

Beyond conventional Runge–Kutta methods in numerical integration of ODEs and DAEs by use of structures and local models[☆]

Laurent O. Jay*

Department of Mathematics, 14 MacLean Hall, The University of Iowa, Iowa City, IA 52242-1419, USA

Received 15 July 2005; received in revised form 10 December 2005

Abstract

There are two parts in this paper. In the first part we consider an overdetermined system of differential-algebraic equations (DAEs). We are particularly concerned with Hamiltonian and Lagrangian systems with holonomic constraints. The main motivation is in finding methods based on Gauss coefficients, preserving not only the constraints, symmetry, symplecticness, and variational nature of trajectories of holonomically constrained Hamiltonian and Lagrangian systems, but also having optimal order of convergence. The new class of (s, s) -Gauss–Lobatto specialized partitioned additive Runge–Kutta (SPARK) methods uses greatly the structure of the DAEs and possesses all desired properties. In the second part we propose a unified approach for the solution of ordinary differential equations (ODEs) mixing analytical solutions and numerical approximations. The basic idea is to consider local models which can be solved efficiently, for example analytically, and to incorporate their solution into a global procedure based on standard numerical integration methods for the correction. In order to preserve also symmetry we define the new class of symmetrized Runge–Kutta methods with local model (SRKLM).

© 2006 Elsevier B.V. All rights reserved.

MSC: 34A45; 65L05; 65L06; 65L80

Keywords: Additivity; Correction; DAEs; Gauss methods; Hamiltonian; Holonomic constraints; Lagrangian; Local model; Runge–Kutta methods; Symmetry; Symplecticness; Variational integrators

1. Introduction

In the first part of this paper we consider an overdetermined system of differential-algebraic equations (DAEs), see Section 2. We are particularly concerned with Hamiltonian and Lagrangian systems with holonomic constraints. The main motivation is in finding methods based on Gauss coefficients, preserving not only the constraints, symmetry, symplecticness, and variational nature of trajectories of holonomically constrained Hamiltonian and Lagrangian systems, but also having optimal order of convergence. When applied to nonstiff ordinary differential equations (ODEs), Gauss methods have maximal order of convergence in the class of RK methods. However, for index 3 DAEs, standard Gauss methods, are either divergent or have very low order of convergence. Gauss methods have thus not been considered

[☆] This material is based upon work supported by the National Science Foundation under Grant no. 9983708.

* Tel.: +1 319 335 0898; fax: +1 319 335 0627.

E-mail addresses: ljay@math.uiowa.edu, na.ljay@na-net.ornl.gov (L.O. Jay).

of much practical interest for the numerical solution of high index DAEs in general. In this paper, we propose some modifications to the application of standard RK methods to index 3 DAEs in order to obtain methods with maximal order of convergence. The modifications that we propose have negligible computational cost. The new class of (s, s) -Gauss–Lobatto specialized partitioned additive Runge–Kutta (SPARK) methods is described in Section 3 and makes great use of the structure of the DAEs. The new schemes are constraint-preserving and symmetric. In Section 4, we show that for Hamiltonian and Lagrangian systems with holonomic constraints these schemes preserve symplecticness of the flow and that they satisfy a discrete variational principle: discrete trajectories are stationary with respect to a discrete action.

In the second part of this paper, we propose a unified approach for the solution of ODEs mixing analytical solutions and numerical approximations. When considering a system of ODEs and a given initial value, ideally one would like to obtain directly and explicitly its analytical solution. This is of course generally not possible. In the absence of an explicit analytical solution, one is generally left with two approximation tools: perturbation techniques and numerical integration methods. Perturbation techniques are primarily based on asymptotic expansions. These techniques require at least the analytical solution of a nearby problem, they are often highly technical and they can be applied only to specific situations. For most systems perturbation techniques are not applicable with ease and one is left to solve the problem numerically. In contrast to perturbation techniques, numerical integration methods do not generally incorporate the use of any analytical solution of a nearby problem even when it is available. One aim of this paper is to reconcile both analytical and numerical approaches by giving unified procedures mixing analytical solutions of local models together with numerical approximations in order to find the solution of the original problem more efficiently. This is an idea analogous to the goal of preconditioning when solving linear systems of equations with iterative methods. Mixing analytical solutions of local models with numerical approximations has some advantages. First of all, for a given standard numerical method it generally reduces the error and thus allows to take larger stepsizes. Secondly it allows the development of multiscale procedures based on hierarchical models. The idea of mixing analytical solutions together with numerical methods is certainly not new, but it has not been much explored and fully exploited in ODEs and DAEs. We note that there has been some renewed interest on exponential methods [6,8,15,16], i.e., on methods using the exact solution of linear ODEs. In this paper, we propose a more general approach applicable to different kind of problems and not limited to linear models. The basic idea is to consider local models which can be solved efficiently, for example analytically, and to incorporate their solution into a global procedure based on standard numerical integration methods for the correction, see Section 5. In order to also preserve symmetry we define the new class of symmetrized Runge–Kutta methods with local model (SRKLM) in Section 6. In Section 7 we give some numerical experiments to illustrate some of the theoretical results.

2. A system of implicit differential-algebraic equations

We consider the following class of systems of implicit DAEs

$$\frac{d}{dt}q(t, y) = v(t, y, z), \quad (1a)$$

$$\frac{d}{dt}p(t, y, z) = f(t, y, z) + r(t, y, \psi), \quad (1b)$$

$$0 = g(t, y). \quad (1c)$$

Differentiating the constraints (1c) once with respect to t leads to

$$0 = g_t(t, y) + g_y(t, y)(q_y(t, y))^{-1}(v(t, y, z) - q_t(t, y)). \quad (1d)$$

In mechanics the quantities q , v , p , f , and r represent, respectively, generalized coordinates, generalized velocities, generalized momenta, generalized forces, and reaction forces due to the constraints (1c) [7,20]. The variable $t \in \mathbb{R}$ is the independent variable, the variables $y \in \mathbb{R}^{n_y}$ and $z \in \mathbb{R}^{n_z}$ are called the *differential* variables, and the variables $\psi \in \mathbb{R}^{n_\psi}$ are called the *algebraic* variables. The latter correspond to Lagrange multipliers when the DAEs are derived from some constrained variational principle [7,20]. We have $q \in \mathbb{R}^{n_y}$, $p \in \mathbb{R}^{n_z}$, $g \in \mathbb{R}^{n_\psi}$, $v \in \mathbb{R}^{n_y}$, $f \in \mathbb{R}^{n_z}$, and $r \in \mathbb{R}^{n_z}$. Some differentiability conditions on the above functions and consistency of the initial values y_0, z_0, ψ_0 at

t_0 are assumed to ensure existence and uniqueness of the solution. In a neighborhood of the solution the following conditions are also supposed to be satisfied

$$q_y, \quad p_z, \quad \text{and} \quad \begin{pmatrix} q_y & -v_z & 0 \\ 0 & p_z & -r_\psi \\ g_y & 0 & 0 \end{pmatrix} \text{ are invertible.} \quad (2)$$

The equations (1) include Hamiltonian and Lagrangian systems with holonomic constraints. We give briefly some definitions and theoretical results related to these systems [1,17]. Hamiltonian systems with holonomic constraints $g(q) = 0$ are formulated for a given Hamiltonian $H(q, p)$ as

$$\frac{dq}{dt} = H_p^T(q, p), \quad (3a)$$

$$\frac{dp}{dt} = -H_q^T(q, p) - g_q^T(q)\psi, \quad (3b)$$

$$0 = g(q). \quad (3c)$$

We suppose usually that g_q is of full row rank and that H_{pp}^T is positive definite. Hamiltonian systems have two important properties. Firstly, the Hamiltonian is invariant along a solution, i.e.,

$$H(q(t), p(t)) = \text{Const.}$$

Secondly, the flow $\phi_\tau : (q(t), p(t)) \mapsto (q(t + \tau), p(t + \tau))$ is *symplectic* on the manifold of constraints

$$V := \{(q, p) \in \mathbb{R}^n \times \mathbb{R}^n \mid g(q) = 0, \quad g_q(q)H_p^T(q, p) = 0\}, \quad (4)$$

i.e., on V the symplectic 2-form

$$\sum_{i=1}^n dq^i \wedge dp^i \quad (5)$$

is preserved by the flow ϕ_τ . Lagrangian systems with holonomic constraints $g(q) = 0$ are formulated for a given Lagrangian $L(q, v)$ as

$$\frac{dq}{dt} = v, \quad (6a)$$

$$\frac{d}{dt}L_v^T(q, v) = L_q^T(q, v) - g_q^T(q)\psi, \quad (6b)$$

$$0 = g(q), \quad (6c)$$

the so-called *Euler–Lagrange equations*. We suppose usually that g_q is of full row rank and that L_{vv}^T is positive definite. Lagrangian systems have two important properties. Firstly, the *action* of the Lagrangian

$$\int_{t_a}^{t_b} L(q(t), v(t)) - g^T(q(t))\psi(t) dt$$

is stationary, this is *Hamilton’s variational principle*. Secondly, the flow may be *reversible* with respect to an involution γ of the variables (q, v) , i.e., $\phi_\tau = \gamma^{-1} \circ \phi_\tau^{-1} \circ \gamma$. Lagrangian systems arise for example in classical mechanics for Lagrangians of the form $L = T - U$ where $T = \frac{1}{2}v^T M(q)v$ the kinetic energy with $M(q)$ a positive definite symmetric mass matrix and U is the potential energy. When $U = U(q)$ is independent of v the flow is reversible with respect to a reflection of the velocities $\gamma : (q, v) \mapsto (q, -v)$. Lagrangian systems and Hamiltonian systems are closely related.

Assuming H_{pp}^T or L_{vv}^T invertible we have the following relations between Lagrangian systems and their Hamiltonian counterpart

$$p^T v = H(q, p) + L(q, v),$$

$$p = L_v^T(q, v),$$

$$v = H_p^T(q, p),$$

$$I_n = H_{pp}^T(q, p)L_{vv}^T(q, v).$$

Hence, properties of Lagrangian systems can be transferred to Hamiltonian systems and vice-versa.

3. Specialized partitioned additive Runge–Kutta (SPARK) methods

For $q \equiv y$ in (1a) and $p \equiv z$ in (1b), the standard application of an s -stage Runge–Kutta (RK) method to the semi-explicit system of index 3 DAEs (1) in Hessenberg form is given as follows [5]:

$$Y_i = y_0 + h \sum_{j=1}^s a_{ij} v(T_j, Y_j, Z_j), \quad i = 1, \dots, s,$$

$$Z_i = z_0 + h \sum_{j=1}^s a_{ij} (f(T_j, Y_j, Z_j) + r(T_j, Y_j, \Psi_j)), \quad i = 1, \dots, s,$$

$$0 = g(T_i, Y_i), \quad i = 1, \dots, s,$$

$$y_1 = y_0 + h \sum_{j=1}^s b_j v(T_j, Y_j, Z_j),$$

$$z_1 = z_0 + h \sum_{j=1}^s b_j (f(T_j, Y_j, Z_j) + r(T_j, Y_j, \Psi_j)).$$

An implicit differential equation such as (1a) is usually treated by applying a standard RK method to

$$q_y(t, y) \frac{dy}{dt} = v(t, y, z) - q_t(t, y),$$

giving

$$Y_i = y_0 + h \sum_{j=1}^s a_{ij} Y_j', \quad i = 1, \dots, s,$$

$$q_y(T_j, Y_j) Y_j' = v(T_j, Y_j, Z_j) - q_t(T_j, Y_j), \quad j = 1, \dots, s,$$

and which requires the computation of the partial derivatives q_y and q_t . For $q \equiv y$ in (1a) and $p \equiv z$ in (1b), the standard $s = 1$ Gauss RK method, based on the implicit midpoint rule for ODEs, reads

$$Y_1 = y_0 + h \frac{1}{2} v(T_1, Y_1, Z_1),$$

$$Z_1 = z_0 + h \frac{1}{2} f(T_1, Y_1, Z_1) + h \frac{1}{2} r(T_1, Y_1, \Psi_1),$$

$$0 = g(T_1, Y_1),$$

$$y_1 = y_0 + h v(T_1, Y_1, Z_1),$$

$$z_1 = z_0 + h f(T_1, Y_1, Z_1) + h r(T_1, Y_1, \Psi_1).$$

Unfortunately, this method is divergent in general even when $r(t, y, \psi)$ is linear in the algebraic variables ψ .

The standard definition of RK methods takes neither advantage of the partitioning and additivity of the system (1), nor of the implicitness of the derivatives. In contrast, we propose hereafter a class of methods based on RK coefficients taking advantage of these structures.

Definition 1. One step of an (s, \bar{s}) -specialized partitioned additive Runge–Kutta (SPARK) method applied to the overdetermined differential-algebraic system (1) with consistent initial values (y_0, z_0) at t_0 and stepsize h is given as follows:

$$q(T_i, Y_i) = q_0 + h \sum_{j=1}^s a_{ij} v(T_j, Y_j, Z_j), \quad i = 1, \dots, s, \tag{7a}$$

$$p(T_i, Y_i, Z_i) = p_0 + h \sum_{j=1}^s \hat{a}_{ij} f(T_j, Y_j, Z_j) + h \sum_{j=0}^{\bar{s}} \tilde{a}_{ij} r(\bar{T}_j, \bar{Y}_j, \Psi_j), \quad i = 1, \dots, s, \tag{7b}$$

$$q(\bar{T}_i, \bar{Y}_i) = q_0 + h \sum_{j=1}^s \bar{a}_{ij} v(T_j, Y_j, Z_j), \quad i = 0, 1, \dots, \bar{s}, \tag{7c}$$

$$0 = g(\bar{T}_i, \bar{Y}_i), \quad i = 0, 1, \dots, \bar{s}, \tag{7d}$$

$$q(t_1, y_1) = q_0 + h \sum_{j=1}^s b_j v(T_j, Y_j, Z_j), \tag{7e}$$

$$p(t_1, y_1, z_1) = p_0 + h \sum_{j=1}^s \hat{b}_j f(T_j, Y_j, Z_j) + h \sum_{j=0}^{\bar{s}} \bar{b}_j r(\bar{T}_j, \bar{Y}_j, \Psi_j), \tag{7f}$$

$$0 = g(t_1, y_1), \tag{7g}$$

$$0 = g_x(t_1, y_1) + g_y(t_1, y_1) q_y^{-1}(t_1, y_1) (v(t_1, y_1, z_1) - q_t(t_1, y_1)), \tag{7h}$$

where

$$q_0 := q(t_0, y_0), \quad p_0 := p(t_0, y_0, z_0), \quad t_1 := t_0 + h, \\ T_i := t_0 + c_i h, \quad i = 1, \dots, s, \quad \bar{T}_i := t_0 + \bar{c}_i h, \quad i = 0, 1, \dots, \bar{s}.$$

We have four sets of coefficients (b_j, a_{ij}) , $(\hat{b}_j, \hat{a}_{ij})$, $(\bar{b}_j, \tilde{a}_{ij})$, $(\bar{b}_j, \bar{a}_{ij})$, and we define

$$c_i := \sum_{j=1}^s a_{ij}, \quad i = 1, \dots, s, \quad \bar{c}_i := \sum_{j=1}^s \bar{a}_{ij}, \quad i = 0, 1, \dots, \bar{s}.$$

SPARK coefficients can be expressed by Butcher-tableaux

$$\begin{array}{c|c} c_i & a_{ij} \\ \hline A & b_j \end{array} \quad \begin{array}{c|c} c_i & \hat{a}_{ij} \\ \hline \hat{A} & \hat{b}_j \end{array} \quad \begin{array}{c|c} c_i & \tilde{a}_{ij} \\ \hline \tilde{A} & \bar{b}_j \end{array} \quad \begin{array}{c|c} \bar{c}_i & \bar{a}_{ij} \\ \hline \bar{A} & \end{array}$$

For example the known (2, 1)-Lobatto IIIA-B SPARK method of order 2 [9,10] (an extension of the Störmer/leap-frog/Verlet/RATTLE/SHAKE methods) has Butcher-tableaux of SPARK coefficients

$$\begin{array}{c|c|c} 0 & 0 & 0 \\ \hline 1 & \frac{1}{2} & \frac{1}{2} \\ \hline A & \frac{1}{2} & \frac{1}{2} \end{array} \quad \begin{array}{c|c|c} 0 & \frac{1}{2} & 0 \\ \hline 1 & \frac{1}{2} & 0 \\ \hline \hat{A} & \frac{1}{2} & \frac{1}{2} \end{array} \quad \begin{array}{c|c|c} 0 & \frac{1}{2} & 0 \\ \hline 1 & \frac{1}{2} & 0 \\ \hline \tilde{A} & \frac{1}{2} & \frac{1}{2} \end{array} \quad \begin{array}{c|c|c} 0 & 0 & 0 \\ \hline 1 & \frac{1}{2} & \frac{1}{2} \\ \hline \bar{A} & \frac{1}{2} & \frac{1}{2} \end{array}$$

An (s, \bar{s}) -SPARK method (7) can be seen as an extension of an s -stage standard partitioned Runge–Kutta (PRK) method for partitioned problems without constraints

$$\frac{d}{dt}y = v(t, y, z), \quad \frac{d}{dt}z = f(t, y, z). \tag{8}$$

A similar application of SPARK methods has been proposed for the numerical solution of mechanical systems in [12], see also [13]. SPARK methods are inspired in part by the partitioned RK methods for semi-explicit index 2 DAEs proposed by Murua in [19]. Note that when the RK matrix A is invertible we can express (7c) as

$$q(\bar{T}_i, \bar{Y}_i) = q_0 + \sum_{j=1}^s \eta_{ij}(q(T_j, Y_j) - q_0),$$

where $\eta := \bar{A}A^{-1}$. Similarly, denoting $v^T := b^T A^{-1}$ we can express (7e) as

$$q(t_1, y_1) = q_0 + \sum_{j=1}^s v_j(q(T_j, Y_j) - q_0).$$

To ensure existence and uniqueness of the SPARK solution, we assume the SPARK coefficients to satisfy the following conditions

$$\bar{a}_{0j} = 0, \quad j = 1, \dots, s, \tag{9a}$$

$$\bar{a}_{\bar{s}j} = b_j, \quad j = 1, \dots, s, \tag{9b}$$

$$\sum_{j=1}^s \bar{a}_{ij}c_j = \sum_{j=1}^s \sum_{k=1}^s \bar{a}_{ij}\hat{a}_{jk} = \sum_{j=1}^s \sum_{k=0}^{\bar{s}} \bar{a}_{ij}\tilde{a}_{jk} = \frac{\bar{c}_i^2}{2}, \quad i = 0, 1, \dots, \bar{s}. \tag{9c}$$

A proof for existence and uniqueness of the SPARK solution can be obtained quite similarly to that of [9, Theorem V.4.1]. The condition (9a) implies that $\bar{c}_0 = 0, \bar{T}_0 = t_0, q(\bar{T}_0, \bar{Y}_0) = q(t_0, y_0)$, and thus $\bar{Y}_0 = y_0$. Therefore, $g(\bar{T}_0, \bar{Y}_0) = 0$ is automatically satisfied since we assume $g(t_0, y_0) = 0$. The condition (9b) implies that $g(t_1, y_1) = 0$ is automatically satisfied since $g(\bar{T}_{\bar{s}}, \bar{Y}_{\bar{s}}) = 0$ from (7d) for $i = \bar{s}, t_1 = \bar{T}_{\bar{s}} = t_0 + h, q(t_1, y_1) = q(\bar{T}_{\bar{s}}, \bar{Y}_{\bar{s}})$, and thus $y_1 = \bar{Y}_{\bar{s}}$.

We are especially interested in extending Gauss RK methods for (8) to corresponding (s, s) -SPARK methods (7) having optimal order of convergence $2s$. The Gauss RK coefficients $\hat{a}_{ij} = a_{ij}, \hat{b}_j = b_j$ can be found, e.g., in [3,6]. For the coefficients \bar{b}_i and \bar{c}_i , we take the coefficients of the $(s + 1)$ -stage Lobatto quadrature formula ($\bar{c}_0 = 0, \bar{c}_s = 1$) of order $2s$, they satisfy

$$\sum_{i=0}^s \bar{b}_i \bar{c}_i^{k-1} = \frac{1}{k}, \quad k = 1, \dots, 2s.$$

The coefficients \bar{a}_{ij} are taken according to

$$\sum_{j=1}^s \bar{a}_{ij}c_j^{k-1} = \frac{\bar{c}_i^k}{k}, \quad i = 0, 1, \dots, s \text{ and } k = 1, \dots, s,$$

and the coefficients \tilde{a}_{ij} are then simply determined by

$$\tilde{a}_{ij} = \bar{b}_j \left(1 - \frac{\bar{a}_{ji}}{\bar{b}_i} \right), \quad i = 1, \dots, s, \quad j = 0, 1, \dots, s.$$

We call these methods (s, s) -Gauss–Lobatto SPARK methods. It can be shown that these methods satisfy the conditions (9) and

$$\tilde{a}_{i0} = \bar{b}_0, \quad i = 1, \dots, s, \quad \tilde{a}_{i\bar{s}} = 0, \quad i = 1, \dots, s.$$

The algebraic variable Ψ_s appears only in (7f) and is thus determined by (7h). The (s, s) -Gauss–Lobatto SPARK methods have optimal order of convergence $2s$ [14]. The $(1, 1)$ -Gauss–Lobatto SPARK method of order 2 is given by

$$\begin{aligned}
 q(T_1, Y_1) &= q_0 + h\frac{1}{2}v(T_1, Y_1, Z_1), \\
 p(T_1, Y_1, Z_1) &= p_0 + h\frac{1}{2}f(T_1, Y_1, Z_1) + h\frac{1}{2}r(t_0, y_0, \Psi_0), \\
 q(t_1, y_1) &= q_0 + hv(T_1, Y_1, Z_1), \\
 0 &= g(t_1, y_1), \\
 p(t_1, y_1, z_1) &= p_0 + hf(T_1, Y_1, Z_1) + h\frac{1}{2}r(t_0, y_0, \Psi_0) + h\frac{1}{2}r(t_1, y_1, \Psi_1), \\
 0 &= g_y(t_1, y_1)q_y^{-1}(t_1, y_1)(v(t_1, y_1, z_1) - q_t(t_1, y_1)) + g_t(t_1, y_1).
 \end{aligned}$$

It corresponds to the following Butcher-tableaux of SPARK coefficients

$$\begin{array}{c|c} \frac{1}{2} & \frac{1}{2} \\ \hline \frac{1}{2} & \frac{1}{2} \\ \hline A & 1 \end{array} \quad \begin{array}{c|c} \frac{1}{2} & \frac{1}{2} \\ \hline \widehat{A} & 1 \end{array} \quad \begin{array}{c|c} \frac{1}{2} & \frac{1}{2} \\ \hline \widetilde{A} & \frac{1}{2} \\ \hline \frac{1}{2} & \frac{1}{2} \end{array} \quad \begin{array}{c|c} 0 & 0 \\ \hline 1 & 1 \\ \hline A & \end{array}$$

The $(2, 2)$ -Gauss–Lobatto SPARK method of order 4 has the following Butcher-tableaux of SPARK coefficients

$$\begin{array}{c|c} \frac{1}{2} - \frac{\sqrt{3}}{6} & \frac{1}{4} & \frac{1}{4} - \frac{\sqrt{3}}{6} \\ \hline \frac{1}{2} + \frac{\sqrt{3}}{6} & \frac{1}{4} + \frac{\sqrt{3}}{6} & \frac{1}{4} \\ \hline A & \frac{1}{2} & \frac{1}{2} \end{array} \quad \begin{array}{c|c} \frac{1}{2} - \frac{\sqrt{3}}{6} & \frac{1}{4} & \frac{1}{4} - \frac{\sqrt{3}}{6} \\ \hline \frac{1}{2} + \frac{\sqrt{3}}{6} & \frac{1}{4} + \frac{\sqrt{3}}{6} & \frac{1}{4} \\ \hline \widehat{A} & \frac{1}{2} & \frac{1}{2} \end{array} \\
 \begin{array}{c|c} \frac{1}{2} - \frac{\sqrt{3}}{6} & \frac{1}{6} & \frac{1}{3} - \frac{\sqrt{3}}{6} & 0 \\ \hline \frac{1}{2} + \frac{\sqrt{3}}{6} & \frac{1}{6} & \frac{1}{3} + \frac{\sqrt{3}}{6} & 0 \\ \hline \widetilde{A} & \frac{1}{6} & \frac{2}{3} & \frac{1}{6} \end{array} \quad \begin{array}{c|c} 0 & 0 & 0 \\ \hline \frac{1}{2} & \frac{1}{4} + \frac{\sqrt{3}}{8} & \frac{1}{4} - \frac{\sqrt{3}}{8} \\ \hline 1 & \frac{1}{2} & \frac{1}{2} \\ \hline A & \end{array}$$

We can define the class of *half-explicit (s, s) -SPARK methods* as having SPARK coefficients satisfying

$$\begin{aligned}
 a_{ij} &= 0, \quad i \leq j, \quad \widehat{a}_{ij} = 0, \quad i < j, \\
 \widetilde{a}_{ij} &= 0, \quad i \leq j, \quad \bar{a}_{ij} = 0, \quad i < j, \quad \widetilde{a}_{i,i-1} \neq 0, \quad \bar{a}_{ii} \neq 0.
 \end{aligned}$$

Assuming $f(t, y, z) = f(t, y)$ in (1b), for half-explicit (s, s) -SPARK methods Eqs. (7b) and (7d) for a given index i form a nonlinear system for Z_i and Ψ_{i-1}

$$\begin{aligned}
 p(T_i, Y_i, Z_i) &= C_i + h\widetilde{a}_{i,i-1}r(\bar{T}_{i-1}, \bar{Y}_{i-1}, \Psi_{i-1}), \\
 0 &= g(\bar{T}_i, D_i + h\bar{a}_{ii}v(T_i, Y_i, Z_i)),
 \end{aligned}$$

where $\bar{Y}_{i-1}, Y_i, C_i,$ and D_i are explicitly known expressions.

4. Symplectic and variational SPARK methods

For Hamiltonian systems with holonomic constraints (3), SPARK methods for which the local numerical flow preserves the symplecticness property are characterized as follows:

Theorem 2. *We consider Hamiltonian systems with holonomic constraints (3). If the SPARK method (7) applied to (3) satisfies*

$$\widehat{b}_i = b_i, \quad i = 1, \dots, s, \tag{10a}$$

$$\widehat{b}_i a_{ij} + b_j \widehat{a}_{ji} - \widehat{b}_i b_j = 0, \quad i, j = 1, \dots, s, \tag{10b}$$

$$\bar{b}_i \bar{a}_{ij} + b_j \bar{a}_{ji} - \bar{b}_i b_j = 0, \quad i = 0, 1, \dots, \bar{s}, \quad j = 1, \dots, s, \tag{10c}$$

then the numerical flow $(q_0, p_0) \mapsto (q_1, p_1)$ preserves the symplectic 2-form (5) on V (4).

The proof is given in [14]. The coefficients of the (s, s) -Gauss–Lobatto SPARK methods defined in the previous Section 3 satisfy the symplecticness conditions (10) since Gauss RK coefficients satisfy (10b), and (10c) is satisfied by definition of the coefficients \bar{a}_{ij} . A direct consequence of Theorem 2 is:

Corollary 3. *We consider Lagrangian systems with holonomic constraints (6). If the SPARK method (7) applied to (6) satisfies (10) then the numerical flow $(q_0, v_0) \mapsto (q_1, v_1)$ preserves the Lagrangian symplectic 2-form*

$$\sum_{i=1}^n \sum_{j=1}^n (L_{v^i q^j}(q, v) dq^i \wedge dq^j + L_{v^i v^j}(q, v) dq^i \wedge dv^j)$$

on $W := \{(q, v) \in \mathbb{R}^n \times \mathbb{R}^n \mid g(q) = 0, \quad g_q(q)v = 0\}$.

Assuming the coefficients (b_j, a_{ij}) and (\widehat{b}_j) to be given, to satisfy the symplecticness conditions (10b) we must have

$$\widehat{a}_{ij} = \widehat{b}_j \left(1 - \frac{a_{ji}}{b_i} \right), \quad i, j = 1, \dots, s \text{ when } b_i \neq 0.$$

Assuming the coefficients $(\bar{b}_j, \bar{a}_{ij})$ and (b_j) to be given, to satisfy the symplecticness conditions (10c) we must have

$$\bar{a}_{ij} = \bar{b}_j \left(1 - \frac{\bar{a}_{ji}}{b_i} \right), \quad i = 1, \dots, s, \quad j = 0, 1, \dots, \bar{s} \text{ when } b_i \neq 0.$$

From the symplecticness condition (10c), the assumption $\bar{a}_{0j} = 0$ (9a) implies $b_j = 0$ or $\bar{a}_{j0} = \bar{b}_0$. We are thus particularly interested in SPARK methods satisfying

$$\bar{a}_{i0} = \bar{b}_0, \quad i = 1, \dots, s.$$

From the symplecticness condition (10c), the assumption $\bar{a}_{\bar{s}j} = b_j$ implies $b_j = 0$ or $\bar{a}_{js} = 0$. We are thus particularly interested in SPARK methods satisfying

$$\bar{a}_{i\bar{s}} = 0, \quad i = 1, \dots, s.$$

From this condition the algebraic variable $\Psi_{\bar{s}}$ appears only in (7f) and is determined by (7h).

For Lagrangian systems with holonomic constraints (6) when the SPARK coefficients satisfy the symplecticness conditions (10), the SPARK method (7) can also be derived from a variational point of view following the ideas introduced by Marsden and West in [18]. The variational property in a backward analysis sense of symplectic PRK integrators was derived in [11]. Note also the nonequivalent derivation of Hairer et al. [6] which would consider V_1, \dots, V_s as independent variables and which would remove the constraints corresponding to (7b). This derivation would be difficult to apply here in the presence of holonomic constraints. Following Marsden and West [18], instead of considering the unknown quantities in Eq. (7) as implicit functions of q_0, v_0 , and h , we consider them as implicit

functions of q_0, q_1 , and h . More precisely, assuming $g(q_0) = 0$ and $g(q_1) = 0$ we implicitly define as functions of q_0, q_1 , and h the quantities $p_0, p_1, v_0, v_1, Q_i, P_i, V_i, F_i$ for $i = 1, \dots, s$ and \bar{Q}_i, R_i, Ψ_i for $i = 0, 1, \dots, \bar{s}$ by (7) except that we replace the equation $g(q_1) = 0$ by $0 = g_q(q_0)v_0$. Formally speaking, we should make a distinction between the solution of (7) and the solution of (7) with the equation $g(q_1) = 0$ replaced by $0 = g_q(q_0)v_0$. In any case the solution to one system is also solution to the other under the assumptions $g(q_0) = 0$ and $g_q(q_0)v_0 = 0$ for the first system of equations and $g(q_0) = 0$ and $g(q_1) = 0$ for the second system of equations. Considering the discrete action

$$A_d(q_0, q_1, h) := h \sum_{i=1}^s b_i L(Q_i, V_i) - h \sum_{i=0}^{\bar{s}} \bar{b}_i \Psi_i^T g(\bar{Q}_i),$$

we can show after some lengthy calculations (see proof of Theorem 4 below) that when the SPARK coefficients satisfy the symplecticness assumptions (10), we have the relations

$$p_0 = -\nabla_1 A_d(q_0, q_1, h), \quad p_1 = \nabla_2 A_d(q_0, q_1, h).$$

Therefore, the discrete Euler–Lagrange equations

$$\nabla_2 A_d(q_{n-1}, q_n, h) + \nabla_1 A_d(q_n, q_{n+1}, h) = 0$$

are satisfied for $n = 1, \dots, N - 1$. This implies stationarity of the total discrete action

$$\sum_{n=1}^N A_d(q_{n-1}, q_n, h) \tag{11}$$

with respect to q_n for $n = 1, \dots, N - 1$. This is nothing else but a discrete version of Hamilton’s principle applied to this sum (11). Therefore a SPARK symplectic integrator is also a variational integrator in this sense. We can state more precisely:

Theorem 4. *For Lagrangian systems with holonomic constraints (6) and a corresponding SPARK method (7), suppose q_0 and q_N to be fixed and consistent. Replace the equations $0 = g(q_{n+1})$ for $n = 0, 1, \dots, N - 1$ by $0 = g_q(q_n)v_n$. If the SPARK coefficients satisfy the symplecticness assumptions (10) then we have a variational integrator in the sense of Marsden and West [18], i.e., we have stationarity of the total discrete action (11) with respect to q_n for $n = 1, \dots, N - 1$.*

Proof. We show now the relations

$$-\nabla_1 A_d(q_0, q_1, h) = p_0, \quad \nabla_2 A_d(q_0, q_1, h) = p_1.$$

We have

$$\begin{aligned} -\frac{\partial A_d}{\partial q_0}(q_0, q_1, h) &= -h \sum_{i=1}^s b_i L_q(Q_i, V_i) \frac{\partial Q_i}{\partial q_0} - h \sum_{i=1}^s b_i L_v(Q_i, V_i) \frac{\partial V_i}{\partial q_0} \\ &\quad + h \sum_{i=0}^{\bar{s}} \bar{b}_i \Psi_i^T \left(g_q(\bar{Q}_i) \frac{\partial \bar{Q}_i}{\partial q_0} \right) + h \sum_{i=0}^{\bar{s}} \bar{b}_i g^T(\bar{Q}_i) \frac{\partial \Psi_i}{\partial q_0} \\ &= -h \sum_{i=1}^s b_i F_i^T \left(I + h \sum_{j=1}^s a_{ij} \frac{\partial V_j}{\partial q_0} \right) - h \sum_{i=1}^s b_i P_i^T \frac{\partial V_i}{\partial q_0} \\ &\quad + h \sum_{i=0}^{\bar{s}} \bar{b}_i \Psi_i^T g_q(\bar{Q}_i) \left(I + h \sum_{j=1}^s \bar{a}_{ij} \frac{\partial V_j}{\partial q_0} \right) + h \sum_{i=0}^{\bar{s}} \bar{b}_i g^T(\bar{Q}_i) \frac{\partial \Psi_i}{\partial q_0} \end{aligned}$$

$$\begin{aligned}
 &= -h \sum_{i=1}^s b_i F_i^T I - h^2 \sum_{i=1}^s \sum_{j=1}^s b_i a_{ij} F_i^T \frac{\partial V_j}{\partial q_0} \\
 &\quad - h \sum_{i=1}^s b_i \left(p_0^T + h \sum_{j=1}^s \widehat{a}_{ij} F_j^T + h \sum_{j=0}^{\bar{s}} \widetilde{a}_{ij} R_j^T \right) \frac{\partial V_i}{\partial q_0} \\
 &\quad - h \sum_{i=0}^{\bar{s}} \bar{b}_i R_i^T I - h^2 \sum_{i=0}^{\bar{s}} \sum_{j=1}^s \bar{b}_i \bar{a}_{ij} R_i^T \frac{\partial V_j}{\partial q_0} + h \sum_{i=0}^{\bar{s}} \bar{b}_i g^T(\bar{Q}_i) \frac{\partial \Psi_i}{\partial q_0} \\
 &= -h \sum_{j=1}^s b_j F_j^T I - h^2 \sum_{i=1}^s \sum_{j=1}^s (b_j a_{ji} + b_i \widehat{a}_{ij}) F_j^T \frac{\partial V_i}{\partial q_0} \\
 &\quad - p_0^T h \sum_{i=1}^s b_i \frac{\partial V_i}{\partial q_0} - h^2 \sum_{i=1}^s \sum_{j=0}^{\bar{s}} b_i \widetilde{a}_{ij} R_j^T \frac{\partial V_i}{\partial q_0} - h \sum_{i=0}^{\bar{s}} \bar{b}_i R_i^T I \\
 &\quad - h^2 \sum_{i=0}^{\bar{s}} \sum_{j=1}^s \bar{b}_i \bar{a}_{ij} R_i^T \frac{\partial V_j}{\partial q_0} + h \sum_{i=0}^{\bar{s}} \bar{b}_i g^T(\bar{Q}_i) \frac{\partial \Psi_i}{\partial q_0}.
 \end{aligned}$$

From $q_1 = q_0 + h \sum_{i=1}^s b_i V_i$ we have

$$0 = I + h \sum_{i=1}^s b_i \frac{\partial V_i}{\partial q_0},$$

hence

$$\begin{aligned}
 -\frac{\partial A_d}{\partial q_0}(q_0, q_1, h) &= -h^2 \sum_{i=1}^s \sum_{j=1}^s (b_j a_{ji} + b_i \widehat{a}_{ij} - b_j b_i) F_j^T \frac{\partial V_i}{\partial q_0} + p_0^T \\
 &\quad - h^2 \sum_{i=0}^{\bar{s}} \sum_{j=1}^s (b_j \widetilde{a}_{ji} + \bar{b}_i \bar{a}_{ij} - \bar{b}_i b_j) R_i^T \frac{\partial V_j}{\partial q_0} + h \sum_{i=0}^{\bar{s}} \bar{b}_i g^T(\bar{Q}_i) \frac{\partial \Psi_i}{\partial q_0}.
 \end{aligned}$$

From $g(\bar{Q}_i) = 0$ and the symplecticness assumptions (10) we obtain the desired result

$$-\frac{\partial A_d}{\partial q_0}(q_0, q_1, h) = p_0^T.$$

Similarly, we have

$$\begin{aligned}
 \frac{\partial A_d}{\partial q_1}(q_0, q_1, h) &= h \sum_{i=1}^s b_i L_q(Q_i, V_i) \frac{\partial Q_i}{\partial q_1} + h \sum_{i=1}^s b_i L_v(Q_i, V_i) \frac{\partial V_i}{\partial q_1} \\
 &\quad - h \sum_{i=0}^{\bar{s}} \bar{b}_i \Psi_i^T \left(g_q(\bar{Q}_i) \frac{\partial \bar{Q}_i}{\partial q_1} \right) - h \sum_{i=0}^{\bar{s}} \bar{b}_i g^T(\bar{Q}_i) \frac{\partial \Psi_i}{\partial q_1}
 \end{aligned}$$

$$\begin{aligned}
&= h \sum_{i=1}^s b_i F_i^T \left(h \sum_{j=1}^s a_{ij} \frac{\partial V_j}{\partial q_1} \right) + h \sum_{i=1}^s b_i P_i^T \frac{\partial V_i}{\partial q_1} \\
&\quad - h \sum_{i=0}^{\bar{s}} \bar{b}_i \Psi_i^T g_q(\bar{Q}_i) \left(h \sum_{j=1}^s \bar{a}_{ij} \frac{\partial V_j}{\partial q_1} \right) - h \sum_{i=0}^{\bar{s}} \bar{b}_i g^T(\bar{Q}_i) \frac{\partial \Psi_i}{\partial q_1} \\
&= h^2 \sum_{i=1}^s \sum_{j=1}^s b_i a_{ij} F_i^T \frac{\partial V_j}{\partial q_1} \\
&\quad + h \sum_{i=1}^s b_i \left(p_0^T + h \sum_{j=1}^s \hat{a}_{ij} F_j^T + h \sum_{j=0}^{\bar{s}} \tilde{a}_{ij} R_j^T \right) \frac{\partial V_i}{\partial q_1} \\
&\quad + h^2 \sum_{i=0}^{\bar{s}} \sum_{j=1}^s \bar{b}_i \bar{a}_{ij} R_i^T \frac{\partial V_j}{\partial q_1} - h \sum_{i=0}^{\bar{s}} \bar{b}_i g^T(\bar{Q}_i) \frac{\partial \Psi_i}{\partial q_1} \\
&= h^2 \sum_{i=1}^s \sum_{j=1}^s b_j a_{ji} F_j^T \frac{\partial V_i}{\partial q_1} \\
&\quad + h \sum_{i=1}^s b_i \left(p_1^T + h \sum_{j=1}^s (\hat{a}_{ij} - \hat{b}_j) F_j^T + h \sum_{j=0}^{\bar{s}} (\tilde{a}_{ij} - \bar{b}_j) R_j^T \right) \frac{\partial V_i}{\partial q_1} \\
&\quad + h^2 \sum_{i=1}^s \sum_{j=0}^{\bar{s}} \bar{b}_j \bar{a}_{ji} R_j^T \frac{\partial V_i}{\partial q_1} - h \sum_{i=0}^{\bar{s}} \bar{b}_i g^T(\bar{Q}_i) \frac{\partial \Psi_i}{\partial q_1} \\
&= h^2 \sum_{i=1}^s \sum_{j=1}^s (b_j a_{ji} + b_i \hat{a}_{ij} - b_i \hat{b}_j) F_j^T \frac{\partial V_i}{\partial q_1} + p_1^T h \sum_{i=1}^s b_i \frac{\partial V_i}{\partial q_1} \\
&\quad + h^2 \sum_{i=1}^s \sum_{j=0}^{\bar{s}} (b_i \tilde{a}_{ij} - b_i \bar{b}_j + \bar{b}_j \bar{a}_{ji}) R_j^T \frac{\partial V_i}{\partial q_1} - h \sum_{i=0}^{\bar{s}} \bar{b}_i g^T(\bar{Q}_i) \frac{\partial \Psi_i}{\partial q_1}.
\end{aligned}$$

From $q_1 = q_0 + h \sum_{i=1}^s b_i V_i$ we have

$$I = h \sum_{i=1}^s b_i \frac{\partial V_i}{\partial q_1},$$

hence from $g(\bar{Q}_i) = 0$ and the symplecticness assumptions (10) we obtain the desired result

$$\frac{\partial A_d}{\partial q_1}(q_0, q_1, h) = p_1^T. \quad \square$$

For (s, s) -Gauss–Lobatto SPARK methods we summarize our findings in the following theorem:

Theorem 5. For the overdetermined differential-algebraic system (1) the (s, s) -Gauss–Lobatto SPARK method (7) is constraint-preserving, symmetric, and of maximal order $2s$, i.e.,

$$y_n - y(t_n) = O(h^{2s}), \quad z_n - z(t_n) = O(h^{2s})$$

for $t_n - t_0 = nh \leq \text{Const}$. For holonomically constrained Hamiltonian systems (3) and Lagrangian systems (6) these methods are also symplectic and variational.

5. Use of local models in mixed analytical/numerical integration of ODEs

We consider a system of ODEs in \mathbb{R}^n

$$\frac{dy}{dt} = f(t, y) \quad (12)$$

with a given initial value $y_0 \in \mathbb{R}^n$ at t_0 . Associated to this system of ODEs (12) we consider in a neighborhood of t_0 and y_0 a local model

$$\frac{dz}{dt} = g(t, z), \quad (13)$$

assumed to be solvable sufficiently accurately and more efficiently than (12), for example by an explicit analytical expression. When no local model problem (13) is associated to (12), we can simply consider by default the trivial local model $dz/dt = 0$, i.e., $g(t, z) \equiv 0$. In this paper we will assume that $z(t)$ can be obtained directly in analytical form. In fact, we can replace the exact values of $z(t)$ by sufficiently accurate approximations provided they do not deteriorate significantly the global procedure. In this paper, we denote by $y(t, r, y_r)$ the exact solution at t of (12) passing through y_r at r . Analogously $z(t, r, z_r)$ denotes the exact solution at t of (13) passing through z_r at r .

The idea that we have in mind is analogous to the goal of preconditioning for the iterative solution of systems linear of equations. Starting from a system of linear equations $Fy = b$ to be solved, the main goal of preconditioning is to find its solution more efficiently by using auxiliary linear systems $Gz = c$ where G is an approximation to F and where the solution of $Gz = c$ can be obtained much more efficiently than the solution of $Fy = b$. Here, the analogue of $Fy = b$ is (12), the analogue of $Gz = c$ is (13), and the analogue of a linear iterative method is given by a numerical integration method.

We want a unified procedure in the following sense:

- its result must reduce to a standard numerical discretization of (12) for the trivial local model $g \equiv 0$;
- its result must reduce to the exact solution of (12) when $g \equiv f$ and (13) is solved exactly for arbitrary initial conditions;
- its order should be at least equal to the standard order of the numerical discretization used, i.e., to the order corresponding to $g \equiv 0$.

A first approach is described as follows. Denoting $y^r(t) := y(t, r, y_r)$ the Groebner–Aleksseev formula reads

$$y^r(t) = z(t, r, y_r) + \int_r^t \partial_3 z(t, s, y^r(s)) (f(s, y^r(s)) - g(s, y^r(s))) ds. \quad (14)$$

This is an integral equation for y^r . For example, for

$$\frac{dy}{dt} = Ay + d(t, y), \quad \frac{dz}{dt} = Az, \quad (15)$$

we have $z(t, r, z_r) = e^{(t-r)A} z_r$ and (14) becomes in this situation

$$y^r(t) = e^{(t-r)A} y_r + \int_r^t e^{(t-s)A} d(s, y^r(s)) ds. \quad (16)$$

For $d(t, y) = b(t)$ independent of y , it corresponds to the well-known variation-of-constants formula. For linear highly oscillatory problems several discretizations based on (16) have been proposed, see, e.g., [6, Chapter XIII]. As a simple example, considering the left rectangle rule to approximate the integral part of (16), one obtains the *Lawson-explicit Euler exponential method*

$$y_{n+1} = e^{h_n A} y_n + h_n e^{h_n A} d(t_n, y_n) = e^{h_n A} (y_n + h_n d(t_n, y_n)), \quad (17)$$

where $h_n := t_{n+1} - t_n$. When $A \equiv 0$ we obtain the standard explicit Euler method. For $d(t, y) \equiv 0$ the numerical solution is exact. The standard order of the method is easily seen to be equal to one.

Instead of considering the Groebner–Aleksseev formula (14) as a starting point to derive methods for solving (12) using (13), we will consider in this paper a different and conceptually simpler approach. On each subinterval $[t_n, t_{n+1}]$ we introduce the correction

$$\delta^n(t) := y^n(t) - z^n(t),$$

where $y^n(t) := y(t, t_n, y_n)$ and $z^n(t) := z(t, t_n, y_n)$. The correction δ^n satisfies the following nonautonomous system of ODEs

$$\frac{d\delta^n}{dt} = f(t, z^n(t) + \delta^n) - g(t, z^n(t)) \quad (18)$$

with initial condition $\delta^n(t_n) = 0$. This initial value problem can be integrated numerically by any one-step numerical integration method such as a Runge–Kutta method. We thus obtain a numerical approximation δ_{n+1} to $\delta^n(t_{n+1})$. We recover a numerical approximation y_{n+1} to $y^n(t_{n+1})$ by taking

$$y_{n+1} := z^n(t_{n+1}) + \delta_{n+1}.$$

For a Runge–Kutta scheme the resulting method is called a *Runge–Kutta method with local model (RKLM)*. Note that this approach is not a defect correction technique [2,21,22]. When applied to (15) and considering the explicit Euler method applied to (18) this procedure leads to what can be called the *explicit Euler for correction exponential method*

$$y_{n+1} = e^{h_n A} y_n + h_n d(t_n, y_n) \quad (19)$$

which is not equivalent to the method (17). This method also has the same properties of reducing to the standard explicit Euler method when $A \equiv 0$, of leading to the exact solution when $d(t, y) \equiv 0$, and of being of order one. Both methods (17) and (19) can also be interpreted as splitting methods, see below.

Unfortunately, even when the underlying RK scheme is symmetric the resulting RKLM is generally not symmetric. For example, consider the midpoint rule applied to (18) and the problem (15) with $d(t, y) \equiv b(t)$, we obtain

$$y_{n+1} = e^{h_n A} y_n + h_n \left(I - \frac{h_n}{2} A \right)^{-1} b \left(t_n + \frac{h_n}{2} \right). \quad (20)$$

Exchanging $y_{n+1} \leftrightarrow y_n$ and $h_n \leftrightarrow -h_n$ we obtain the adjoint method

$$y_{n+1}^* = e^{h_n A} y_n + h_n e^{h_n A} \left(I + \frac{h_n}{2} A \right)^{-1} b \left(t_n + \frac{h_n}{2} \right) \quad (21)$$

which is clearly a different method. In this paper, we will show how symmetry can still be preserved for an underlying RK scheme using an approach based on integrating the correction ODEs (18), see the symmetrized Runge–Kutta methods with local model (SRKLM) of Section 6. An approach related to methods based on correction is to consider *splitting methods*. For example one rewrites the system of ODEs (12) as

$$\frac{dy}{dt} = g(t, y) + d(t, y),$$

where $d(t, y) := f(t, y) - g(t, y)$ and solve for g and d separately and not necessarily with identical methods. For example one can consider the order 1 splitting

$$(t_n + h_n, y_{n+1}) := (G_{h_n} \circ D_{h_n})(t_n, y_n),$$

where $D_h(t, v) := (t + h, d_h(t, v))$ with $d_h(t, v)$ an approximation to $u(t + h, t, v)$ the exact solution at $t + h$ of $du/dt = d(t, u)$, and $G_h(t, u) := (t, g_h(t, u))$ with $g_h(t, u)$ an approximation to $v(t + h, t, u)$ the exact solution at $t + h$ of $dv/dt = g(t, v)$. Taking $g(t, y) = Ay$, $d_h(t, v) := v + hd(t, v)$ and $g_h(t, u) := v(t + h, t, u)$ in (15) leads to the method (17). Similarly, considering the order 1 splitting

$$(t_n + h_n, y_{n+1}) := (D_{h_n} \circ G_{h_n})(t_n, y_n)$$

for (15) one obtains the method (19). To obtain methods of order 2 one can consider the Strang splitting

$$(t_n + h_n, y_{n+1}) := (D_{h_n/2} \circ G_{h_n} \circ D_{h_n/2})(t_n, y_n)$$

with d_h and g_h symmetric approximations. We will not discuss splitting methods in this paper, see, e.g., [6] for an introduction to splitting methods.

6. Symmetrized Runge–Kutta methods with local model (SRKLM)

As mentioned before, a RKLM based on the application of a standard symmetric RK scheme to the correction ODEs (18) is generally not symmetric. In this section, we propose some new methods based on Runge–Kutta coefficients and correction ODEs preserving the symmetry of the underlying scheme. To simplify the notation we assume that $n = 0$, we consider the interval $[t_0, t_1]$, and we denote the stepsize by $h := t_1 - t_0$. For the numerical procedure an initial/input value y_0 at t_0 is supposed to be given. We will define below a procedure to obtain the numerical value y_1 at t_1 . We define $z_0(t)$ and $z_1(t)$ as the exact solutions of (13) satisfying $z_0(t_0) = y_0$ and $z_1(t_1) = y_1$, respectively. First, let us consider the application of a Runge–Kutta method to the correction ODEs (18) with initial condition $\delta(t_0) := y_0 - z_0(t_0) = 0$. We obtain

$$\Delta_i = h \sum_{j=1}^s a_{ij}(f(T_j, z_0(T_j) + \Delta_j) - g(T_j, z_0(T_j))), \quad i = 1, \dots, s,$$

$$\delta_1 = h \sum_{j=1}^s b_j(f(T_j, z_0(T_j) + \Delta_j) - g(T_j, z_0(T_j))),$$

where $T_j := t_0 + c_j h$. Rewritten in terms of the original y -variable $y = z + \delta$ we obtain the *forward* value y_1^+

$$Y_i^+ = z_0(T_i) + h \sum_{j=1}^s a_{ij}(f(T_j, Y_j^+) - g(T_j, z_0(T_j))), \quad i = 1, \dots, s, \tag{22a}$$

$$y_1^+ = z_0(t_1) + h \sum_{j=1}^s b_j(f(T_j, Y_j^+) - g(T_j, z_0(T_j))). \tag{22b}$$

To define y_1 we will need $z_1(t)$ which in turn depends on y_1 through the relation $z_1(t_1) = y_1$. To simplify the discussion we suppose for an instant that y_1 and therefore $z_1(t)$ are implicitly given. We can define the *backward* value y_0^-

$$Y_i^- = z_1(T_i) - h \sum_{j=1}^s (b_j - a_{ij})(f(T_j, Y_j^-) - g(T_j, z_1(T_j))), \quad i = 1, \dots, s, \tag{22c}$$

$$y_0^- = z_1(t_0) - h \sum_{j=1}^s b_j(f(T_j, Y_j^-) - g(T_j, z_1(T_j))). \tag{22d}$$

Now we need an extra condition to determine y_1 . We take

$$y_1 - y_1^+ = y_0 - y_0^-, \tag{22e}$$

and we call the resulting method (22) a *symmetrized Runge–Kutta method with local model (SRKLM)*. This definition is motivated by Theorem 6 below. One important point is that an SRKLM method is symmetric when the underlying RK coefficients are symmetric.

For example consider the problem (15) with $d(t, y) \equiv b(t)$ and the RK coefficients of the midpoint rule, Eq. (22e) of the midpoint SRKLM gives

$$y_1 - e^{hA} y_0 - h \left(I - \frac{h}{2} A \right)^{-1} b \left(t_0 + \frac{h}{2} \right) = y_0 - e^{-hA} y_1 + h \left(I + \frac{h}{2} A \right)^{-1} b \left(t_0 + \frac{h}{2} \right)$$

leading to

$$y_1 = e^{hA} y_0 + h e^{hA} (I + e^{hA})^{-1} \left(\left(I + \frac{h}{2} A \right)^{-1} + \left(I - \frac{h}{2} A \right)^{-1} \right) b \left(t_0 + \frac{h}{2} \right)$$

and which differs from (20) and (21).

Theorem 6. (1) For $|h| \leq h_0$ with $h_0 > 0$ sufficiently small there exists a unique SRKLM solution;

(2) The order of a SRKLM method is equal to the standard order of the underlying standard RK method;

(3) When $g \equiv 0$ we have $y_1 = y_1^+$, $y_0^- = y_0$, $Y_i^- = Y_i^+$ for $i = 1, \dots, s$, and y_1 is simply the result of the standard RK method with coefficients (b_j, a_{ij}, c_i) applied to (12);

(4) When $g \equiv f$ we have $y_1 = y_1^+ = y(t_1)$, $y_0^- = y_0$, $Y_i^- = Y_i^+ = y(T_i)$ for $i = 1, \dots, s$ where $y(t)$ is the exact solution of (12) with initial value $y(t_0) = y_0$;

(5) If the RK coefficients (b_j, a_{ij}, c_i) are symmetric, i.e.,

$$c_i = 1 - c_{s+1-i}, \quad i = 1, \dots, s,$$

$$a_{s+1-i, s+1-j} + a_{ij} = b_j = b_{s+1-j}, \quad i, j = 1, \dots, s$$

then the corresponding SRKLM method is symmetric.

Proof. To prove the first assertion we define

$$F(h, y_1) := \frac{1}{2}(y_1 - y_1^+ - (y_0 - y_0^-)).$$

We have

$$F(0, y_0) = 0$$

since for $h = 0$ the solution is simply given by $y_1 = y_1^+ = y_0^- = Y_i^- = Y_i^+ = y_0$. We also have

$$\frac{\partial F}{\partial y_1}(0, y_0) = I$$

since $y_0^- = y_0 + O(h)$ and both y_1^+ , y_0 do not depend on y_1 . Hence, by the implicit function theorem, we have existence and uniqueness of the SRKLM solution for $|h|$ sufficiently small.

For the order of a SRKLM method, by assumption we have $y_1^+ - y(t_1, t_0, y_0) = O(h^{p+1})$ and $y_0^- - y(t_0, t_1, y_1) = O(h^{p+1})$ where p is the order of the underlying standard RK method. We want to estimate $y_1 - y(t_1, t_0, y_0)$. We have

$$y_1 - y(t_1, t_0, y_0) = y_1 - y_1^+ + y_1^+ - y(t_1, t_0, y_0) = y_0 - y_0^- + O(h^{p+1})$$

by (22e). We rewrite

$$y_0 - y_0^- = y_0 - y(t_0, t_1, y_1) + y(t_0, t_1, y_1) - y_0^- = y_0 - y(t_0, t_1, y_1) + O(h^{p+1}).$$

Since $y_0 = y(t_0, t_1, y(t_1, t_0, y_0))$, by a simple Taylor series expansion with respect to the third argument we get

$$\begin{aligned} y_0 - y(t_0, t_1, y_1) &= y(t_0, t_1, y_1 + (y(t_1, t_0, y_0) - y_1)) - y(t_0, t_1, y_1) \\ &= \partial_{3y}(t_0, t_1, y_1)(y(t_1, t_0, y_0) - y_1) + O(\|y(t_1, t_0, y_0) - y_1\|^2) \\ &= y(t_1, t_0, y_0) - y_1 \\ &\quad + O(h\|y(t_1, t_0, y_0) - y_1\| + \|y(t_1, t_0, y_0) - y_1\|^2) \end{aligned}$$

since

$$\partial_{3y}(t_0, t_1, y_1) = I + O(h).$$

Therefore, collecting the above estimates we obtain

$$y_1 - y(t_1, t_0, y_0) = O(h\|y_1 - y(t_1, t_0, y_0)\| + \|y_1 - y(t_1, t_0, y_0)\|^2 + h^{p+1}),$$

leading to $y_1 - y(t_1, t_0, y_0) = O(h^{p+1})$.

For the third assertion, when $g \equiv 0$, we have $z_0(t) \equiv y_0$, $z_1(t) \equiv y_1$, and the SRKLM method (22) reads

$$Y_i^+ = y_0 + h \sum_{j=1}^s a_{ij} f(T_j, Y_j^+), \quad i = 1, \dots, s, \quad y_1^+ = y_0 + h \sum_{j=1}^s b_j f(T_j, Y_j^+),$$

$$Y_i^- = y_1 - h \sum_{j=1}^s (b_j - a_{ij}) f(T_j, Y_j^-), \quad i = 1, \dots, s, \quad y_0^- = y_1 - h \sum_{j=1}^s b_j f(T_j, Y_j^-),$$

$$y_1 - y_1^+ = y_0 - y_0^-.$$

The last three equations are equivalent to

$$Y_i^- = y_0^- + h \sum_{j=1}^s a_{ij} f(T_j, Y_j^-) \quad i = 1, \dots, s, \quad y_1 = y_0^- + h \sum_{j=1}^s b_j f(T_j, Y_j^-),$$

$$y_0 - y_0^- = \frac{h}{2} \sum_{j=1}^s b_j (f(T_j, Y_j^-) - f(T_j, Y_j^+)).$$

Clearly the equalities $Y_i^- = Y_i^+$ for $i = 1, \dots, s$, $y_0^- = y_0$, and $y_1 = y_1^+$ are satisfied by the solution to the above equations. The solution is thus simply the result of the standard RK method with coefficients (b_j, a_{ij}, c_i) applied to (12).

For the fourth assertion, when $g \equiv f$ we have

$$Y_i^+ = z_0(T_i) + h \sum_{j=1}^s a_{ij} (f(T_j, Y_j^+) - f(T_j, z_0(T_j))), \quad i = 1, \dots, s,$$

$$y_1^+ = z_0(t_1) + h \sum_{j=1}^s b_j (f(T_j, Y_j^+) - f(T_j, z_0(T_j))),$$

$$Y_i^- = z_1(T_i) - h \sum_{j=1}^s (b_j - a_{ij}) (f(T_j, Y_j^-) - f(T_j, z_1(T_j))), \quad i = 1, \dots, s,$$

$$y_0^- = z_1(t_0) - h \sum_{j=1}^s b_j (f(T_j, Y_j^-) - f(T_j, z_1(T_j))),$$

$$y_1 - y_1^+ = y_0 - y_0^-.$$

Defining $y(t)$ as the exact solution of (12) passing through y_0 at t_0 , it can be easily checked that the solution to the above equations is given by $y_1 = y_1^+ = y(t_1)$, $y_0^- = y_0$, and $Y_i^- = Y_i^+ = y(T_i)$ for $i = 1, \dots, s$, and that $z_0(t) = z_1(t) = y(t)$ is satisfied.

Finally, it remains to prove the assertion on symmetry. Exchanging $y_1 \leftrightarrow y_0$ and $h \leftrightarrow -h$ and $t_0 \leftrightarrow t_1$ in (22) we obtain the adjoint SRKLM equations

$$\tilde{Y}_i^+ = \tilde{z}_1(\tilde{T}_i) - h \sum_{j=1}^s a_{ij} (f(\tilde{T}_j, \tilde{Y}_j^+) - g(\tilde{T}_j, \tilde{z}_1(\tilde{T}_j))), \quad i = 1, \dots, s,$$

$$\tilde{y}_1^+ = \tilde{z}_1(t_1) - h \sum_{j=1}^s b_j (f(\tilde{T}_j, \tilde{Y}_j^+) - g(\tilde{T}_j, \tilde{z}_1(\tilde{T}_j))),$$

$$\tilde{Y}_i^- = z_0(\tilde{T}_i) + h \sum_{j=1}^s (b_j - a_{ij}) (f(\tilde{T}_j, \tilde{Y}_j^-) - g(\tilde{T}_j, z_0(\tilde{T}_j))), \quad i = 1, \dots, s,$$

$$\tilde{y}_0^- = z_0(t_0) + h \sum_{j=1}^s b_j (f(\tilde{T}_j, \tilde{Y}_j^-) - g(\tilde{T}_j, z_0(\tilde{T}_j))),$$

$$y_0 - \tilde{y}_1^+ = \tilde{y}_1 - \tilde{y}_0^-,$$

where $\tilde{T}_i := t_0 + (1 - c_i)h$ for $i = 1, \dots, s$. By symmetry of the nodes c_i we have $\tilde{T}_i = T_{s+1-i}$ for $i = 1, \dots, s$. Using symmetry of the RK coefficients a_{ij} and of the weights b_j we obtain

$$\tilde{Y}_{s+1-i}^- = z_0(T_i) + h \sum_{j=1}^s (b_{s+1-j} - a_{s+1-i, s+1-j}) (f(T_j, \tilde{Y}_{s+1-j}^-) - g(T_j, z_0(T_j)))$$

$$= z_0(T_i) + h \sum_{j=1}^s a_{ij} (f(T_j, \tilde{Y}_{s+1-j}^-) - g(T_j, z_0(T_j))), \quad i = 1, \dots, s,$$

$$\tilde{y}_0^- = z_0(t_0) + h \sum_{j=1}^s b_{s+1-j} (f(T_j, \tilde{Y}_{s+1-j}^-) - g(T_j, z_0(T_j))),$$

$$= z_0(t_0) + h \sum_{j=1}^s b_j (f(T_j, \tilde{Y}_{s+1-j}^-) - g(T_j, z_0(T_j))),$$

$$\tilde{Y}_{s+1-i}^+ = \tilde{z}_1(T_i) - h \sum_{j=1}^s a_{s+1-i, s+1-j} (f(T_j, \tilde{Y}_{s+1-j}^+) - g(T_j, \tilde{z}_1(T_j)))$$

$$= \tilde{z}_1(T_i) - h \sum_{j=1}^s (b_j - a_{ij}) (f(T_j, \tilde{Y}_{s+1-j}^+) - g(T_j, \tilde{z}_1(T_j))), \quad i = 1, \dots, s,$$

$$\tilde{y}_1^+ = \tilde{z}_1(t_1) - h \sum_{j=1}^s b_{s+1-j} (f(T_j, \tilde{Y}_{s+1-j}^+) - g(T_j, \tilde{z}_1(T_j)))$$

$$= \tilde{z}_1(t_1) - h \sum_{j=1}^s b_j (f(T_j, \tilde{Y}_{s+1-j}^+) - g(T_j, \tilde{z}_1(T_j))),$$

$$\tilde{y}_1 - \tilde{y}_0^- = y_0 - \tilde{y}_1^+.$$

It can be checked that the solution of the adjoint SRKLM equations is given by

$$\tilde{Y}_{s+1-i}^+ = Y_i^-, \quad i = 1, \dots, s, \quad \tilde{y}_1^+ = y_0^-,$$

$$\tilde{Y}_{s+1-i}^- = Y_i^+, \quad i = 1, \dots, s, \quad \tilde{y}_0^- = y_1^+, \quad \tilde{y}_1 = y_1,$$

where $Y_i^-, Y_i^+, y_0^-, y_1^+, y_1$ is the SRKLM solution of (22), and $\tilde{z}_1(t) = z_1(t)$. Since the adjoint method satisfies $\tilde{y}_1 = y_1$, the SRKLM method with symmetric RK coefficients is therefore symmetric. \square

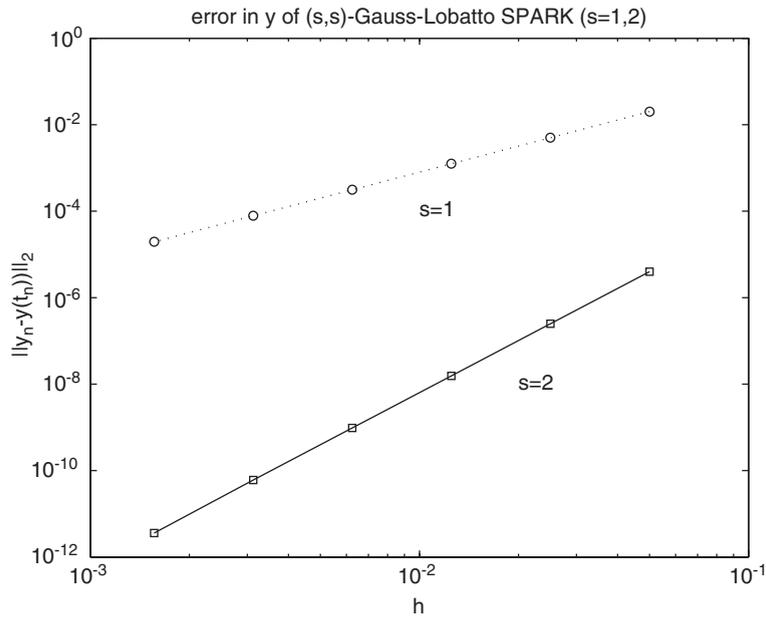


Fig. 1. Global error in y at $t_n = 1$ of (s, s) -Gauss–Lobatto SPARK methods ($s = 1, 2$) applied with various constant step sizes h to the test problem (23).

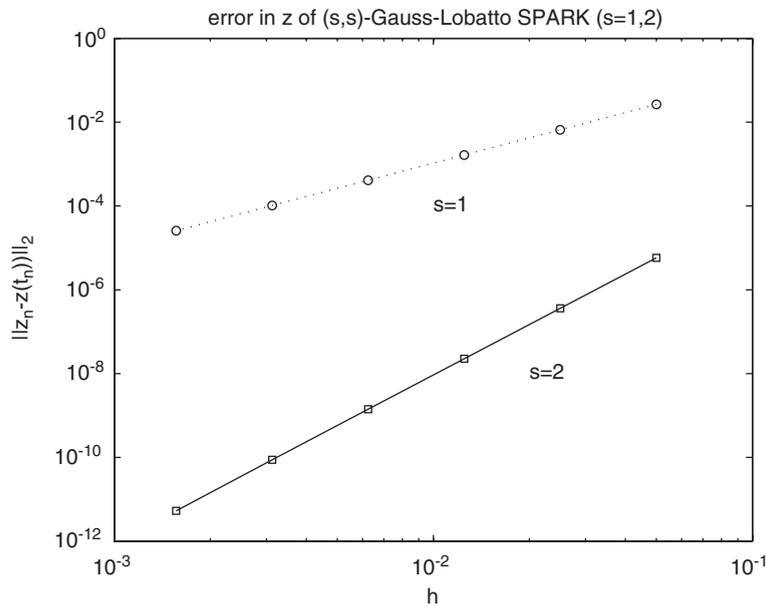


Fig. 2. Global error in z at $t_n = 1$ of (s, s) -Gauss–Lobatto SPARK methods ($s = 1, 2$) applied with various constant step sizes h to the test problem (23).

7. Numerical experiments

To illustrate Theorem 5, we have applied (s, s) -Gauss–Lobatto SPARK methods with constant step size h to the following system of index 3 DAEs

$$\begin{pmatrix} y'_1 \\ y'_2 \end{pmatrix} = \begin{pmatrix} 2z_1 \\ -z_2 \end{pmatrix}, \tag{23a}$$

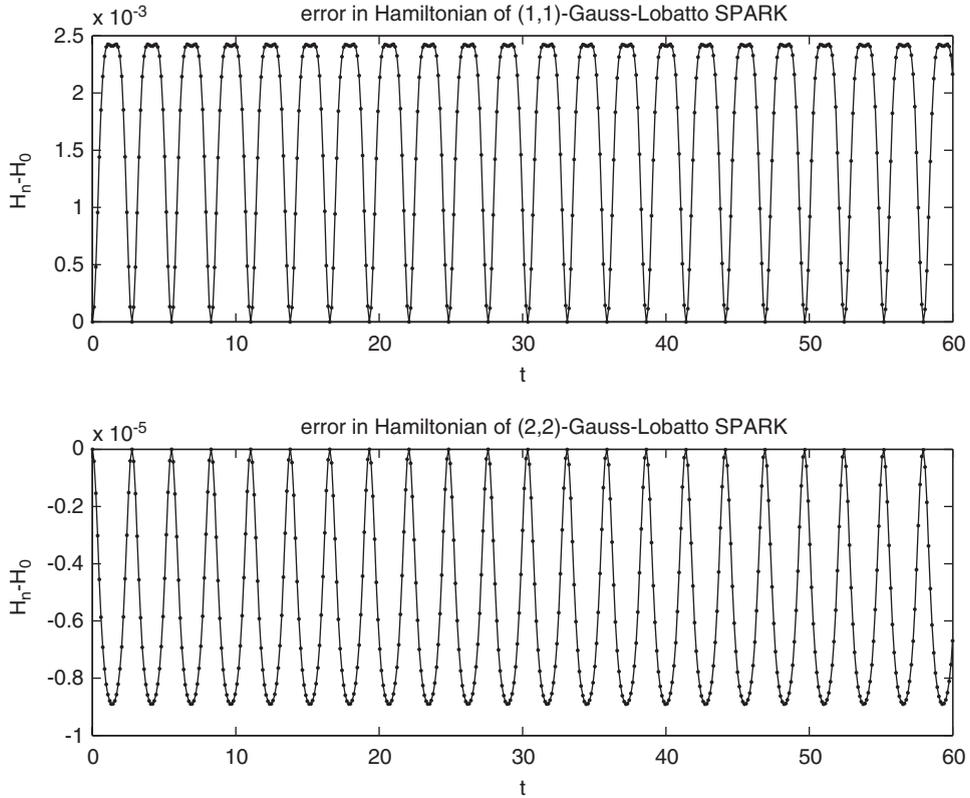


Fig. 3. Error in Hamiltonian of (s, s) -Gauss–Lobatto SPARK methods ($s = 1, 2$) applied with constant stepsize $h = 0.12$ to the test problem (24).

$$\begin{pmatrix} z_1' \\ z_2' \end{pmatrix} = \begin{pmatrix} 2y_1y_2z_1z_2 - y_1z_1z_2 \\ z_1 - y_1z_2^3 \end{pmatrix} + \begin{pmatrix} y_1y_2\psi_1^2 \\ -\sqrt{y_1}\psi_1 \end{pmatrix}, \tag{23b}$$

$$0 = y_1y_2^2 - 1. \tag{23c}$$

For the initial conditions $y_1(0) = y_2(0) = z_1(0) = z_2(0) = 1$ at $t_0 = 0$ the exact solution to this test problem is given by $y_1(t) = z_1(t) = e^{2t}$, $y_2(t) = z_2(t) = e^{-t}$, $\psi_1(t) = e^t$. We have plotted in Figs. 1 and 2 the global errors for the y - and z -components at $t_n = 1$ with respect to various constant stepsizes h . Logarithmic scales have been used so that a curve appears as a straight line of slope k whenever the leading term of the global error is of order k , i.e., when $\|y_n - y(t_n)\| = O(h^k)$. For the (s, s) -Gauss–Lobatto SPARK methods with $s = 1, 2$ of order $2s = 2, 4$ we observe straight lines of slope $2s = 2, 4$ thus confirming the orders of convergence predicted by Theorem 5.

As a second test problem, we consider the motion of a particle of mass m and electric charge e under the influence of an electric field $(0, 0, E)^T$ and a magnetic field $(0, 0, B)^T$ and restricted to a sphere of radius R [4, Problem 7.16]. This system can be described in term of Cartesian coordinates $(q_1, q_2, q_3)^T$ and generalized momenta $(p_1, p_2, p_3)^T$ with a nonseparable Hamiltonian

$$H = \frac{1}{2m}((p_1 + m\omega q_2)^2 + (p_2 - m\omega q_1)^2 + p_3^2) - eEq_3 \tag{24a}$$

with $\omega := eB/(2mc)$ and holonomic constraint

$$\sqrt{q_1^2 + q_2^2 + q_3^2} - R = 0. \tag{24b}$$

We choose the parameters

$$m = 1, \quad \omega = 1, \quad R = 1, \quad eE = 1,$$

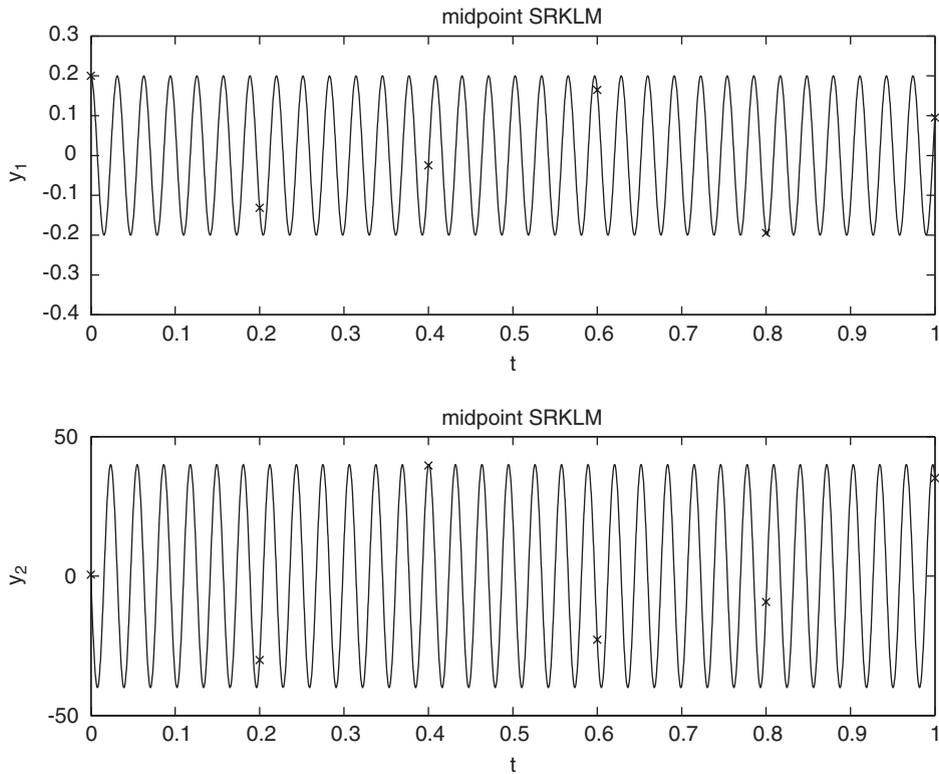


Fig. 4. Exact (–) and numerical (x) solution of harmonic oscillator problem for the midpoint SRKLM method.

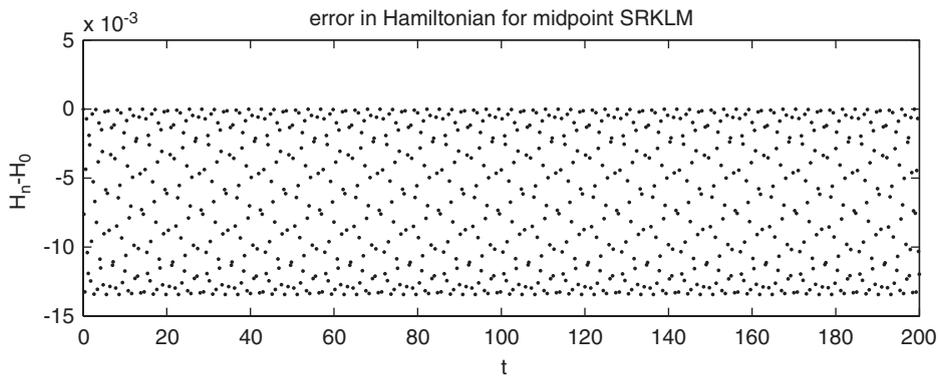


Fig. 5. Error in Hamiltonian of harmonic oscillator problem for the midpoint SRKLM method.

and initial conditions

$$q_1(0) = 0.2, \quad q_2(0) = 0.2, \quad q_3(0) = \sqrt{0.92}, \quad p_1(0) = 1, \quad p_2(0) = -1, \quad p_3(0) = 0.$$

In Fig. 3, we plot the Hamiltonian error of (s, s) -Gauss–Lobatto SPARK methods ($s = 1, 2$) applied with constant stepsize $h = 0.12$ to this system. As expected for a symplectic integrator, we observe that the Hamiltonian error remains bounded and small over long-time intervals.

To illustrate the applicability of SRKLM methods we consider the basic example of the linear harmonic oscillator

$$y_1' = y_2, \quad y_2' = -\omega^2 y_1,$$

with parameter $\omega = \rho\sqrt{1 + (\delta/\rho)^2}$ where $\rho = 200$ and $\delta = 1$, corresponding to a period $T = 2\pi/\omega \approx 0.03141553384 \dots$. We take the initial conditions $y_1(0) = 0.2$, $y_2(0) = 0.5$. The midpoint SRKLM method, based on the midpoint rule, is applied to this problem with stepsize $h = 0.2 \gg T$ using the local model

$$z_1' = z_2, \quad z_2' = -\rho^2 z_1.$$

In Fig. 4, we plot the exact solution and the numerical solution obtained at a few points on the time interval $[0, 1]$. We see that the numerical solution jumps over several periods without losing track of the phase of the solution. In Fig. 5, we plot the error in the Hamiltonian $H(y_1, y_2) = (\omega^2 y_1^2 + y_2^2)/2$ of the harmonic oscillator for the midpoint SRKLM method on the long-time interval $[0, 200]$. We observe that the Hamiltonian error remains bounded and small.

References

- [1] V.I. Arnold, *Mathematical Methods of Classical Mechanics*, Graduate Texts in Mathematics, second ed., vol. 60, Springer, New York, 1989.
- [2] W. Auzinger, R. Frank, H.J. Stetter, Vienna contributions to the development of RK-methods, *Appl. Numer. Math.* 22 (1996) 35–49.
- [3] J.C. Butcher, *Numerical Methods for Ordinary Differential Equations*, second ed., Wiley, Chichester, 2003.
- [4] P. Choquard, *Mécanique analytique*, Cahiers mathématiques de l'Ecole Polytechnique Fédérale de Lausanne, vol. 1, Presses Polytechniques et Universitaires Romandes, 1992.
- [5] E. Hairer, C. Lubich, M. Roche, *The Numerical Solution of Differential-Algebraic Systems by Runge–Kutta Methods*, Lecture Notes in Mathematics, vol. 1409, Springer, Berlin, 1989.
- [6] E. Hairer, C. Lubich, G. Wanner, *Geometric Numerical Integration*, Computational Mathematics, vol. 31, Springer, Berlin, 2002.
- [7] E.J. Haug, *Computer aided kinematics and dynamics of mechanical systems*, Basic Methods, vol. I, Allyn and Bacon, Boston, USA, 1989.
- [8] M. Hochbruck, A. Ostermann, Exponential Runge–Kutta methods for parabolic problems, *Appl. Numer. Math.* 53 (2005) 323–339.
- [9] L.O. Jay, *Runge–Kutta type methods for index three differential-algebraic equations with applications to Hamiltonian systems*, Ph.D. Thesis, Department of Mathematics, University of Geneva, Switzerland, 1994.
- [10] L.O. Jay, Symplectic partitioned Runge–Kutta methods for constrained Hamiltonian systems, *SIAM J. Numer. Anal.* 33 (1996) 368–387.
- [11] L.O. Jay, *Lagrangian integration with symplectic methods*, Technical Report, AHPCRC Preprint 97-009, AHPCRC, University of Minnesota, Minneapolis, USA, 1997.
- [12] L.O. Jay, Structure preservation for constrained dynamics with super partitioned additive Runge–Kutta methods, *SIAM J. Sci. Comput.* 20 (1998) 416–446.
- [13] L.O. Jay, Iterative solution of nonlinear equations for SPARK methods applied to DAEs, *Numer. Algorithms* 31 (2002) 171–191.
- [14] L.O. Jay, *Symplectic specialized partitioned additive Runge–Kutta methods for conservative systems with holonomic constraints*, Technical Report, Department of Mathematics, University of Iowa, USA, 2005, in progress.
- [15] S. Koikari, Rooted tree analysis of Runge–Kutta methods with exact treatment of linear terms, *J. Comput. Appl. Math.* 177 (2005) 427–523.
- [16] S. Krogstad, Generalized integrating factor methods for stiff PDEs, *J. Comp. Phys.* 203 (2005) 72–88.
- [17] J.E. Marsden, T.S. Ratiu, *Introduction to Mechanics and Symmetry*, Texts in Applied Mathematics, vol. 17, Springer, New York, 1994.
- [18] J.E. Marsden, M. West, Discrete mechanics and variational integrators, *Acta Numerica* (2001) 357–514.
- [19] A. Murua, *Partitioned Runge–Kutta methods for semi-explicit differential-algebraic systems of index 2*, Technical Report, EHU-KZAA-IKT-196, University of the Basque Country, 1996.
- [20] P.J. Rabier, W.C. Rheinboldt, *Nonholonomic Motion of Rigid Mechanical Systems from a DAE Viewpoint*, SIAM, Philadelphia, 2000.
- [21] H.J. Stetter, The defect correction principle and discretization methods, *Numer. Math.* 29 (1978) 425–443.
- [22] P.E. Zadunaisky, On the estimation of errors propagated in the numerical integration of ordinary differential equations, *Numer. Math.* 27 (1976) 21–39.