \leq C^{(j)}h^{q+j} \quad (j = 1, \dots, k)$$

by induction, where $C^{(j)} := (2L_S C_Y)^j C$ ($C_Y := \sup_{t \in [0, h]} \|\nabla V(y(t))\|_{\mathcal{X}}$). We assume that $\sup_{t \in [0, h]} \|y^{(j-1)}(t) - y(t)\|_{\mathcal{X}} \leq C^{(j-1)}h^{q+j-1}$ holds so that

$$\begin{aligned} C_S^{(j-1)} &:= \sup_{t \in [0, h]} \|S(y^{(j-1)}(t))\|_{\mathcal{X}} \\ &\leq \|S(y_0)\|_{\mathcal{X}} + L_S \left(C^{(j-1)}h^{q+j-1} + \sup_{t \in [0, h]} \|y(t) - y_0\|_{\mathcal{X}} \right) \\ &< \infty, \end{aligned}$$

where $\|S(y_0)\|_{\mathcal{X}}$ denotes the operator norm of $S(y_0)$ with respect to the inner product of \mathcal{X} . Then, we see

$$\begin{aligned} & \sup_{t \in [0, h]} \|y^{(j)}(t) - y(t)\|_{\mathcal{X}} \\ = & \sup_{t \in [0, h]} \left\| \int_0^t \left(M y^{(j)}(r) + S(y^{(j-1)}(r)) \nabla V(y^{(j)}(r)) - M y(r) - S(y(r)) \nabla V(y(r)) \right) dr \right\|_{\mathcal{X}} \\ \leq & h \|M\|_{\mathcal{X}} \sup_{t \in [0, h]} \|y^{(j)}(t) - y(t)\|_{\mathcal{X}} \\ & + \sup_{t \in [0, h]} \left\| \int_0^t \left(S(y^{(j-1)}(r)) \nabla V(y^{(j)}(r)) - S(y^{(j-1)}(r)) \nabla V(y(r)) \right) dr \right\|_{\mathcal{X}} \\ & + \sup_{t \in [0, h]} \left\| \int_0^t \left(S(y^{(j-1)}(r)) \nabla V(y(r)) - S(y(r)) \nabla V(y(r)) \right) dr \right\|_{\mathcal{X}} \\ \leq & h \left(\|M\|_{\mathcal{X}} + C_S^{(j-1)} \|Q\|_{\mathcal{X}} \right) \sup_{t \in [0, h]} \|y^{(j)}(t) - y(t)\|_{\mathcal{X}} + L_S C^{(j-1)} C_y h^{q+j}. \end{aligned}$$

Since we assume that h is sufficiently small, in particular, $h \leq (2\|M\|_{\mathcal{X}} + 2M_S^{(j-1)}\|Q\|_{\mathcal{X}})^{-1}$ holds, we see

$$\sup_{t \in [0, h]} \|y^{(j)}(t) - y(t)\|_{\mathcal{X}} \leq 2L_S C^{(j-1)} C_y h^{q+j}.$$

Thus, by induction, we obtain $\sup_{t \in [0, h]} \|y^{(j)}(t) - y(t)\|_{\mathcal{X}} \leq C^{(j)} h^{q+j}$ ($j = 1, \dots, k$). Since it implies $\|y^{(k)}(h) - y(h)\|_{\mathcal{X}} \leq C^{(k)} h^{q+k}$, the theorem holds. \square

Remark 3 *In the above and subsequent theorems, we assume the global Lipschitz continuity of S for the sake of brevity. However, this can be relaxed to the local Lipschitz continuity.*

Remark 4 *As can be seen from the proof above, even if \mathcal{X} is infinite-dimensional, the same theorem holds if $\|S(y_0)\|_{\mathcal{X}}$, $\|M\|_{\mathcal{X}}$, and $\|Q\|_{\mathcal{X}}$ are bounded.*

4 Application to the SAV approach

The scheme defined in Definition 2 generally requires solving linear equations that involve matrix exponential functions, which is computationally somewhat more expensive than the usual linearly implicit schemes from [16]. However, combining it to the SAV system (9) produces a scheme that is computationally extremely inexpensive. Specifically, the main part of each iteration is the product of a matrix exponential function and a vector $O(s)$ times, which can be computed in parallel. The computational cost of the scheme can be the same level as that of explicit exponential integrators (see Appendix 3). When an efficient implementation of the matrix exponential function is available, the proposed scheme in this section overwhelms the linearly implicit scheme defined in Definition 1 in terms of computational efficiency (see Fig. 2).

4.1 Simple case (9)

The SAV system (9) can be rewritten as

$$\frac{d}{dt} \begin{bmatrix} u \\ r \end{bmatrix} = \begin{bmatrix} J & J\phi(u) \\ -(J\phi(u))^* & \langle \phi(u), J\phi(u) \rangle \end{bmatrix} \begin{bmatrix} Lu \\ 2r \end{bmatrix}.$$

Recall that $J \in \mathcal{L}(\mathcal{V})$, $L \in \mathcal{L}(\mathcal{V})$, $\phi: \mathcal{V} \rightarrow \mathcal{V}$ is a function defined as $\phi(u) = \nabla E(u) / (2\sqrt{E(u) + \alpha})$, and this system has the quadratic invariant $V(u, r) = \frac{1}{2}\langle Lu, u \rangle + r^2 - \alpha$. Since the skew-symmetry of J implies $\langle \phi(u), J\phi(u) \rangle = 0$, it can be rewritten as

$$\frac{d}{dt} \begin{bmatrix} u \\ r \end{bmatrix} = \begin{bmatrix} JL & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} u \\ r \end{bmatrix} + \begin{bmatrix} J\phi(u) \\ -(J\phi(u))^* \end{bmatrix} \nabla V(u, r), \quad (12)$$

which is a special case of (10) with

$$y = \begin{bmatrix} u \\ r \end{bmatrix}, \quad M = \begin{bmatrix} JL & \\ & 0 \end{bmatrix}, \quad S(u) = \begin{bmatrix} J\phi(u) \\ -(J\phi(u))^* \end{bmatrix}, \quad Q = \begin{bmatrix} L \\ 2 \end{bmatrix}.$$

Note that the linear part $\dot{y} = My$ preserves V :

$$\frac{d}{dt} V(y(t)) = \langle My, \dot{y} \rangle_{\mathcal{X}} = \langle Lu, JLu \rangle = 0.$$

Therefore, we can apply the scheme defined by Definition 2 for (12). In this case, the linear equation system (11) reads

$$\begin{bmatrix} U_i^{(k)} \\ R_i^{(k)} \end{bmatrix} = \begin{bmatrix} \exp(c_i h JL) & \\ & 1 \end{bmatrix} \begin{bmatrix} u_0 \\ r_0 \end{bmatrix} \\ + h \sum_{j \in [s]} a_{ij} \begin{bmatrix} \exp((c_i - c_j) h JL) & \\ & 1 \end{bmatrix} \begin{bmatrix} J\phi(U_j^{(k-1)}) \\ -(J\phi(U_j^{(k-1)}))^* \end{bmatrix} \begin{bmatrix} LU_j^{(k)} \\ 2R_j^{(k)} \end{bmatrix},$$

where $U_i^{(k)} \in \mathcal{V}$ and $R_i^{(k)} \in \mathbb{R}$ are the inner stages with respect to u and r , respectively. This linear equation, expressed element by element, is as follows:

$$U_i^{(k)} = \exp(c_i h JL) u_0 + 2h \sum_{j \in [s]} a_{ij} R_j^{(k)} \exp((c_i - c_j) h JL) J\phi(U_j^{(k-1)}), \\ R_i^{(k)} = r_0 - h \sum_{j \in [s]} a_{ij} \langle J\phi(U_j^{(k-1)}), LU_j^{(k)} \rangle.$$

By substituting the first equation into the second equation, we obtain a linear equation

$$R_i^{(k)} = r_0 - h \sum_{\ell \in [s]} a_{i\ell} \left\langle J\phi \left(U_\ell^{(k-1)} \right), L \exp \left(c_\ell h J L \right) u_0 \right\rangle - 2h^2 \sum_{j \in [s]} \left(\sum_{\ell \in [s]} a_{i\ell} a_{\ell j} \left\langle J\phi \left(U_\ell^{(k-1)} \right), L \exp \left((c_\ell - c_j) h J L \right) J\phi \left(U_j^{(k-1)} \right) \right\rangle \right) R_j^{(k)}$$

only on $\{R_i^{(k)}\}_{i=1}^s$. To reduce the number of computations of the matrix exponential function, we use Lemma 4, i.e., the relation $L \exp(tJL) = (\exp(-tJL))^* L$, and obtain

$$R_i^{(k)} = r_0 + h \sum_{\ell \in [s]} a_{i\ell} \left\langle \psi_\ell^{(k-1)}, Lu_0 \right\rangle - 2h^2 \sum_{j \in [s]} \left(\sum_{\ell \in [s]} a_{i\ell} a_{\ell j} \left\langle \psi_\ell^{(k-1)}, L\psi_j^{(k-1)} \right\rangle \right) R_j^{(k)},$$

where $\psi_j^{(k-1)} := \exp(-c_j h J L) J\phi \left(U_j^{(k-1)} \right)$. By introducing $R^{(k)} = [R_1^{(k)} \dots R_s^{(k)}]^T$, the linear equation can be simplified to

$$\left(I_s + 2h^2 A(A \circ \Psi) \right) R^{(k)} = r_0 \mathbf{1}_s - h A v, \tag{13}$$

where $I_s \in \mathbb{R}^{s \times s}$ denotes the identity matrix, $A = (a_{ij}) \in \mathbb{R}^{s \times s}$, \circ denotes the Hadamard product, $\Psi \in \mathbb{R}^{s \times s}$ is defined by $\Psi_{i,j} = \left\langle \psi_i^{(k-1)}, L\psi_j^{(k-1)} \right\rangle$, $\mathbf{1}_s \in \mathbb{R}^s$ denotes all one vector of size s , and $v \in \mathbb{R}^s$ is defined by $v_i = \left\langle \psi_i^{(k-1)}, Lu_0 \right\rangle$.

Another note on the implementation should address the fact that $r_1^{(k)}$ can be computed without $\left\{ U_i^{(k)} \right\}_{i=1}^s$:

$$r_1^{(k)} = r_0 + h \sum_{i \in [s]} b_i \left\langle J\phi \left(U_i^{(k-1)} \right), LU_i^{(k)} \right\rangle = r_0 + \sum_{i,j \in [s]} b_i \omega_{ij} \left(R_j^{(k)} - r_0 \right),$$

where ω_{ij} denotes the (i, j) element of the inverse matrix of A . Consequently, we propose the following scheme.

Definition 3

Step 0 Prepare $U_i^{(0)} \approx u(c_i h)$ and set $k = 1$.

Step 1 Compute $\{\psi_i^{(k-1)}\}_{i=1}^s$ by

$$\psi_i^{(k-1)} = \exp(-c_i h J L) J\phi \left(U_i^{(k-1)} \right). \tag{14}$$

Solve the linear equation (13) to obtain $R^{(k)}$. If some criteria hold, go to Step 2. Otherwise, compute $\{U_i^{(k)}\}_{i=1}^s$ by

$$U_i^{(k)} = \exp(c_i h J L) \left(u_0 + 2h \sum_{j \in [s]} a_{ij} R_j^{(k)} \psi_j^{(k-1)} \right), \quad (15)$$

set $k = k + 1$ and repeat Step 1.

Step 2 Output

$$u_1^{(k)} = \exp(h J L) \left(u_0 + 2h \sum_{j \in [s]} b_j R_j^{(k)} \psi_j^{(k-1)} \right),$$

$$r_1^{(k)} = r_0 + \sum_{i, j \in [s]} b_i \omega_{ij} \left(R_j^{(k)} - r_0 \right).$$

Note that, in many cases, the computation of (14) and (15) is the most computationally expensive part. However, efficient algorithms for computing the matrix exponential and vector products are known. Moreover, since the computation of $\psi_1^{(k-1)}, \psi_2^{(k-1)}, \dots, \psi_s^{(k-1)}$ by using (14) can be done independently, they can be trivially parallelized. Similarly, the computation of $U_i^{(k)}$'s by (15) can be parallelized. See Appendix 3 for more details on the computational cost.

The scheme defined by Definition 3 has a unique solution if and only if the linear equation (13) has a unique solution, which is satisfied for small step size h . The step size restriction depends largely on the nature of J and L (see Remark 5).

Theorem 8 Suppose that the step size h satisfies

$$h \max \left\{ 1, \exp \left(\frac{c_{\max} h}{2} \lambda_{\max}(LJ - JL) \right) \right\} < \frac{1}{\|A\|_2 \|J\| \sqrt{2} \|L\| \sup_i \left\| \phi \left(U_i^{(k-1)} \right) \right\|}, \quad (16)$$

where $\lambda_{\max}(LJ - JL)$ denotes the maximum eigenvalue of $LJ - JL$, $c_{\max} = \max_i c_i$, and $\|A\|_2$ denotes the spectral norm of $A \in \mathbb{R}^{s \times s}$. Then, the linear (13) has a unique solution.

Proof The linear (13) has a unique solution if and only if the matrix $I_s + 2h^2 A$ ($A \circ \Psi$) is nonsingular. Then, $\|2h^2 A (A \circ \Psi)\|_2 < 1$ is a sufficient condition for the solution's unique existence.

Since we have $\|A \circ \Psi\|_2 \leq \|A\|_2 \|\Psi\|_2$ (cf. [31]), we evaluate the norm $\|\Psi\|_2$:

$$\begin{aligned} \|\Psi\|_2 &= \sup_{x \in \mathbb{R}^s} \frac{\sum_{i,j \in [s]} x_i x_j \langle \psi_i^{(k-1)}, L \psi_j^{(k-1)} \rangle}{\|x\|_2^2} \\ &\leq \|L\| \sup_{x \in \mathbb{R}^s} \frac{\sum_{i,j \in [s]} x_i x_j \|\psi_i^{(k-1)}\| \|\psi_j^{(k-1)}\|}{\|x\|_2^2} \\ &= \|L\| \max_i \|\psi_i^{(k-1)}\|^2 \\ &= \|L\| \max_i \left\| \exp(-c_i h J L) J \phi \left(U_i^{(k-1)} \right) \right\|^2 \\ &\leq \|L\| \|J\|^2 \left(\sup_{t \in [0, c_{\max} h]} \left\| \exp(-t J L) \right\| \right)^2 \max_i \left\| \phi \left(U_i^{(k-1)} \right) \right\|^2. \end{aligned}$$

Note that $\|\exp(t(-JL))\| \leq \exp(t\omega(-JL))$ holds for any $t \geq 0$ (cf. [32]), where $\omega(-JL)$ denotes the numerical abscissa of the matrix $-JL$. Since $\omega(-JL) = \lambda_{\max} \left(\frac{1}{2} (-JL + (-JL)^*) \right) = \lambda_{\max} \left(\frac{1}{2} (LJ - JL) \right)$ holds, we have

$$\sup_{t \in [0, c_{\max} h]} \|\exp(-t J L)\|^2 \leq \max \{1, \exp(c_{\max} h \lambda_{\max}(LJ - JL))\},$$

which implies the theorem. □

Remark 5 *The step size restriction (16) may seem severe. Of course, this is true when $\lambda_{\max}(LJ - JL)$ is positive, and the above argument only shows the unique existence for very small step sizes. However, when L and J commute, we have $\lambda_{\max}(LJ - JL) = 0$, and the step size restriction is rather mild. The commutativity of L and J is sometimes assumed in the literature (cf. [33], see also Sect. 5).*

Theorem 6 implies the following theorem showing the discrete conservation law with respect to the modified invariant V .

Theorem 9 (Conservation law) *The solution $(u_1^{(k)}, r_1^{(k)})$ of the scheme defined by Definition 3 satisfies $V(u_1^{(k)}, r_1^{(k)}) = V(u_0, r_0)$.*

Theorem 7 implies the following theorem showing the accuracy of the scheme defined by Definition 3.

Theorem 10 (Accuracy) *Assume the following conditions:*

- (A1) $U_i^{(0)}$ satisfies $\|U_i^{(0)} - u(c_i h)\| \leq Ch^q$ for each $i \in [s]$.
- (A2) $\phi: \mathcal{V} \rightarrow \mathcal{V}$ is Lipschitz continuous.
- (A3) The base Runge–Kutta method is of order p .

Then, the numerical solution $(u_1^{(k)}, r_1^{(k)})$ of the scheme defined by Definition 3 satisfies

$$\begin{aligned} \|u_1^{(k)} - u(h)\| &\leq C'h^{\min\{p, q+k-1\}+1}, \\ |r_1^{(k)} - r(h)| &\leq C'h^{\min\{p, q+k-1\}+1} \end{aligned}$$

for sufficiently small $h > 0$, i.e., the scheme with k iteration is of order $\min\{p, q+k-1\}$. Here, C' is a constant depending only on the exact solution u , k , J , L , ϕ , A , b , and the constant C in (A1).

Proof It is sufficient to prove

$$\left\| \begin{bmatrix} J\phi(u_1) \\ -(J\phi(u_1))^* \end{bmatrix} - \begin{bmatrix} J\phi(u_2) \\ -(J\phi(u_2))^* \end{bmatrix} \right\|_{\mathcal{X}} = \|J\phi(u_1) - J\phi(u_2)\|. \quad (17)$$

In general, for any $v \in \mathcal{V}$, we see

$$\left\| \begin{bmatrix} v \\ -(v)^* \end{bmatrix} \right\|_{\mathcal{X}} = \sup_{w \in \mathcal{V}, \rho \in \mathbb{R}} \frac{\sqrt{\rho^2 \|v\|^2 + ((v, w))^2}}{\sqrt{\|w\|^2 + \rho^2}} = \|v\| \sup_{w \in \mathcal{V}, \rho \in \mathbb{R}} \frac{\sqrt{\rho^2 + \|w\|^2}}{\sqrt{\|w\|^2 + \rho^2}} = \|v\|,$$

which proves (17). \square

Remark 6 The scheme defined by Definition 3 includes the schemes proposed in [27] as special cases. In [27], they focus on the cases $U_i^{(0)}$ computed by the extrapolation or $U_i^{(0)} = u_0$. In addition, the intended implementation is also somewhat different: for example, their algorithm requires $ks(s-1)$ computations of the product of the matrix exponential functions and vectors, while it is reduced to $(2k-1)s$ in Definition 3 by introducing $\psi^{(k-1)}$.

Moreover, the consequences of Theorem 10 are consistent with the accuracy confirmed numerically in [27]. For example, from numerical experiments in [27, Remark 3.6], the authors predict that the scheme has fourth-order accuracy when the base RK method is a fourth-order Gauss method, $U_i^{(0)} = u_0$ (i.e., $q = 1$), and $k = 4$.

4.2 Multiple scalar auxiliary variables

As shown in [15], when E is unbounded, two scalar auxiliary variables are needed. The scheme described in the previous section can be extended to this case.

Let us consider the Hamiltonian system (7) with

$$H(u) = \frac{1}{2} \langle Lu, u \rangle + E_L(u) - E_U(u),$$

where $E_X: \mathcal{V} \rightarrow \mathbb{R}$ are bounded from below, i.e., $\alpha_X := -\inf E_X(u) + \epsilon_X$ is finite ($X \in \{L, U\}$). Then, by introducing $r_X := \sqrt{E_X(u) + \alpha_X}$, $\phi_X(u) := \nabla E_X(u) / (2\sqrt{E_X(u) + \alpha_X})$, and $V(u, r_L, r_U) = \frac{1}{2}\langle Lu, u \rangle + r_L^2 - r_U^2$, we obtain

$$\frac{d}{dt} \begin{bmatrix} u \\ r_L \\ r_U \end{bmatrix} = \begin{bmatrix} JL & & \\ & 0 & \\ & & 0 \end{bmatrix} \begin{bmatrix} u \\ r_L \\ r_U \end{bmatrix} + \begin{bmatrix} J\phi_L & J\phi_U \\ -(J\phi_L)^* & 0 & \langle \phi_L, J\phi_U \rangle \\ -(J\phi_U)^* & -\langle \phi_L, J\phi_U \rangle & 0 \end{bmatrix} \nabla V(u, r_L, r_U), \tag{18}$$

which is a special case of (10) ($\phi_X(u)$ is abbreviated to ϕ_X).

Therefore, we can apply the scheme defined by Definition 2. In addition, by introducing $\psi_i^X := \exp(-c_i h JL) J\phi_X(U_i^{(k-1)})$, the techniques to reduce the computational cost in Sect. 4.1 can also be used in this case (see Appendix 1 for details). Then, the linear equation system with respect to $R_L^{(k)} = [R_{L,1}^{(k)} \dots R_{L,s}^{(k)}]^T$ and $R_U^{(k)} = [R_{U,1}^{(k)} \dots R_{U,s}^{(k)}]^T$ can be written as

$$\begin{bmatrix} I_s + 2h^2 A (A \circ \Psi^{LL}) & 2hA\Phi - 2h^2 A (A \circ \Psi^{LU}) \\ 2hA\Phi + 2h^2 A (A \circ \Psi^{UL}) & I_s - 2h^2 A (A \circ \Psi^{UU}) \end{bmatrix} \begin{bmatrix} R_L^{(k)} \\ R_U^{(k)} \end{bmatrix} = \begin{bmatrix} r_{L,0} \mathbf{1}_s - hA v^L \\ r_{U,0} \mathbf{1}_s - hA v^U \end{bmatrix} \tag{19}$$

where $\Psi^{LL}, \Psi^{LU}, \Psi^{UL}, \Psi^{UU} \in \mathbb{R}^{s \times s}$ are defined by $\Psi_{i,j}^{XY} = \langle \psi_i^X, L\psi_j^Y \rangle$ for $X, Y \in \{L, U\}$ (note that $\Psi^{LU} = (\Psi^{UL})^T$ holds), $\Phi \in \mathbb{R}^{s \times s}$ is a diagonal matrix defined by $\Phi_{ii} = \langle \phi_L(U_i^{(k-1)}), J\phi_U(U_i^{(k-1)}) \rangle$, and $v^L, v^U \in \mathbb{R}^s$ are defined by $v_i^X = \langle \psi_i^X, Lu_0 \rangle$. Consequently, we obtain the following scheme:

Definition 4

Step 0 Prepare $U_i^{(0)} \approx u(c_i h)$ and set $k = 1$.

Step 1 Compute $\{\psi_i^X\}_{i=1}^s$ by

$$\psi_i^X = \exp(-c_i h JL) J\phi_X(U_i^{(k-1)}).$$

Solve the linear equation system (19) to obtain $R_L^{(k)}$ and $R_U^{(k)}$. If some criteria hold, go to Step 2. Otherwise, compute $\{U_i^{(k)}\}_{i=1}^s$ by

$$U_i^{(k)} = \exp(c_i h JL) \left(u_0 + 2h \sum_{j \in [s]} a_{ij} (R_{L,j}^{(k)} \psi_j^L - R_{U,j}^{(k)} \psi_j^U) \right),$$

set $k = k + 1$ and repeat Step 1.

Step 2 Output

$$\begin{aligned} u_1^{(k)} &= \exp(hJL) \left(u_0 + 2h \sum_{j \in [s]} b_j \left(R_{L,j}^{(k)} \psi_j^L - R_{U,j}^{(k)} \psi_j^U \right) \right), \\ r_{X,1}^{(k)} &= r_{X,0} + \sum_{i,j \in [s]} b_i \omega_{ij} \left(R_{X,j}^{(k)} - r_0 \right). \end{aligned}$$

The scheme defined by Definition 4 has properties similar to the scheme defined by Definition 3. Here, we only provide the results, and the proof is given in Appendix 2. In particular, for the uniqueness and existence, only the case $LJ = JL$ is presented here, while the general case is presented in Theorem 14.

Theorem 11 Suppose that $LJ = JL$ holds. If h satisfies

$$h < \frac{\sqrt{1 + \frac{4\|L\|}{C_\phi}} - 1}{4\|A\|_2\|J\|\|L\|}, \quad C_\phi := \max_{i \in [s], X \in \{L, U\}} \left\| \phi_X \left(U_i^{(k-1)} \right) \right\|^2,$$

the linear equation (19) has a unique solution.

Theorem 12 (Conservation law) The solution $(u_1^{(k)}, r_{L,1}^{(k)}, r_{U,1}^{(k)})$ of the scheme defined by Definition 4 satisfies $V(u_1^{(k)}, r_{L,1}^{(k)}, r_{U,1}^{(k)}) = V(u_0, r_{L,0}, r_{U,0})$.

Theorem 13 (Accuracy) Assume the following conditions:

- (A1) $U_i^{(0)}$ satisfies $\|U_i^{(0)} - u(c_i h)\| \leq Ch^q$ for each $i \in [s]$.
- (A2) $\phi_X: \mathcal{V} \rightarrow \mathcal{V}$ is Lipschitz continuous and bounded.
- (A3) The base Runge–Kutta method is of order p .

Then, the numerical solution $(u_1^{(k)}, r_{L,1}^{(k)}, r_{U,1}^{(k)})$ of the scheme defined by Definition 4 satisfies

$$\begin{aligned} \left\| u_1^{(k)} - u(h) \right\| &\leq C' h^{\min\{p, q+k-1\}+1}, \\ \left| r_{L,1}^{(k)} - r_L(h) \right| &\leq C' h^{\min\{p, q+k-1\}+1}, \\ \left| r_{U,1}^{(k)} - r_U(h) \right| &\leq C' h^{\min\{p, q+k-1\}+1} \end{aligned}$$

for a sufficiently small $h > 0$, i.e., the scheme with k iteration is of order $\min\{p, q+k-1\}$. Here, C' is a constant depending only on the exact solution u , k , J , L , ϕ_X , A , b , and the constant C in (A1).

5 Numerical experiments

We employ the sixth-order Gauss method as the base RK method. The initial approximation $U_i^{(0)} \approx u(c_i h)$ is computed by the Lawson transformation and the continuous explicit RK (CERK) method with order 1, 2, . . . , 5 ($q = 2, 3, \dots, 6$). All numerical experiments are performed in Julia (Version 1.8.0) on a PC with Apple M1 Ultra and 128GB RAM. The numerical experiments in this paper are not parallelized. We leave it to future work.

5.1 Modified Korteweg–de Vries equation

Let us consider the modified Korteweg–de Vries (mKdV) equation (on $\mathbb{S} := \mathbb{R}/L\mathbb{Z}$):

$$u_t = -\partial_x \frac{\delta \mathcal{H}}{\delta u}, \quad \mathcal{H}(u) = \int \left(-\frac{1}{2}(u_x)^2 + \frac{1}{2}u^4 \right) dx.$$

We employ an energy-preserving spatial discretization

$$\dot{u}_k = -\delta_x \left(\delta_x^2 u_k + 2(u_k)^3 \right), \quad H(u) = -\frac{1}{2} \sum_{k=1}^N (\delta_x u_k)^2 \Delta x + \frac{1}{2} \sum_{k=1}^N (u_k)^4 \Delta x, \quad (20)$$

where δ_x denotes the Fourier-spectral difference operator (see, e.g., [34] for details on the difference operator). The inner product is defined as $\langle v, w \rangle = \sum_{k=1}^N v_k w_k \Delta x$. Then, $\nabla H(u) = \delta_x^2 u_k + 2(u_k)^3$ holds.

Then, $E(u) := \frac{1}{2} \sum_{k=1}^N (u_k)^4 \Delta x$ is bounded from below so that we can use the scheme defined by Definition 3. Note that the product of the matrix exponential function $e^{t \delta_x^3}$ and a vector can be computed by the FFT. In numerical experiments, we use the exact solution $\text{dn}(x - (2 - m)t \mid m)$ with the spatial period $L = 2K(m)$ and temporal period $T = L/|2 - m|$, where dn is one of the Jacobi elliptic functions and K denotes the complete elliptic integral of the first kind. We choose the parameters $m = 0.1$ ($L \approx 3.225, T \approx 1.697$) and $N = 16$.

Figure 1 summarizes the relative errors. They decrease along the reference lines drawn based on Theorem 10. Figure 2 compares the computational cost of the proposed scheme with the scheme defined by [16] (see Appendix 4 for details) and an explicit Lawson method (corresponds to a seven stage sixth-order explicit Runge–Kutta method [35]). Note that all schemes have sixth order. The proposed scheme is more efficient than the scheme defined by [16]. However, the proposed scheme is less efficient than the explicit Lawson method for such a short time interval. As shown in the next example, the proposed scheme is more efficient than the explicit Lawson method for a long time interval (this advantage is typical in the comparison between conservative and non-conservative schemes).

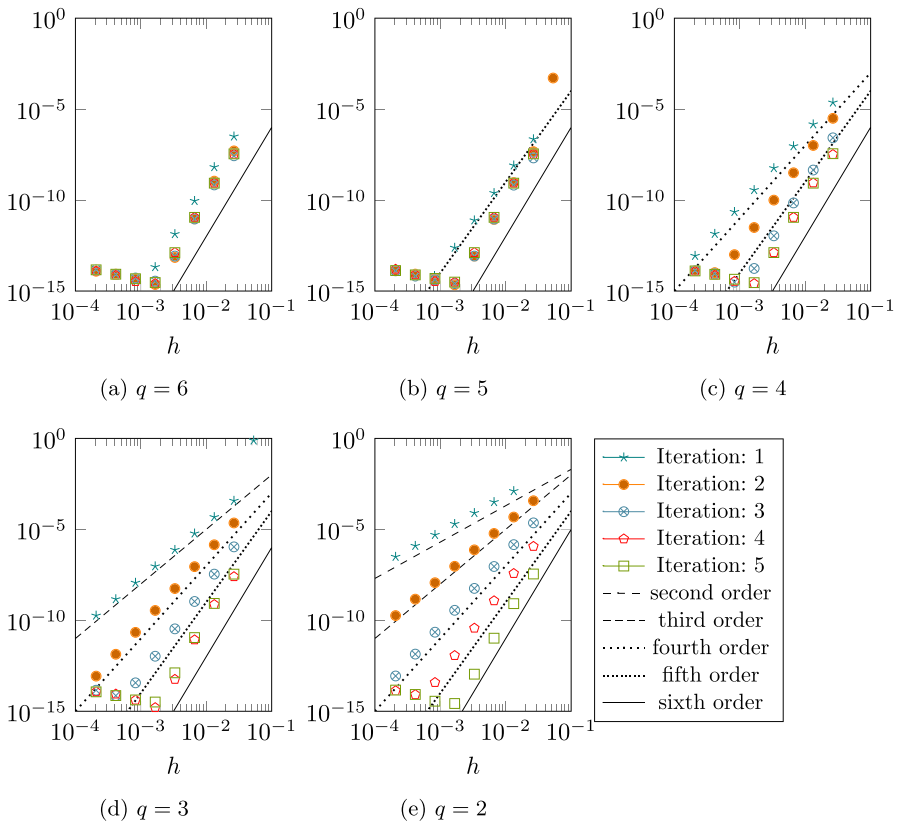


Fig. 1 Relative errors of the numerical solution for the mKdV equation

5.2 Korteweg–de Vries equation

Let us consider the Korteweg–de Vries (KdV) equation (on $\mathbb{S} := \mathbb{R}/L\mathbb{Z}$):

$$u_t = \partial_x \frac{\delta \mathcal{H}}{\delta u}, \quad \mathcal{H}(u) = \int \left(\frac{1}{2} (u_x)^2 - u^3 \right) dx.$$

We employ an energy-preserving spatial discretization

$$\dot{u}_k = \delta_x \left(-\delta_x^2 u_k - 3(u_k)^2 \right), \quad H(u) = \frac{1}{2} \sum_{k=1}^N (\delta_x u_k)^2 \Delta x + \sum_{k=1}^N (u_k)^3 \Delta x. \quad (21)$$

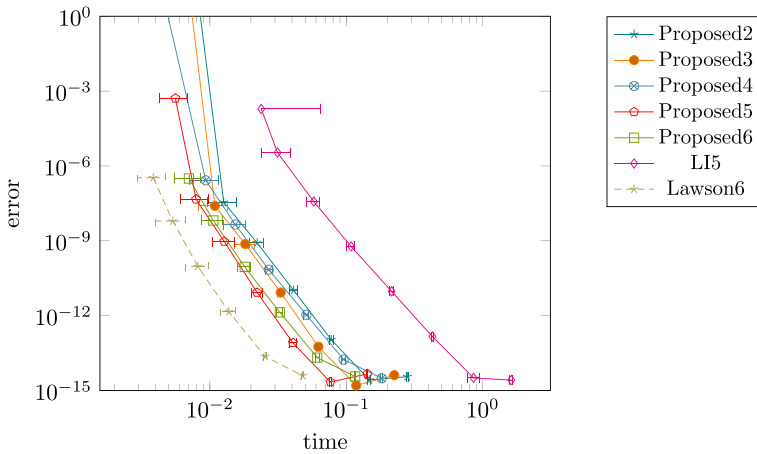


Fig. 2 Comparison of the computational efficiency of several schemes for the mKdV equation. In the legend, “proposed*i*” denotes the sixth-order scheme defined by Definition 3 with $q = i$ and $k = 6 - i$, “LI5” denotes the sixth-order scheme defined by [16] (1) with $q = 5$ and $k = 2$, and “Lawson6” denotes a seven stage sixth-order explicit Lawson method. Note that, the proposed schemes and the scheme defined by [16] are applied to the SAV systems (12) and (9), respectively, and the explicit Lawson method is directly applied to the system (20). Each scheme is computed 30 times, and the mean values and the standard deviations of the computational times are shown

Since $\sum_{k=1}^N (u_k)^3 \Delta x$ is unbounded, we employ the scheme defined by Definition 4. We define

$$E_L(u) := \sum_{k=1}^N \left((u_k)^4 - (u_k)^3 \right) \Delta x, \quad E_U(u) := \sum_{k=1}^N (u_k)^4 \Delta x$$

similarly to [15]. In the numerical experiments below, we use the exact solution $2m \operatorname{cn}(x - ct \mid m)^2$, where $c = 4(2m - 1)$ and cn is a Jacobi elliptic function. We choose the parameters $m = 0.1$ ($L = 2K(m) \approx 3.225$, $T = L/|c| \approx 1.008$) and $N = 64$.

Figure 3 summarizes the relative errors. They decrease along the reference lines.

To confirm the long-term behavior of the proposed scheme, we conducted the numerical experiments over 32 periods. Figure 4 shows the evolution of relative errors of the invariant H , modified invariant V , and $I(u) := \frac{1}{2} \sum (u_k)^2 \Delta x$ corresponding to the another invariant $\mathcal{I}(u) := \int (\frac{1}{2}u^2) dx$ of the KdV equation.

As for the modified invariant V , as expected, it is very well preserved. The relative errors for the original invariants H and I are also small.

There, the proposed scheme ($q = 6$ and $k = 1$) with $h = T/64 \approx 0.01575$ is applied to the SAV system (12), and the explicit Lawson method with $h = T/512 \approx 0.001968$ is directly applied to the system (21). We employ the small step size for the Lawson method because numerical solutions computed with $h = T/256$ diverged. Even in the step size employed in the figure, the behavior of the relative errors of H suggests that it is expected to diverge after a slightly long run. In the current setup,

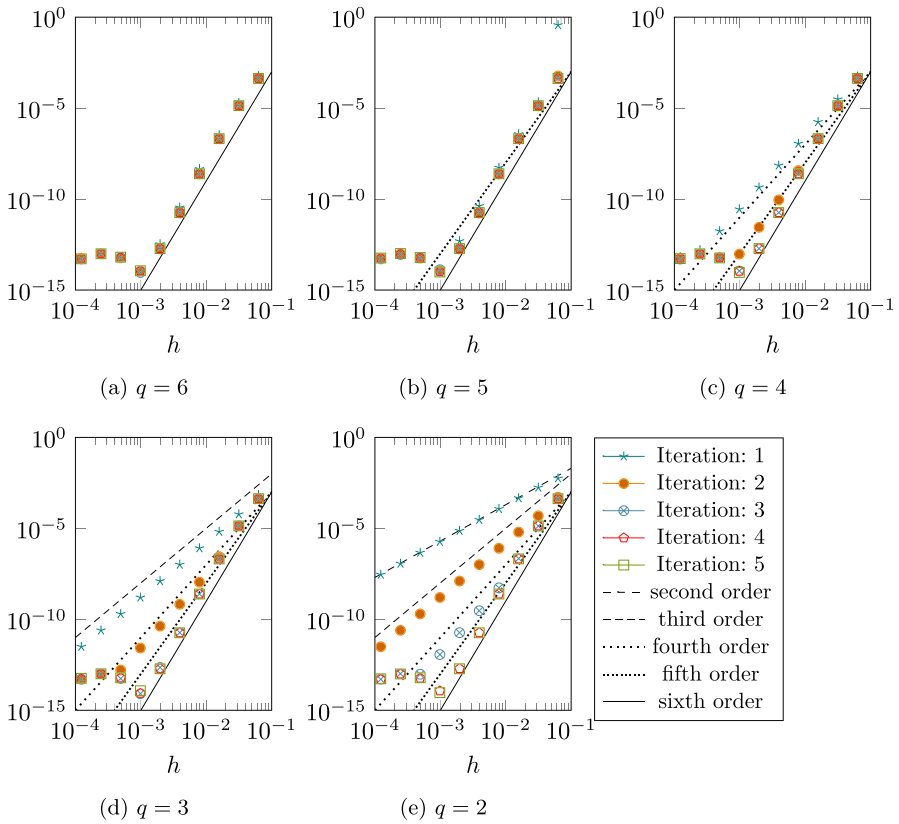


Fig. 3 Relative errors of the numerical solution for the KdV equation

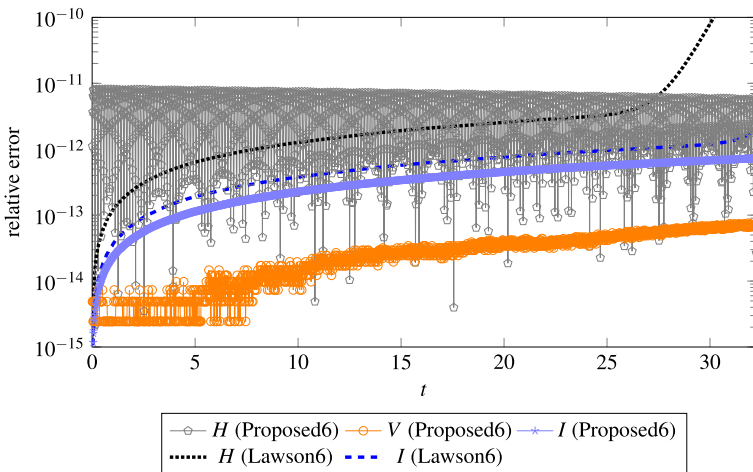


Fig. 4 Evolution of the relative errors of invariants of the KdV equation. In the figure, H , V , and I denote the invariant, the modified invariant, and the other invariant, respectively. The proposed scheme employed the step size $h = T/64 \approx 0.01575$, and the Lawson method employed the step size $h = T/512 \approx 0.001968$

the relative errors of numerical solutions itself at the final time are approximately 5.39×10^{-7} for the proposed scheme and 8.33×10^{-7} for the Lawson method. The computational time of the proposed scheme is 0.678 s, and that of the Lawson method is 0.882 s (they are mean values of 30 computations, with a standard deviation of about 0.01 for both methods).

5.3 Sine-Gordon equation

Let us consider the sine-Gordon (sG) equation (on $\mathbb{S} := \mathbb{R}/L\mathbb{Z}$):

$$u_{tt} - u_{xx} + \sin u = 0.$$

By introducing a new variable $v := u_t$, we obtain the following first-order system:

$$\frac{\partial}{\partial t} \begin{bmatrix} u \\ v \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix} \begin{bmatrix} \frac{\delta \mathcal{H}}{\delta u} \\ \frac{\delta \mathcal{H}}{\delta v} \end{bmatrix}, \quad \mathcal{H}(u, v) = \int \left(\frac{1}{2} (u_x)^2 + \frac{1}{2} v^2 - \cos u \right) dx.$$

We employ an energy-preserving spatial discretization

$$\begin{aligned} \frac{d}{dt} \begin{bmatrix} u_k \\ v_k \end{bmatrix} &= \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix} \begin{bmatrix} v_k \\ -\delta_x^2 u_k + \sin u_k \end{bmatrix}, \\ H(u, v) &= \sum_{k=1}^N \left(\frac{1}{2} (\delta_x u_k)^2 + \frac{1}{2} (v_k)^2 - \cos u_k \right) \Delta x. \end{aligned}$$

Since $E(u) = \sum_{k=1}^N \cos u_k \Delta x$ is bounded, we employ the scheme defined by Definition 3.

In this case, since we employ the Fourier spectral difference δ_x which can be diagonalized by the Fourier transform matrix F , i.e., $\delta_x = F^* (i\Lambda) F$ ($\Lambda \in \mathbb{R}^{N \times N}$ is a diagonal matrix), we have

$$\begin{aligned} (JL)^{2m} &= \begin{bmatrix} F & \\ & F \end{bmatrix}^* \begin{bmatrix} (-1)^m \Lambda^{2m} & \\ & (-1)^m \Lambda^{2m} \end{bmatrix} \begin{bmatrix} F & \\ & F \end{bmatrix}, \\ (JL)^{2m+1} &= \begin{bmatrix} F & \\ & F \end{bmatrix}^* \begin{bmatrix} & (-1)^m \Lambda^{2m} \\ (-1)^{m+1} \Lambda^{2m+2} & \end{bmatrix} \begin{bmatrix} F & \\ & F \end{bmatrix} \end{aligned}$$

for all $m = 0, 1, \dots$. Therefore, we have

$$\exp(tJL) = \begin{bmatrix} F & \\ & F \end{bmatrix}^* \begin{bmatrix} \cos(t\Lambda) & f_{sG}(t, \Lambda) \\ -\Lambda \sin(t\Lambda) & \cos(t\Lambda) \end{bmatrix} \begin{bmatrix} F & \\ & F \end{bmatrix},$$

where $f_{sG}(t, \Lambda)$ is a diagonal matrix defined by $(f_{sG}(t, \Lambda))_{ii} = \Lambda_{ii}^{-1} \sin(t\Lambda_{ii})$ if $\Lambda_{ii} \neq 0$ and $(f_{sG}(t, \Lambda))_{ii} = t$ otherwise.

In numerical experiments, we use the exact solution [36]

$$u(t, x) = 4 \arctan \left(\gamma^2 \operatorname{cn}(\beta_1 x | k_1) \operatorname{cn}(\beta_2 t | k_2) \right),$$

where $k_1^2 = \frac{\gamma^2}{1+\gamma^2} \left(1 + \frac{1}{\beta_1(1+\gamma^2)} \right)$, $k_2^2 = \frac{\gamma^2}{1+\gamma^2} \left(1 - \frac{1}{\beta_2(1+\gamma^2)} \right)$, $\beta_2^2 = \beta_1^2 + \frac{1-\gamma^2}{1+\gamma^2}$, β_1 and γ are parameters (we choose $\gamma = 0.1$ and $\beta_1 = 1$ in the numerical experiments below). This solution has the spatial period $L = 4K(k_1^2)$ and the temporal period $T = 4K(k_2^2)$. In the numerical experiment, we choose $N = 16$.

Figure 5 summarizes the relative errors. The errors achieved higher orders of convergence than expected by Theorem 10. From the numerical results, the order of accuracy seems to be $\min\{p, q + 2k - 1\}$. In fact, this can be proved by using the specific form of the sine-Gordon equation (see Appendix 5). There, we show that the order of accuracy with respect to u is $\min\{p, q + 2k - 1\}$, while that with respect to v and r is $\min\{p, q + 2k - 2\}$.

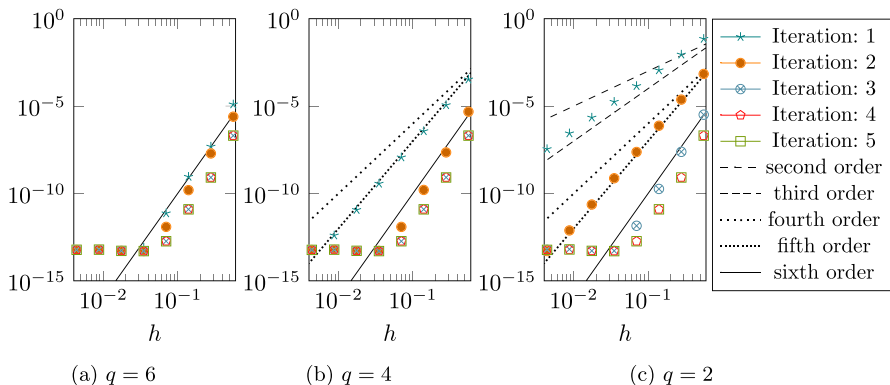


Fig. 5 Relative errors of the numerical solution with respect to u for the sG equation

Appendix 1. Derivation of the scheme defined by Definition 4

When we apply the scheme defined by (2) to the equation (18), the linear equation system (11) reads

$$\begin{aligned}
 U_i^{(k)} &= \exp(c_i h J L) u_0 \\
 &\quad + 2h \sum_{j \in [s]} a_{ij} \exp((c_i - c_j) h J L) \left(R_{L,j}^{(k)} J \phi_L \left(U_j^{(k-1)} \right) - R_{U,j}^{(k)} J \phi_U \left(U_j^{(k-1)} \right) \right), \\
 R_{L,i}^{(k)} &= r_{L,0} - h \sum_{j \in [s]} a_{ij} \left\langle J \phi_L \left(U_j^{(k-1)} \right), L U_j^{(k)} \right\rangle \\
 &\quad - 2h \sum_{j \in [s]} a_{ij} R_{U,j}^{(k)} \left\langle \phi_L \left(U_j^{(k-1)} \right), J \phi_U \left(U_j^{(k-1)} \right) \right\rangle, \\
 R_{U,i}^{(k)} &= r_{U,0} - h \sum_{j \in [s]} a_{ij} \left\langle J \phi_U \left(U_j^{(k-1)} \right), L U_j^{(k)} \right\rangle \\
 &\quad - 2h \sum_{j \in [s]} a_{ij} R_{L,j}^{(k)} \left\langle \phi_L \left(U_j^{(k-1)} \right), J \phi_U \left(U_j^{(k-1)} \right) \right\rangle.
 \end{aligned}$$

By introducing $\psi_i^X = \exp(-c_i h J L) J \phi_X \left(U_i^{(k-1)} \right)$, we see

$$\begin{aligned}
 \left\langle J \phi_X \left(U_i^{(k-1)} \right), L U_i^{(k)} \right\rangle &= \left\langle \psi_i^X, \left(\exp(c_i h J L) \right)^* L U_i^{(k)} \right\rangle \\
 &= \left\langle \psi_i^X, L \exp(-c_i h J L) U_i^{(k)} \right\rangle \\
 &= \left\langle \psi_i^X, L \left(u_0 + 2h \sum_{j \in [s]} a_{ij} \left(R_{L,j}^{(k)} \psi_j^L - R_{U,j}^{(k)} \psi_j^U \right) \right) \right\rangle \\
 &= \left\langle \psi_i^X, L u_0 \right\rangle + 2h \sum_{j \in [s]} a_{ij} \left(R_{L,j}^{(k)} \left\langle \psi_i^X, L \psi_j^L \right\rangle - R_{U,j}^{(k)} \left\langle \psi_i^X, L \psi_j^U \right\rangle \right).
 \end{aligned}$$

Therefore, we obtain the linear equation system

$$\begin{aligned}
 R_{L,i}^{(k)} &= r_{L,0} - h \sum_{\ell \in [s]} a_{i\ell} \langle \psi_\ell^L, L u_0 \rangle - 2h^2 \sum_{j \in [s]} \left(\sum_{\ell \in [s]} a_{i\ell} a_{\ell j} \langle \psi_\ell^L, L \psi_j^L \rangle \right) R_{L,j}^{(k)} \\
 &\quad - 2h \sum_{j \in [s]} \left(a_{ij} \langle \phi_L \left(U_j^{(k-1)} \right), J \phi_U \left(U_j^{(k-1)} \right) \rangle - h \sum_{\ell \in [s]} a_{i\ell} a_{\ell j} \langle \psi_\ell^L, L \psi_j^U \rangle \right) R_{U,j}^{(k)}, \\
 R_{U,i}^{(k)} &= r_{U,0} - h \sum_{\ell \in [s]} a_{i\ell} \langle \psi_\ell^U, L u_0 \rangle + 2h^2 \sum_{j \in [s]} \left(\sum_{\ell \in [s]} a_{i\ell} a_{\ell j} \langle \psi_\ell^U, L \psi_j^U \rangle \right) R_{U,j}^{(k)} \\
 &\quad - 2h \sum_{j \in [s]} \left(a_{ij} \langle \phi_L \left(U_j^{(k-1)} \right), J \phi_U \left(U_j^{(k-1)} \right) \rangle + h \sum_{\ell \in [s]} a_{i\ell} a_{\ell j} \langle \psi_\ell^U, L \psi_j^L \rangle \right) R_{L,j}^{(k)}.
 \end{aligned}$$

Appendix 2. Proofs of theorems in 4.2

Theorem 11 is a corollary of the following theorem.

Theorem 14 *Suppose that the step size h satisfies*

$$h + 2h^2 \|A\|_2 \|J\| \|L\| \max \left\{ 1, \exp \left(\frac{c_{\max} h}{2} \lambda_{\max}(LJ - JL) \right) \right\} < \frac{1}{2 \|A\|_2 \|J\| C_\phi}.$$

Then, the linear equation (19) has a unique solution.

Proof Since the coefficient matrix of the linear equation (19) can be written as

$$I_{2s} + 2h \begin{bmatrix} & A\Phi \\ A\Phi & \end{bmatrix} + 2h^2 \begin{bmatrix} A(A \circ \Psi^{LL}) & -A(A \circ \Psi^{LU}) \\ A(A \circ \Psi^{UL}) & -A(A \circ \Psi^{UU}) \end{bmatrix},$$

it is sufficient to evaluate the spectral norms of the second and the third terms. The second term can be evaluated as

$$\left\| \begin{bmatrix} & A\Phi \\ A\Phi & \end{bmatrix} \right\|_2 = \|A\Phi\|_2 \leq \|A\|_2 \max_i \left| \langle \phi_L(U_i^{(k-1)}), J\phi_U(U_i^{(k-1)}) \rangle \right| \leq \|A\|_2 \|J\| C_\phi.$$

The third term can be evaluated as

$$\begin{aligned} \left\| \begin{bmatrix} A(A \circ \Psi^{LL}) & -A(A \circ \Psi^{LU}) \\ A(A \circ \Psi^{UL}) & -A(A \circ \Psi^{UU}) \end{bmatrix} \right\| &= \left\| \begin{bmatrix} A & \\ & A \end{bmatrix} \left(\begin{bmatrix} A & -A \\ A & -A \end{bmatrix} \circ \begin{bmatrix} \Psi^{LL} & \Psi^{LU} \\ \Psi^{UL} & \Psi^{UU} \end{bmatrix} \right) \right\|_2 \\ &\leq \|A\|_2 (2\|A\|_2) \left\| \begin{bmatrix} \Psi^{LL} & \Psi^{LU} \\ \Psi^{UL} & \Psi^{UU} \end{bmatrix} \right\|_2, \end{aligned}$$

and

$$\left\| \begin{bmatrix} \Psi^{LL} & \Psi^{LU} \\ \Psi^{UL} & \Psi^{UU} \end{bmatrix} \right\|_2 \leq \|L\| \|J\|^2 C_\phi \max \left\{ 1, \exp \left(\frac{c_{\max} h}{2} \lambda_{\max}(LJ - JL) \right) \right\}$$

holds, which can be shown in a manner similar to the proof of Theorem 8. \square

Next, we prove Theorem 13.

Proof of Theorem 13 It is sufficient to prove that the skew-symmetric matrix function

$$S_{\text{MSAV}}(u) := \begin{bmatrix} & J\phi_L(u) & J\phi_U(u) \\ -(J\phi_L(u))^* & 0 & \langle \phi_L(u), J\phi_U(u) \rangle \\ -(J\phi_U(u))^* & -\langle \phi_L(u), J\phi_U(u) \rangle & 0 \end{bmatrix}$$

is Lipschitz continuous.

For any $v_1, v_2 \in \mathcal{V}$ and $\rho \in \mathbb{R}$, we have

$$\begin{aligned} \left\| \begin{bmatrix} v_1 & v_2 \\ -v_1^* & \rho \\ -v_2^* & -\rho \end{bmatrix} \right\| &\leq \left\| \begin{bmatrix} v_1 & v_2 \\ -v_1^* \\ -v_2^* \end{bmatrix} \right\| + \left\| \begin{bmatrix} O & \\ & \rho \end{bmatrix} \right\| \\ &\leq \sqrt{\|v_1\|^2 + \|v_2\|^2} + |\rho| \\ &\leq \|v_1\| + \|v_2\| + |\rho|. \end{aligned}$$

Therefore, we see

$$\begin{aligned} \|S_{\text{MSAV}}(u_1) - S_{\text{MSAV}}(u_2)\| &\leq \|J\phi_L(u_1) - J\phi_L(u_2)\| + \|J\phi_U(u_1) - J\phi_U(u_2)\| \\ &\quad + |\langle \phi_L(u_1), J\phi_U(u_1) \rangle - \langle \phi_L(u_2), J\phi_U(u_2) \rangle| \\ &\leq \|J\|\|\phi_L(u_1) - \phi_L(u_2)\| + \|J\|\|\phi_U(u_1) - \phi_U(u_2)\| \\ &\quad + \frac{1}{2}|\langle \phi_L(u_1) - \phi_L(u_2), J(\phi_U(u_1) + \phi_U(u_2)) \rangle| \\ &\quad + \frac{1}{2}|\langle \phi_L(u_1) + \phi_L(u_2), J(\phi_U(u_1) - \phi_U(u_2)) \rangle| \\ &\leq \|J\| \left(1 + \frac{\|\phi_U(u_1) + \phi_U(u_2)\|}{2} \right) \|\phi_L(u_1) - \phi_L(u_2)\| \\ &\quad + \|J\| \left(1 + \frac{\|\phi_L(u_1) + \phi_L(u_2)\|}{2} \right) \|\phi_U(u_1) - \phi_U(u_2)\|, \end{aligned}$$

which shows the Lipschitz continuity of S_{MSAV} under the assumption (A2). □

Appendix 3. Computational efficiency of the scheme defined by Definition 3

In this section, we evaluate the computational efficiency of the scheme defined by Definition 3 in terms of the number of the computation of the product of matrix exponential and vector. To this end, we consider the explicit Lawson method and the four types of the initial approximation. Before evaluating the computational efficiency, we summarize the number of stages of the explicit Runge–Kutta and continuous explicit Runge–Kutta (CERK) methods achieving the order p in Table 1.

Let us list the four types of initial approximation. To achieve the order p , we consider the Gauss method with $s = \lceil p/2 \rceil$ stages as the base Runge–Kutta method. Then, to minimize the number of iterations in Definition 3, we choose $k = p - q + 1$. Recall that Steps 1 and 2 in Definition 3 require the computation of the product of matrix exponential and vector $(2k - 1)s + 1$ times in sequential computation and $2k$ times in parallel computation. Below, N_S and N_P denote the number of the computation of the product of matrix exponential and vector in sequential and parallel computation, respectively.

Table 1 The number of stages of the explicit RK and CERK methods achieving the order p

Methods\order	1	2	3	4	5	6	7	8
Explicit RK	1	2	3	4	6	7	9	11
Continuous explicit RK	1	2	4	6	8	11	15	–

The written number is the minimum number of stages achieving the order p (cf. [37–39]) except for seventh-order CERK, where we refer to the 15 stage method by Verner [39]. (Eighth-order CERK is omitted since we do not use it later.)

- Proposed (I): the simplest initial approximation $U_i^{(0)} = u_0$. In this case, since $q = 1$ holds, we choose $k = p$. Then, $N_S = (2p - 1)s + 1$ and $N_P = 2p$ hold.
- Proposed (II): an accurate initial approximation using continuous explicit Runge–Kutta methods with the Lawson transformation. In this case, since $q = p$ holds, we choose $k = 1$. Then, $N_S = 2(s^* - 1) + 2s + 1$ and $N_P = 2(s^* - 1) + 3$ hold, where s^* is the number of stages of the continuous explicit Runge–Kutta method achieving the order $p - 1$ (see Table 1).
- Proposed (III): initial approximation using extrapolation of inner stages in the previous step ($p \geq 2$). The inner stages of the Gauss method are $(s + 1)$ th order approximations, the order of approximation is at most $(s + 1)$ th order. Therefore, let us construct the extrapolation satisfying $q = s + 1$ and $k = p - s$. However, since $U_k^{(i)}$ is not computed at the last iteration in Definition 3, the situation is a bit complicated. First, when $p = 2$, we have $s = 1$ and $k = 1$. In this case, to construct a second-order approximation, we employ the extrapolation using the input u_0 and u_1 in the previous step. Therefore, $N_S = N_P = 2$ hold. Next, when $p = 3$, we have $s = 2$ and $k = 1$. In this case, to construct a third-order approximation, we employ the extrapolation using the input u_0 , u_1 , and one of the inner stages in the previous step. Therefore, $N_S = 4$ and $N_P = 3$ hold. Finally, when $p \geq 4$, we have $k \geq 2$. In this case, we can use $U_i^{(k-1)}$'s which are computed in the previous step. Therefore, $N_S = (2p - 2s - 1)s + 1$ and $N_P = 2(p - s)$ hold.
- Proposed (IV): initial approximation using extrapolation of several previous steps. In this case, by using sufficiently many previous steps, we can construct the extrapolation satisfying $q = p$ so that $k = 1$. Then, we have $N_S = s + 1$ and $N_P = 2$ hold.

We compare the computational efficiency of the proposed methods with that of the explicit Lawson method in Table 2. Note that the explicit Lawson method requires the product of the matrix exponential and vector $2s - 1$ times.

As shown in Table 2, the computational efficiency of the proposed method largely depends on the choice of the initial approximation. In particular, the extrapolation methods (III) and (IV) are much more efficient than the other methods. However, as reported in [27], the extrapolation method (III) causes the instability in the numerical solution. We expect worse stability for the extrapolation method (IV). In this sense, we believe the methods (I) and (II) (numerically tested in Sect. 5) are better to use in practice. As summarized in Table 2, the proposed method (I) is better in parallel computation, while the proposed method (II) is better in sequential computation.

Table 2 The number of the computation of the product of matrix exponential and vector

Methods\order	1	2	3	4	5	6	7	8
Explicit Lawson	1	3	5	7	11	13	17	21
Proposed (Is)	2	4	11	15	28	34	53	61
Proposed (Ip)	2	4	6	8	10	12	14	16
Proposed (IIs)	2	3	7	11	17	21	29	37
Proposed (IIp)	2	3	5	9	13	17	23	31
Proposed (IIIs)	–	2	4	7	10	16	21	29
Proposed (IIIp)	–	2	3	4	4	6	6	8
Proposed (IVs)	–	2	3	3	4	4	5	5
Proposed (IVp)	–	2	2	2	2	2	2	2

The symbol (Is) and (Ip) denote the proposed method (I) in sequential and parallel computation, respectively. The other symbols have the same meaning

Appendix 4. Application of the scheme by [16] to the SAV system (18)

In this section, we consider the scheme by [16] (i.e., Definition 1) applied to the SAV system (18).

In this case, the (3) reads

$$\begin{bmatrix} U_i^{(k)} \\ R_{L,i}^{(k)} \\ R_{U,i}^{(k)} \end{bmatrix} = \begin{bmatrix} u_0 \\ r_{L,0} \\ r_{U,0} \end{bmatrix} + h \sum_{j \in [s]} a_{ij} \begin{bmatrix} JLU_j^{(k)} + 2R_{L,j}^{(k)}J\hat{\phi}_{L,j} - 2R_{U,j}^{(k)}J\hat{\phi}_{U,j} \\ -\langle J\hat{\phi}_{L,j}, LU_j^{(k)} \rangle - 2R_{U,j} \langle \hat{\phi}_{L,j}, J\hat{\phi}_{U,j} \rangle \\ -\langle J\hat{\phi}_{U,j}, LU_j^{(k)} \rangle - 2R_{L,j} \langle \hat{\phi}_{L,j}, J\hat{\phi}_{U,j} \rangle \end{bmatrix}, \quad (22)$$

where $\hat{\phi}_{X,j}^{(k-1)} = \phi_X(U_j^{(k-1)})$ for $X \in \{L, U\}$. Let $U^{(k)}$ denote the vector $[U_1^{(k)} \dots U_s^{(k)}]^T$. Then, the first equation in the above equation can be written as

$$(I_s \otimes I_V - hA \otimes JL)U^{(k)} = \mathbf{1}_s \otimes u_0 + 2h \sum_{j \in [s]} A_j \otimes J\hat{\phi}_{L,j}^{(k-1)} - 2h \sum_{j \in [s]} A_j \otimes J\hat{\phi}_{U,j}^{(k-1)},$$

where \otimes denotes the Kronecker product, $I_s \in \mathbb{R}^{s \times s}$ is the identity matrix, $I_V \in \mathcal{L}(V)$ is the identity operator, $A \in \mathbb{R}^{s \times s}$ is the matrix with the entries a_{ij} , $A_j = [a_{1j} \dots a_{sj}]^T \in \mathbb{R}^s$ is the vector, and $\mathbf{1}_s \in \mathbb{R}^s$ is the vector with the entries 1. Since the linear equation in the form

$$(I_s \otimes I_V - hA \otimes JL)x = w \otimes v$$

can be efficiently solved (cf. [40, IV.8]), it is computationally inexpensive to compute

$$\begin{bmatrix} \tilde{u}_1 \\ \vdots \\ \tilde{u}_s \end{bmatrix} = (I_s \otimes I_Y - hA \otimes JL)^{-1} (\mathbf{1}_s \otimes u_0), \tag{23}$$

$$\begin{bmatrix} \tilde{\phi}_{j1}^X \\ \vdots \\ \tilde{\phi}_{js}^X \end{bmatrix} = (I_s \otimes I_Y - hA \otimes JL)^{-1} (A_j \otimes J\hat{\phi}_{X,j}^{(k-1)}). \tag{24}$$

By using this, $U_i^{(k)}$ can be written as

$$U_i^{(k)} = \tilde{u}_i + 2h \sum_{j \in [s]} R_{L,j}^{(k)} \tilde{\phi}_{ji}^L - 2h \sum_{j \in [s]} R_{U,j}^{(k)} \tilde{\phi}_{ji}^U. \tag{25}$$

Then, by using the above equation, the second and the third equations in (22) can be simplified into

$$\begin{bmatrix} I_s + 2h^2 A \tilde{\Phi}_{LL} & -2h^2 A \tilde{\Phi}_{LU} + 2h A \Phi \\ 2h^2 A \tilde{\Phi}_{UL} + 2h A \Phi & I_s - 2h^2 A \tilde{\Phi}_{UU} \end{bmatrix} \begin{bmatrix} R_L^{(k)} \\ R_U^{(k)} \end{bmatrix} = \begin{bmatrix} r_{L,0} \mathbf{1}_s \\ r_{U,0} \mathbf{1}_s \end{bmatrix} - h \begin{bmatrix} A v_L \\ A v_U \end{bmatrix}, \tag{26}$$

where $R_X^{(k)} = [R_{X,1}^{(k)} \dots R_{X,s}^{(k)}]^T$, $\tilde{\Phi}_{XY} \in \mathbb{R}^{s \times s}$ is defined by $(\tilde{\Phi}_{XY})_{ij} = \langle L \hat{\phi}_{X,i}, \tilde{\phi}_{ji}^Y \rangle$, $\Phi \in \mathbb{R}^{s \times s}$ is a diagonal matrix defined by $\Phi_{ii} = \langle \hat{\phi}_{L,i}, J \hat{\phi}_{U,i} \rangle$, $v_X \in \mathbb{R}^s$ is defined by $(v_X)_i = \langle L \hat{\phi}_{X,i}, \tilde{u}_i \rangle$. Consequently, we obtain the following scheme:

Definition 5

Step 0 Prepare $U_i^{(0)} \approx u(c_i h)$ and set $k = 1$.

Step 1 Compute \tilde{u}_i and $\tilde{\phi}_{ij}^X$ for $i, j \in [s]$ by (23) and (24). Solve the linear equation system (26) to obtain $R_L^{(k)}$ and $R_U^{(k)}$. Compute $\{U_i^{(k)}\}_{i=1}^s$ by (25). If some criteria hold, go to Step 2. Otherwise, set $k = k + 1$ and repeat Step 1.

Step 2 Output

$$\begin{bmatrix} u_1 \\ r_{L,1} \\ r_{U,1} \end{bmatrix} = \begin{bmatrix} u_0^{(k)} \\ r_{L,0}^{(k)} \\ r_{U,0}^{(k)} \end{bmatrix} + h \sum_{j \in [s]} b_j \begin{bmatrix} JLU_j^{(k)} + 2R_{L,j}^{(k)} J\hat{\phi}_{L,j} - 2R_{U,j}^{(k)} J\hat{\phi}_{U,j} \\ -\langle J\hat{\phi}_{L,j}, LU_j^{(k)} \rangle - 2R_{U,j} \langle \hat{\phi}_{L,j}, J\hat{\phi}_{U,j} \rangle \\ -\langle J\hat{\phi}_{U,j}, LU_j^{(k)} \rangle - 2R_{L,j} \langle \hat{\phi}_{L,j}, J\hat{\phi}_{U,j} \rangle \end{bmatrix}.$$

Appendix 5. The error behavior of the proposed scheme for the sine-Gordon equation

In this section, we consider the ODE

$$\begin{cases} \dot{u}(t) = v(t), \\ \dot{v}(t) = \delta_x^2 u(t) + 2r(t)\phi(u(t)), \\ \dot{r}(t) = -\langle \phi(u(t)), v(t) \rangle \end{cases} \tag{27}$$

given by a spatial discretization of the sine-Gordon equation.

Theorem 15 *Let (u, v, r) be the solution of the ODE (27). Assume the following conditions:*

- (A1) $U_i^{(0)}$ satisfies $\|U_i^{(0)} - u(c_i h)\| \leq Ch^q$ for each $i \in [s]$.
- (A2) α in the definition of ϕ satisfies $\alpha \geq 2N \Delta x$.
- (A3) The base Runge–Kutta method is of order p .

Then, the solution $(u_1^{(k)}, v_1^{(k)}, r_1^{(k)})$ of the scheme defined by Definition 3 satisfies

$$\begin{aligned} \|u_1^{(k)}(h) - u(h)\| &\leq C'h^{\min\{p,q+2k-1\}+1}, \\ \|v_1^{(k)}(h) - v(h)\| &\leq C'h^{\min\{p,q+2k-2\}+1}, \\ \|r_1^{(k)}(h) - r(h)\| &\leq C'h^{\min\{p,q+2k-2\}+1} \end{aligned}$$

for a sufficiently small step size $h > 0$. Here, C' is a constant depending only on the exact solution (u, v, r) , k, N, A, b , and the constant C in (A1).

Proof In this case, the coupled system considered in the proof of Theorem 7 reads

$$\begin{cases} \dot{u}^{(i)}(t) = v^{(i)}(t), \\ \dot{v}^{(i)}(t) = \delta_x^2 u^{(i)}(t) + 2r^{(i)}(t)\phi(u^{(i-1)}(t)), \\ \dot{r}^{(i)}(t) = -\langle \phi(u^{(i-1)}(t)), v^{(i)}(t) \rangle \end{cases}$$

for $i = 1, 2, \dots, k$. According to the proof of Theorem 7, it is sufficient to show that $\sup_{t \in [0,h]} \|u^{(i)}(t) - u(t)\| \leq C''h^{q+2k}$, $\sup_{t \in [0,h]} \|v^{(i)}(t) - v(t)\| \leq C''h^{q+2k-1}$ and $\sup_{t \in [0,h]} |r^{(i)}(t) - r(t)| \leq C''h^{q+2k-1}$.

First, we introduce several bounds with respect to ϕ . Since $E(u^{(i)}(t)) = \sum_{k=1}^N \cos(u_k^{(i)}(t)) \Delta x \geq -N \Delta x$ and $\|\sin(u^{(i)}(t))\| = \sqrt{\sum_{k=1}^N |\sin(u_k^{(i)})| \Delta x} \leq \sqrt{N \Delta x}$ hold, we see

$$\|\phi(u^{(i)}(t))\| = \frac{1}{2\sqrt{E(u^{(i)}(t)) + \alpha}} \|\sin(u^{(i)}(t))\| \leq \frac{\sqrt{N \Delta x}}{2\sqrt{\alpha - N \Delta x}},$$

which implies $\sup_{t \in [0, h]} \|\phi(u^{(i)}(t))\| \leq 1/2$ by choosing $\alpha = 2N\Delta x$. In addition, under the setting $\alpha = 2N\Delta x$, for any $v, w \in \mathbb{R}^N$, we see

$$\begin{aligned} & \|\phi(v) - \phi(w)\| \\ &= \left\| \frac{1}{2\sqrt{E(v) + \alpha}} \sin(v) - \frac{1}{2\sqrt{E(w) + \alpha}} \sin(w) \right\| \\ &\leq \frac{1}{2\sqrt{E(v) + \alpha}} \|\sin(v) - \sin(w)\| + \left| \frac{1}{2\sqrt{E(v) + \alpha}} - \frac{1}{2\sqrt{E(w) + \alpha}} \right| \|\sin(w)\| \\ &\leq \frac{1}{\sqrt{N\Delta x}} \|v - w\| + \frac{\sqrt{N\Delta x}}{2} \left| \frac{E(v) - E(w)}{\sqrt{(E(v) + \alpha)(E(w) + \alpha)}(\sqrt{E(v) + \alpha} + \sqrt{E(w) + \alpha})} \right| \\ &\leq \frac{1}{\sqrt{N\Delta x}} \|v - w\| + \frac{1}{4N\Delta x} \left| \sum_{k=1}^N (\cos v_k - \cos w_k) \Delta x \right| \\ &= \frac{1}{\sqrt{N\Delta x}} \|v - w\| + \frac{1}{2N\Delta x} \left| \sum_{k=1}^N \sin \frac{v_k + w_k}{2} \sin \frac{v_k - w_k}{2} \Delta x \right| \\ &\leq \frac{1}{\sqrt{N\Delta x}} \|v - w\| + \frac{1}{2N\Delta x} \left\| \sin \frac{v + w}{2} \right\| \left\| \sin \frac{v - w}{2} \right\| \\ &\leq \frac{5}{4\sqrt{N\Delta x}} \|v - w\|. \end{aligned}$$

Therefore, the map ϕ is Lipschitz continuous with the Lipschitz constant $L_\phi := 5/(4\sqrt{N\Delta x})$.

Let us introduce $E_u^{(i)} := \sup_{t \in [0, h]} \|u^{(i)}(t) - u(t)\|$, $E_v^{(i)} := \sup_{t \in [0, h]} \|v^{(i)}(t) - v(t)\|$ and $E_r^{(i)} := \sup_{t \in [0, h]} |r^{(i)}(t) - r(t)|$. Then, by introducing $C_v := \sup_{t \in [0, h]} \|v(t)\|$ and $C_r := \sup_{t \in [0, h]} |r(t)|$, we see

$$\begin{aligned} E_u^{(i)} &= \sup_{t \in [0, h]} \left\| \int_0^t (v^{(i)}(\tau) - v(\tau)) \, d\tau \right\| \\ &\leq h E_v^{(i)}, \\ E_v^{(i)} &= \sup_{t \in [0, h]} \left\| \int_0^t (\delta_x^2 u^{(i)}(\tau) + 2r^{(i)}(\tau)\phi(u^{(i-1)}(\tau)) - \delta_x^2 u(\tau) - 2r(\tau)\phi(u(\tau))) \, d\tau \right\| \\ &\leq h \|\delta_x^2\| E_u^{(i)} + 2 \sup_{t \in [0, h]} \left\| \int_0^t (r^{(i)}(\tau) - r(\tau)) \phi(u^{(i-1)}(\tau)) \, d\tau \right\| \\ &\quad + 2 \sup_{t \in [0, h]} \left\| \int_0^t r(\tau) (\phi(u^{(i-1)}(\tau)) - \phi(u(\tau))) \, d\tau \right\| \\ &\leq h \|\delta_x^2\| E_u^{(i)} + h E_r^{(i)} + 2h C_r L_\phi E_u^{(i-1)}, \\ E_r^{(i)} &= \sup_{t \in [0, h]} \left| \int_0^t (-\langle \phi(u^{(i-1)}(\tau)), v^{(i)}(\tau) \rangle + \langle \phi(u(\tau)), v(\tau) \rangle) \, d\tau \right| \\ &\leq \frac{h}{2} E_v^{(i)} + h C_v L_\phi E_u^{(i-1)}. \end{aligned} \tag{28}$$

The inequalities on $E_v^{(i)}$ and $E_r^{(i)}$ imply

$$\left(1 - \frac{h^2}{2}\right) E_v^{(i)} \leq h \|\delta_x^2\| E_u^{(i)} + hL_\phi (2C_r + hC_v) E_u^{(i-1)}.$$

Therefore, when $h \leq 1$, we see

$$E_v^{(i)} \leq 2h \|\delta_x^2\| E_u^{(i)} + 2hL_\phi (2C_r + hC_v) E_u^{(i-1)}. \tag{29}$$

This inequality and $E_u^{(i)} \leq hE_v^{(i)}$ imply

$$\left(1 - 2h \|\delta_x^2\|\right) E_u^{(i)} \leq 2h^2L_\phi (2C_r + hC_v) E_u^{(i-1)}.$$

Therefore, when $h \leq (4\|\delta_x^2\|)^{-1}$, we see

$$E_u^{(i)} \leq 4h^2L_\phi (2C_r + hC_v) E_u^{(i-1)},$$

which implies $E_u^{(i)} \leq C^{(k)}h^{q+2k}$, where $C^{(k)} = (4L_\phi (2C_r + C_v))^k$. This fact and inequalities (29) and (28) show the theorem. □

As indicated in the proof above, the same argument holds true for the nonlinear Klein–Gordon equation with a general potential function when ϕ is Lipschitz continuous.

Acknowledgements We would like to thank Editage (www.editage.com) for English language editing. The author would also like to thank the anonymous referees for their helpful comments.

Author contribution S.S. wrote the whole manuscript.

Funding Open Access funding provided by The University of Tokyo. The author is supported by a Japan Society for the Promotion of Science (JSPS) Grant-in-Aid for Scientific Research (B) (No. 20H01822) and a JSPS Grant-in-Aid for Early-Career Scientists (No. 22K13955).

Data Availability The code and data that support the findings of this paper are available from the corresponding author upon reasonable request.

Declarations

Conflict of interest The author declares no competing interests.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article’s Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article’s Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

References

- Celledoni, E., Grimm, V., McLachlan, R.I., McLaren, D., O’Neale, D., Owren, B., Quispel, G.R.: Preserving energy resp. dissipation in numerical PDEs using the “average vector field” method. *J. Comp. Phys.* **231**, 6770–6789 (2012)
- Gonzalez, O.: Time integration and discrete Hamiltonian systems. *J. Nonlinear Sci.* **6**, 449–467 (1996)
- McLachlan, R.I., Quispel, G.R.W., Robidoux, N.: Unified approach to Hamiltonian systems, Poisson systems, gradient systems with Lyapunov functions or first integrals. *Phys. Rev. Lett.* **81**, 2399–2403 (1998)
- McLachlan, R.I., Quispel, G.R.W., Robidoux, N.: Geometric integration using discrete gradients. *Philos. Trans. R. Soc. Lond. A Math. Phys. Eng. Sci.* **357**, 1021–1045 (1999)
- Furihata, D.: Finite difference schemes for $\partial u/\partial t = (\partial/\partial x)^\alpha \delta G/\delta u$ that inherit energy conservation or dissipation property. *J. Comput. Phys.* **156**(1), 181–205 (1999). <https://doi.org/10.1006/jcph.1999.6377>
- Furihata, D., Mori, M.: General derivation of finite difference schemes by means of a discrete variation (in Japanese). *Trans. Japan Soc. Indust. Appl.* **8**(3), 317–340 (1998)
- Furihata, D., Matsuo, T.: Discrete variational derivative method—a structure-preserving numerical method for partial differential equations. CRC Press, Boca Raton (2011)
- Cohen, D., Hairer, E.: Linear energy-preserving integrators for Poisson systems. *BIT* **51**(1), 91–101 (2011). <https://doi.org/10.1007/s10543-011-0310-z>
- Besse, C.: A relaxation scheme for the nonlinear Schrödinger equation. *SIAM J. Numer. Anal.* **42**(3), 934–952 (2004). <https://doi.org/10.1137/S0036142901396521>
- Zhang, F., Pérez-García, V.M., Vázquez, L.: Numerical simulation of nonlinear Schrödinger systems: a new conservative scheme. *Appl. Math. Comput.* **71**(2–3), 165–177 (1995)
- Matsuo, T., Furihata, D.: Dissipative or conservative finite-difference schemes for complex-valued nonlinear partial differential equations. *J. Comput. Phys.* **171**(2), 425–447 (2001). <https://doi.org/10.1006/jcph.2001.6775>
- Dahlby, M., Owren, B.: A general framework for deriving integral preserving numerical methods for PDEs. *SIAM J. Sci. Comput.* **33**(5), 2318–2340 (2011). <https://doi.org/10.1137/100810174>
- Yang, X., Han, D.: Linearly first- and second-order, unconditionally energy stable schemes for the phase field crystal model. *J. Comput. Phys.* **330**, 1116–1134 (2017). <https://doi.org/10.1016/j.jcp.2016.10.020>
- Shen, J., Xu, J., Yang, J.: The scalar auxiliary variable (SAV) approach for gradient flows. *J. Comput. Phys.* **353**, 407–416 (2018). <https://doi.org/10.1016/j.jcp.2017.10.021>
- Kemmochi, T., Sato, S.: Scalar auxiliary variable approach for conservative/dissipative partial differential equations with unbounded energy functionals. *BIT* **62**, 903–930 (2022)
- Sato, S., Miyatake, Y., Butcher, J.C.: High-order linearly implicit schemes conserving quadratic invariants. *Appl. Numer. Math.* **187**, 71–88 (2023). <https://doi.org/10.1016/j.apnum.2023.02.005>
- Hochbruck, M., Ostermann, A.: Exponential integrators. *Acta Numer.* **19**, 209–286 (2010). <https://doi.org/10.1017/S0962492910000048>
- Celledoni, E., Cohen, D., Owren, B.: Symmetric exponential integrators with an application to the cubic Schrödinger equation. *Found. Comput. Math.* **8**, 303–317 (2008)
- Mei, L., Huang, L., Huang, S.: Exponential integrators with quadratic energy preservation for linear Poisson systems. *J. Comput. Phys.* **387**, 446–454 (2019)
- Li, Y.-W., Wu, X.: Exponential integrators preserving first integrals or Lyapunov functions for conservative or dissipative systems. *J. Sci. Comput.* **38**(3), 1876–1895 (2016)
- Mei, L., Huang, L., Wu, X.: Energy-preserving exponential integrators of arbitrarily high order for conservative or dissipative systems with highly oscillatory solutions. *J. Comput. Phys.* **442**, 110429 (2021). <https://doi.org/10.1016/j.jcp.2021.110429>
- Li, L.: A new symmetric linearly implicit exponential integrator preserving polynomial invariants or Lyapunov functions for conservative or dissipative systems. *J. Comput. Phys.* **449**, 110800 (2022). <https://doi.org/10.1016/j.jcp.2021.110800>
- Jiang, C., Wang, Y., Cai, W.: A linearly implicit energy-preserving exponential integrator for the nonlinear Klein-Gordon equation. *J. Comput. Phys.* **419**, 109690 (2020)
- Cui, J., Xu, Z., Wang, Y., Jiang, C.: Mass- and energy-preserving exponential Runge-Kutta methods for the nonlinear Schrödinger equation. *Appl. Math. Lett.* **112**, 106770 (2021). <https://doi.org/10.1016/j.aml.2020.106770>

25. Xu, Z., Cai, W., Song, Y., Wang, Y.: Explicit high-order energy-preserving exponential time differencing method for nonlinear Hamiltonian PDEs. *Appl. Math. Comput.* **404**, 126208 (2021)
26. Fu, Y., Hu, D., Xu, Z.: High-order explicit conservative exponential integrator schemes for fractional Hamiltonian PDEs. *Appl. Numer. Math.* **172**, 315–331 (2022). <https://doi.org/10.1016/j.apnum.2021.10.011>
27. Jiang, C., Cui, J., Qian, X., Song, S.: High-order linearly implicit structure-preserving exponential integrators for the nonlinear Schrödinger equation. *J. Sci. Comput.* **90**(1), 66–27 (2022). <https://doi.org/10.1007/s10915-021-01739-x>
28. Lawson, J.D.: Generalized Runge-Kutta processes for stable systems with large Lipschitz constants. *SIAM J. Numer. Anal.* **4**(3), 372–380 (1967)
29. Cooper, G.J.: Stability of Runge-Kutta methods for trajectory problems. *IMA J. Numer. Anal.* **7**(1), 1–13 (1987). <https://doi.org/10.1093/imanum/7.1.1>
30. Butcher, J.C.: B-series—algebraic analysis of numerical methods. Springer Series in Computational Mathematics, vol. 55, p. 310. Springer, Cham (2021). <https://doi.org/10.1007/978-3-030-70956-3>
31. Horn, R.A., Johnson, C.R.: Topics in matrix analysis, p. 607. Cambridge University Press, Cambridge (1991). <https://doi.org/10.1017/CBO9780511840371>
32. Trefethen, L.N., Embree, M.: Spectra and pseudospectra, p. 606. Princeton University Press, Princeton, NJ (2005)
33. Shen, X., Leok, M.: Geometric exponential integrators. *J. Comput. Phys.* **382**, 27–42 (2019). <https://doi.org/10.1016/j.jcp.2019.01.005>
34. Fornberg, B.: A practical guide to pseudospectral methods. Cambridge Monographs on Applied and Computational Mathematics, vol. 1, p. 231. Cambridge University Press, Cambridge (1996). <https://doi.org/10.1017/CBO9780511626357>
35. Butcher, J.C.: On Runge-Kutta processes of high order. *J. Austral. Math. Soc.* **4**, 179–194 (1964)
36. Marchesoni, F.: Exact solutions of the sine-Gordon equation with periodic boundary conditions. *Progr. Theoret. Phys.* **77**(4), 813–824 (1987). <https://doi.org/10.1143/PTP.77.813>
37. Zhang, D.K.: Discovering new Runge-Kutta methods using unstructured numerical search (2019)
38. Owren, B., Zennaro, M.: Continuous explicit Runge-Kutta methods. In: Computational Ordinary Differential Equations (London, 1989). *Inst. Math. Appl. Conf. Ser. New Ser.*, vol. 39, pp. 97–105. Oxford Univ. Press, New York, ??? (1992)
39. Verner, J.H.: Differentiable interpolants for high-order Runge-Kutta methods. *SIAM J. Numer. Anal.* **30**(5), 1446–1466 (1993). <https://doi.org/10.1137/0730075>
40. Hairer, E., Wanner, G.: Solving ordinary differential equations. II, Stiff and Differential-algebraic Problems. Springer Series in Computational Mathematics, vol. 14, p. 614. Springer, Berlin (2010). <https://doi.org/10.1007/978-3-642-05221-7>

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