# The numerical influence of additional parameters of inertia representations for quaternion-based rigid body dynamics

**Xiaoming Xu**[1] · **Jiahui Luo**[1] · **Zhigang Wu**[1]

**Abstract** Different inertia representations can lead to different formulations of the differential-algebraic equations for the quaternion-based rigid body dynamics. In this paper, the inertia representations are classified into $\alpha$-type and $\gamma$-type, according to the additional parameters in the kinetic energy. These two types of representations and the corresponding parameters $\alpha$ and $\gamma$ are theoretically equivalent if the constraint $q^T q = 1$ is satisfied exactly. Nevertheless, the error estimation demonstrates that they can present entirely different numerical features in simulation and suggests that the parameter $\gamma$ can be used to optimize the numerical performance of the integrations in simulation. To further verify the numerical difference between the inertia representations of $\alpha$-type and $\gamma$-type, the corresponding modified Hamilton's equations are discretized by the IMS (implicit midpoint scheme), EMS (energy–momentum preserving scheme) and Gauss–Lobatto SPARK methods. Numerical performance for the examples of the spinning symmetrical top is shown to result from the comprehensive effect of the discretization schemes including the distribution of discretized points and the convergence order, the inertia representations and their combinations. Numerical results further suggest that the integrations of $\gamma$-type are superior to those of $\alpha$-type and the optimized values of $\gamma$ can be used to achieve better numerical accuracy, convergence speed and stability.

**Keywords** Unit quaternion · Rigid body dynamics · Mass matrix · Singularity · High accuracy · Conserving integrations

## 1 Introduction

The kinematics and dynamics of rigid bodies constitute an important part of the simulation of multibody system. Many kinds of coordinates can be used to represent the rotational motion of a rigid body, such as Euler angles and unit quaternions [1]. As a non-minimal representation, the unit quaternion has found new attraction in recent years because of its

---

✉ X. Xu
xxm020201@163.com

1 School of Aeronautics and Astronautics, Sun Yat-sen University, Guangzhou 510006, P.R. China

simplicity in the mathematical formulation and the ability to avoid the singularity which may occur when using Euler angles.

A unit quaternion describes the three-dimensional rotational motion using four real parameters, with a unit length constraint, which essentially extends the motion equations by a holonomic constraint and yields a set of differential-algebraic equations (DAEs) of Index 3, instead of ordinary differential equations (ODEs). In 1980s, the unit quaternion was investigated in detail by Nikravesh and Haug et al. [2–5] for the dynamic analysis of a three-dimensional constrained mechanical system. Later, the quaternion-based rigid body dynamics was developed rapidly owing to many theoretical researches and applications about unit quaternions [6–22]. Some of these researches focused on various equivalent formulations for rotational motion in the Lagrangian or Hamiltonian framework [6–14], and the others concentrated on the numerical integrations in terms of unit quaternions for the simulation of rigid body dynamics, especially the conserving integrations [11, 15–20].

In this previous work, an essential process is to derive the dynamic description of rigid body rotation, which leads to different inertia representations in terms of unit quaternions, embodied by the mass matrix in the kinetic energy. In rigid body dynamics, the quaternion-based kinetic energy is commonly derived from the standard quadratic form of the angular velocity, which generally results in a singular mass matrix. To avoid the singularity, Betsch and Siebert [11] proposed an augmented mass matrix when deriving the conserving numerical integration in terms of unit quaternions for rigid body dynamics. Besides, similar approaches [12–16] were proposed to describe the motion equations, which also lead to a non-singular mass matrix. In this work, an additional parameter denoted as $\alpha$ in this paper, is generally added in the original formulation of kinetic energy. Nielsen and Krenk [16] developed the quaternion-based momentum scheme based on the non-singular mass matrix, and suggested that the additional parameter served as a multiplier on the kinematic constraint, and better convergence characteristics could be achieved by choosing this parameter somewhat larger than the inertial moments in numerical simulation. Recently, a modified inertia representation [23] was proposed for the quaternion-based rigid body dynamics, in which a new parameter, denoted $\gamma$ in this paper, was introduced in the kinetic energy. The parameters $\alpha$ and $\gamma$ are mathematically equivalent if the unit length constraint is satisfied exactly. Nevertheless, the two parameters are different in discretization, and numerical results demonstrate that the parameter $\gamma$ is superior to $\alpha$ because it can be used to improve the numerical accuracy of integration in simulation.

In this paper, the numerical performance of different inertia representations is investigated at great length for the quaternion-based rigid body dynamics. In Sect. 2, we first derive the Hamilton equations for the rigid body dynamics in terms of unit quaternions, and classify these equations into two types: the first is based on the augmented formulation of kinetic energy, denoted as $\alpha$-type, and the other is according to the modified inertia representation, denoted as $\gamma$-type. In Sect. 3, error estimation demonstrates that the Hamiltonians in $\alpha$-type and $\gamma$-type are essentially different from the point of view of discretization, and it suggests that the parameter $\gamma$ is expected to have great influence on the discretization error of kinetic energy. Based on the error estimation, two predetermined values $\gamma_m$ and $\gamma_h$ are recognized to reduce numerical errors for the integrations of $\gamma$-type.

Section 4 develops several kinds of integrations, especially the specialized partitioned additive Runge–Kutta methods (SPARK) [24] for rigid body dynamics in terms of unit quaternions, according to the Hamilton equations of $\alpha$-type and $\gamma$-type. Section 5 systematically investigates the numerical accuracy and stability for these integrations. Numerical results demonstrate that most of the integrations of $\gamma$-type assigned with $\gamma = \gamma_m$ or $\gamma_h$ can present impressively better numerical accuracy than those of $\alpha$-type, especially for the

Gauss SPARK methods. All the integrations of $\gamma$-type can present better convergence speed and stronger robustness for the newton iteration than those of $\alpha$-type. A large amount of numerical comparisons reveal that the numerical performance of the integrations is extensively influenced by the inertia representations, discretization schemes and their combinations. The findings of this paper suggest that numerical integrations in terms of unit quaternions should be constructed with an appropriate inertia representation in order to obtain better numerical accuracy, convergence speed and stability.

## 2 Kinematic description

The unit quaternion can be considered as a four-parameter vector [16]:

$$\boldsymbol{q}^T = \begin{bmatrix} q_0 & q_1 & q_2 & q_3 \end{bmatrix} \tag{1}$$

with an algebraic constraint,

$$\boldsymbol{q}^T \boldsymbol{q} - 1 = 0. \tag{2}$$

An important application of unit quaternions is that the three-dimensional orthogonal matrix can be expressed in terms of unit quaternions:

$$\boldsymbol{R}(\boldsymbol{q}) = \begin{bmatrix} q_0^2 + q_1^2 - q_2^2 - q_3^2 & 2(q_1 q_2 - q_0 q_3) & 2(q_1 q_3 + q_0 q_2) \\ 2(q_1 q_2 + q_0 q_3) & q_0^2 - q_1^2 + q_2^2 - q_3^2 & 2(q_2 q_3 - q_0 q_1) \\ 2(q_1 q_3 - q_0 q_2) & 2(q_2 q_3 + q_0 q_1) & q_0^2 - q_1^2 - q_2^2 + q_3^2 \end{bmatrix}, \tag{3}$$

where $\boldsymbol{R}$ satisfies the orthogonal condition $\boldsymbol{R}^T \boldsymbol{R} = \boldsymbol{R} \boldsymbol{R}^T = \boldsymbol{I}_3$ and $\boldsymbol{I}_3$ is the 3-dimensional identity matrix. In rigid body dynamics, $\boldsymbol{R}$ is named the rotation matrix; it maps from the space-fixed coordinates $\boldsymbol{X}$ to the body-fixed coordinate $\boldsymbol{x}$, i.e.,

$$\boldsymbol{x} = \boldsymbol{R} \boldsymbol{X}. \tag{4}$$

Differentiating (4), we can derive the following motion equations:

$$\dot{\boldsymbol{R}} = \boldsymbol{R} \hat{\boldsymbol{\Omega}}, \tag{5}$$

where the superscript 'dot' denotes the derivation of variables with respect to time, and the matrix $\hat{\boldsymbol{\Omega}}$ is in the form of

$$\hat{\boldsymbol{\Omega}} = \begin{bmatrix} 0 & -\Omega_3 & \Omega_2 \\ \Omega_3 & 0 & -\Omega_1 \\ -\Omega_2 & \Omega_1 & 0 \end{bmatrix}, \tag{6}$$

where $\boldsymbol{\Omega} = [\Omega_1 \ \Omega_2 \ \Omega_3]^T$ is the angular velocity vector. Substituting (3) into (5) leads to the quaternion-based motion equations [11]

$$\boldsymbol{\Omega} = 2 \boldsymbol{L}(\boldsymbol{q}) \dot{\boldsymbol{q}}, \tag{7}$$

where

$$\boldsymbol{L}(\boldsymbol{q}) = \begin{bmatrix} -q_1 & q_0 & q_3 & -q_2 \\ -q_2 & -q_3 & q_0 & q_1 \\ -q_3 & q_2 & -q_1 & q_0 \end{bmatrix}. \tag{8}$$

Noting that $q^T q = 1$, a direct calculation reveals that

$$L(q)q = 0, \qquad L(q)\dot{q} = -L(\dot{q})q, \qquad L(\dot{q})\dot{q} = 0, \tag{9}$$

$$L(q)L^T(q) = I_3, \tag{10}$$

$$L(q)^T L(q) = I_4 - qq^T, \tag{11}$$

where $I_n$ denotes the $n$-dimensional identity matrix. These identity relations are useful in later discussion.

## 3 The rotational motion equations in the Hamiltonian framework

In the Hamiltonian description of conservative system, the motion equations can be expressed in a unified form [25]:

$$\begin{aligned}
\dot{q} &= H_p(q, p), \\
\dot{p} &= -H_q(q, p) - g_q(q)\lambda, \\
0 &= g(q), \\
0 &= g_q(q)^T H_p(q, p),
\end{aligned} \tag{12}$$

where the third and fourth equations denote the constraint conditions. The Hamiltonian

$$H = T(q, p) + V(q) \tag{13}$$

consists of the kinetic energy $T$ and the potential energy $V$, and the abbreviations are defined as $H_p = \partial H / \partial p^T$, $H_q = \partial H / \partial q^T$ and $g_q(q) = \partial g / \partial q^T$. $\lambda$ is the Lagrange multiplier which preserves the path of quaternions satisfying $g(q) = 0$, and $\mathbf{p}$ denotes the generalized momentum. For the quaternion-based rigid body rotational motion, we have $g(q) = q^T q - 1$ and $g_q(q) = 2q$.

Suppose that the body-fixed coordinate axes are aligned along the principal axes of the inertia of the rigid body, and $J = \operatorname{diag}(I_1, I_2, I_3)$ denotes the inertia matrix whose principal elements are three principal moments of inertia. Combined with (7), the kinetic energy of the rigid body rotation can be expressed as

$$T = \boldsymbol{\Omega} J \boldsymbol{\Omega} / 2 = 2\dot{q}^T L(q)^T J L(q)\dot{q}. \tag{14}$$

Based on a Legendre transformation, the generalized momentum can be derived as

$$p = \partial T / \partial \dot{q}^T = M\dot{q}, \tag{15}$$

where

$$M = 4L(q)^T J L(q) \tag{16}$$

is defined as the mass matrix. Multiplying (15) with $q^T$ from the left and recalling $L(q)q = 0$, we have

$$q^T p = p^T q = 0, \tag{17}$$

which is equivalent to the constraint $0 = g_q(q)H_p(q, p)$ for unit quaternions. Multiplying (15) with $L(q)$ from the left and recalling $L(q)L(q)^T = I_3$, we have

$$L(q)p = 4L(q)L(q)^T JL(q)\dot{q} = 4JL(q)\dot{q}. \tag{18}$$

Noting that $\Omega = 2L(q)\dot{q}$ and $l = J\Omega$ where $l$ denotes the local angular momentum, it can be derived that

$$l = J\Omega = 2JL(q)\dot{q} = \frac{1}{2}L(q)p = -\frac{1}{2}L(p)q, \tag{19}$$

where the local angular momentum can be expressed as a function of $q$ and $p$. Substitute $\Omega = J^{-1}l$ into (14) and we can derive the kinetic energy in the Hamiltonian framework:

$$T = l^T J^{-1}l/2 = \frac{1}{8}p^T L(q)^T J^{-1}L(q)p. \tag{20}$$

Substituting (20) into the first two equations in (12), we can derive the differential parts of the Hamilton's equations in the following form:

$$\dot{q} = \frac{1}{4}L(q)^T J^{-1}L(q)p,$$
$$\dot{p} = -\frac{1}{4}L(p)^T J^{-1}L(p)q - \partial V/\partial q - 2\lambda q, \tag{21}$$

where the multiplier $\lambda$ is used to preserve the constraint $q^T q - 1 = 0$.

### 3.1 The inertia representations of $\alpha$-type and $\gamma$-type

The numerical application of the Hamilton's equations (21) is complicated because of the singularity of the matrix $\frac{1}{4}L(q)^T J^{-1}L(q)$. To avoid the singularity, researchers [11–16] introduced an additional term $\alpha^{-1}(p^T q)^2$ into the original formulation of kinetic energy. Generally, the augmented formulation of the kinetic energy can be expressed as

$$T = \frac{1}{8}p^T L(q)^T J^{-1}L(q)p + \frac{1}{8}\alpha^{-1}(p^T q)^2. \tag{22}$$

Because of $p^T q = 0$, the term $\frac{1}{8}\alpha^{-1}(p^T q)^2$ has no influence on the value of the kinetic energy. Substituting (22) into (12), we can derive the Hamilton's equations of $\alpha$-type:

$$\dot{q} = \frac{1}{4}L(q)^T J^{-1}L(q)p + \frac{1}{4}\alpha^{-1}qq^T p,$$
$$\dot{p} = -\frac{1}{4}L(p)^T J^{-1}L(p)q - \frac{1}{4}\alpha^{-1}(p^T q)p - \partial V/\partial q - 2\lambda q, \tag{23}$$

where the multiplier $\lambda$ is used to preserve the constraint $q^T q - 1 = 0$. The corresponding matrix $\frac{1}{4}L(q)^T J^{-1}L(q) + \frac{1}{4}\alpha^{-1}qq^T$ in (23) is non-singular if only $\alpha^{-1} \neq 0$. Numerical results demonstrate that the conserving integrations of $\alpha$-type can present good long-time behavior in simulations [11, 16].

Noting that $l = -\frac{1}{2}L(p)q$, $L(q)^T L(q) + qq^T = I_4$ and $q^T p = 0$, we have

$$l^T l = \frac{1}{4}p^T p. \tag{24}$$

Let us reformulate the kinetic energy in the following form:

$$T = \boldsymbol{l}^T \boldsymbol{J}^{-1} \boldsymbol{l}/2 = \gamma^{-1} \boldsymbol{l}^T \boldsymbol{l}/2 + \boldsymbol{l}^T \boldsymbol{J}_\gamma^* \boldsymbol{l}/2, \tag{25}$$

where $\gamma$ is an arbitrary constant and $\gamma \neq 0$, and $\boldsymbol{J}_\gamma^* = \boldsymbol{J}^{-1} - \gamma^{-1} \boldsymbol{I}_3$. Substituting (19) and (24) into (25) yields a new formulation:

$$T = \frac{1}{8} \gamma^{-1} \boldsymbol{p}^T \boldsymbol{p} + \frac{1}{8} \boldsymbol{p}^T \boldsymbol{L}(\boldsymbol{q})^T \boldsymbol{J}_\gamma^* \boldsymbol{L}(\boldsymbol{q}) \boldsymbol{p}, \tag{26}$$

where the kinetic energy is split into two components. The first is a square term of the quaternion momentum, whose magnitude is adjusted by the parameter $\gamma$. The other is a quadratic form in the derivatives with quadratic coefficients in the quaternion parameters. Equation (26) is the modified inertia representation presented by Xu and Zhong [23]. Substituting (26) into (12) yields the Hamilton's equations of $\gamma$-type:

$$\begin{aligned}
\dot{\boldsymbol{q}} &= \frac{1}{4} \gamma^{-1} \boldsymbol{p} + \frac{1}{4} \boldsymbol{L}(\boldsymbol{q})^T \boldsymbol{J}_\gamma^* \boldsymbol{L}(\boldsymbol{q}) \boldsymbol{p}, \\
\dot{\boldsymbol{p}} &= -\frac{1}{4} \boldsymbol{L}(\boldsymbol{p})^T \boldsymbol{J}_\gamma^* \boldsymbol{L}(\boldsymbol{p}) \boldsymbol{q} - \partial V/\partial \boldsymbol{q} - 2\lambda \boldsymbol{q},
\end{aligned} \tag{27}$$

where the multiplier $\lambda$ is used to preserve the constraint $\boldsymbol{q}^T \boldsymbol{q} - 1 = 0$.

The inertia representations of $\alpha$ and $\gamma$-type are mathematically equivalent under certain condition. To be specific, it can be derived that

$$\begin{aligned}
T &= \frac{1}{8} \gamma^{-1} \boldsymbol{p}^T \boldsymbol{p} + \frac{1}{8} \boldsymbol{p}^T \boldsymbol{L}(\boldsymbol{q})^T \left( \boldsymbol{J}^{-1} - \gamma^{-1} \boldsymbol{I}_3 \right) \boldsymbol{L}(\boldsymbol{q}) \boldsymbol{p} \\
&= \frac{1}{8} \alpha^{-1} \boldsymbol{p}^T \left[ \boldsymbol{L}(\boldsymbol{q})^T \boldsymbol{L}(\boldsymbol{q}) + \boldsymbol{q} \boldsymbol{q}^T \right] \boldsymbol{p} + \frac{1}{8} \boldsymbol{p}^T \boldsymbol{L}(\boldsymbol{q})^T \left( \boldsymbol{J}^{-1} - \alpha^{-1} \boldsymbol{I}_3 \right) \boldsymbol{L}(\boldsymbol{q}) \boldsymbol{p} \\
&= \frac{1}{8} \boldsymbol{p}^T \boldsymbol{L}(\boldsymbol{q})^T \boldsymbol{J}^{-1} \boldsymbol{L}(\boldsymbol{q}) \boldsymbol{p} + \frac{1}{8} \alpha^{-1} \left( \boldsymbol{p}^T \boldsymbol{q} \right)^2
\end{aligned} \tag{28}$$

by considering $\gamma = \alpha$ and $\boldsymbol{I}_4 = \boldsymbol{L}(\boldsymbol{q})^T \boldsymbol{L}(\boldsymbol{q}) + \boldsymbol{q} \boldsymbol{q}^T$. The mathematical equivalence implies that the parameters $\alpha$ and $\gamma$ theoretically have no influence on the value of the kinetic energy if the constraint $\boldsymbol{q}^T \boldsymbol{q} - 1 = 0$ is satisfied exactly. In the following, the two representations and the corresponding parameters $\alpha$ and $\gamma$ will be systematically studied both in theoretical analysis and numerical simulation, and the numerical performance of the two parameters $\alpha$ and $\gamma$ would be investigated intensively with several different discretized schemes and convergence orders.

## 3.2 The discretization error estimation

Suppose that $\boldsymbol{q} = \boldsymbol{q}(t)$ and $\boldsymbol{p} = \boldsymbol{p}(t)$ denote the real solutions of rotational motion of a single body, and that $\boldsymbol{q}^+ = \boldsymbol{q}^+(t)$ and $\boldsymbol{p}^+ = \boldsymbol{p}^+(t)$ are their interpolation approximation solutions. We can define the real angular momentum and its approximation in the following form:

$$\boldsymbol{l} = \boldsymbol{l}(\boldsymbol{q}, \boldsymbol{p}), \qquad \boldsymbol{l}^+ = \boldsymbol{l}(\boldsymbol{q}^+, \boldsymbol{p}^+). \tag{29}$$

Then the local discretization error of $\boldsymbol{q}$ and $\boldsymbol{l}$ can be defined, respectively, as

$$\delta \boldsymbol{q} = \boldsymbol{q}^+(t) - \boldsymbol{q}(t), \qquad \delta \boldsymbol{p} = \boldsymbol{p}^+(t) - \boldsymbol{p}(t) \tag{30}$$

and

$$\delta l = l^+ - l = \frac{1}{2}\delta\big(L(q)p\big) = \frac{1}{2}(\delta L)p + \frac{1}{2}L\delta p + \frac{1}{2}\delta L\delta p, \tag{31}$$

where $\delta L = L(q^+) - L(q)$. The constraint $q^T q = 1$ is satisfied approximately as

$$q^{+T}q^+ = 1 + 2q^T\delta q + \delta q^T\delta q. \tag{32}$$

Considering the inertia representation of $\gamma$-type, the discretization error of kinetic energy can be formally expressed as

$$\delta T_\gamma = T\big(q^+, p^+, \gamma\big) - T(q, p, \gamma). \tag{33}$$

Consider Eq. (26) and substitute (29) and (31) into (33), and then the discretization error of kinetic energy of $\gamma$-type can be expressed as

$$\delta T_\gamma = k_s\gamma^{-1} + \delta T_0, \tag{34}$$

where the slope is defined as

$$k_s = \frac{1}{8}\big(p^{+T}p^+ - 4l^{+T}l^+\big), \tag{35}$$

and the intercept is of the form of

$$\delta T_0 = \delta l^T J^{-1}l + \delta l^T J^{-1}\delta l/2. \tag{36}$$

Note that $\delta T_0$ is the discretization error of the original formulation of kinetic energy presented by (20). Similarly, we can derive the discretization error of kinetic energy of $\alpha$-type, which is

$$\delta T_\alpha = T\big(q^+, p^+, \alpha\big) - T(q, p, \alpha)$$
$$= k_a\alpha^{-1} + \delta T_0 \tag{37}$$

with the slope

$$k_a = 8^{-1}\big(p^{+T}q^+\big)^2. \tag{38}$$

Equations (34) and (37) reveal the linear relation between the energy error and $1/\gamma$ or $1/\alpha$. This relation can be exploited to reduce the numerical errors of integrations by choosing the appropriate $\gamma$ or $\alpha$ to cancel out the intercept term $\delta T_0$, as done in [23]. For this to work, the slope should ideally be of the same convergence order as $\delta T_0$. However, this is not always the case for the inertia representation of $\alpha$-type. More specifically, substituting $q^+ = q + \delta q$ and $p^+ = p + \delta p$ into (38) and recalling $q^T p = 0$ yield

$$k_a = 8^{-1}\big(q^T\delta p + p^T\delta q + \delta p^T\delta q\big)^2. \tag{39}$$

When the error terms $\delta q$ and $\delta p$ are of the same convergence order (i.e., $O(\delta q) = O(\delta p)$), Eq. (39) means $k_a$ is generally a second-order small quantity with respect to $\delta q$ and $\delta p$. In this case, the slope $k_a$ vanishes more quickly than the term $\delta T_0$ when the discretization errors $\delta q$ and $\delta p$ tend to be zero. Accordingly, the corresponding parameter $\alpha$ has nearly no influence on the discretization errors.

On the other hand, for the inertia representation of $\gamma$-type, the slope $k_s$ can be expanded as

$$k_s = \left(\frac{1}{4}\delta p^T p - \delta l^T l\right) + \frac{1}{2}\left(\frac{1}{4}\delta p^T \delta p - \delta l^T \delta l\right), \tag{40}$$

where the second term on the right side is a second-order small quantity. Ignoring the second-order small quantity and considering $L(q)^T L(q) + qq^T = I_4$ and $p^T q = 0$, we can obtain

$$k_s \approx -\frac{1}{4}p^T(\delta L)^T L p, \tag{41}$$

where the error term $\delta L$ is of the same order with the discretization error $\delta q$. One can expect

$$k_s \propto O\big(\|\delta q\|\big) \tag{42}$$

for any $q$ and $p$, where $\|\delta q\| = \sqrt{\delta q^T \delta q}$ denotes the 2-norm of $\delta q$. This implies that the slope $k_s$ is of the same order as $\delta T_0$, and the discretization error of kinetic energy will be greatly influenced by the parameter $\gamma$.

We note that it is not always straightforward to verify (42) when we investigate the integration schemes. Instead, a more practical condition can be used. Substituting $l^+ = -\frac{1}{2}L(p^+)q^+$ into (35) and considering that$(p^{+T}q^+)^2 = O[(\delta q)^2, (\delta p)^2]$, the slope $k_s$ can be reformulated as

$$k_s + O\big[(\delta q)^2, (\delta p)^2\big] = \frac{1}{8}p^{+T}\big(I_4 - L\big(q^+\big)^T L\big(q^+\big) + q^+ q^{+T}\big)p^+. \tag{43}$$

This suggests that the interpolation function should satisfy a necessary but not sufficient condition:

$$L\big(q^+\big)^T L\big(q^+\big) + q^+ q^{+T} = I_4 + O(\delta q) \tag{44}$$

for $q^+ = q^+(t)$, to preserve $k_s \propto O(\delta q)$.

Condition (44) means that, for $k_s$ to be exploitable, we actually want the constraint on the unit quaternion to be satisfied only to the same order approximately, but not to higher order or exactly. As we will show in the following, Condition (44) can easily be verified, especially for the Gauss SPARK methods.

### 3.3 The optimal value of parameter $\gamma$

Suppose that the integrations satisfy (42) at every discretized points. The error estimation in (34) suggests a linear relationship between $\delta T$ and $\gamma^{-1}$. Let $\delta T = 0$ in (34), and it gives the optimal value of the parameter $\gamma$ in the following form:

$$\gamma_{\text{opt}} = -\frac{k_s}{\delta T_0} \approx \frac{1}{4}\frac{p^T(\delta L)^T L p}{\delta l^T J^{-1} l}, \tag{45}$$

where higher-order terms are neglected. This is a reasonable value for the parameter $\gamma$ to obtain a smaller discretization error and is expected to improve the numerical accuracy for the integrations of $\gamma$-type. Unfortunately, the quantities $k_s$ and $\delta T_0$ are unknown during the simulation. Substituting $l = -\frac{1}{2}Lp$ and $\delta l \approx \frac{1}{2}(\delta L)p + \frac{1}{2}L\delta p$ into (45), yields

$$\gamma_{\text{opt}} \approx \frac{\sum_{i=1}^3 w_i}{\sum_{i=1}^3 w_i I_i^{-1}}, \tag{46}$$

where $\boldsymbol{l} = [l_1, l_2, l_3]^T$, $w_i = l_i \delta l_i$ $(i = 1, 2, 3)$ and the error term $\delta \boldsymbol{p}$ is neglected. Let $w_1 = w_2 = w_3 = 1$, and it gives the harmonic average of the three principal moments:

$$\gamma_h = \frac{3}{I_1^{-1} + I_2^{-1} + I_3^{-1}}. \tag{47}$$

Let $w_1 = I_1$, $w_2 = I_2$ and $w_3 = I_3$, and it gives the arithmetic average of the three principal moments:

$$\gamma_m = \frac{I_1 + I_2 + I_3}{3}. \tag{48}$$

The values $\gamma_m$ and $\gamma_h$ are recommended as two reasonable values of $\gamma$ for integrations of $\gamma$-type.

In addition, the two optimal values $\gamma_m$ and $\gamma_h$ present very similar numerical performance when the three principal moments (i.e., $I_1$, $I_2$ and $I_3$) have close values. Nevertheless, the values $\gamma_m$ and $\gamma_h$ become very different if there are large differences among the three principal moments. For instance, suppose that $I_3 \to 0$, and that $I_1$ and $I_2$ remain unchanged. We thus obtain

$$\gamma_h \to 0, \qquad \gamma_m \to \frac{I_1 + I_2}{3}. \tag{49}$$

Because $\delta T_\gamma = k_s \gamma^{-1} + \delta T_0 \to \infty$ when $\gamma \to 0$, Eq. (49) implies a potential risk of serious accuracy loss in simulation for the integrations of $\gamma$-type with $\gamma = \gamma_h$. On the other hand, it can be derived that $T = (I_1 \Omega_1^2 + I_2 \Omega_2^2)/2 = (I_1^{-1} l_1^2 + I_2^{-1} l_2^2)/2$ by setting $I_3 = 0$. It means that $l_3$ has no contribution to the kinetic energy nor to the discretization error if $I_3 = 0$. Consequently, substituting $[w_1, w_2, w_3] = [1, 1, 0]$ and $[w_1, w_2, w_3] = [I_1, I_2, 0]$ into (46) and we can derive the correct optimal values of the parameter $\gamma$ for the case $I_3 \to 0$ as follows:

$$\tilde{\gamma}_h = \gamma_h(I_3 \to 0) = \frac{2}{I_1^{-1} + I_2^{-1}}, \qquad \tilde{\gamma}_m = \gamma_m(I_3 \to 0) = \frac{I_1 + I_2}{2}. \tag{50}$$

Compared to (49), it reveals that the value $\gamma_h$ in (47) tends to magnify the numerical influence of the smallest value of the three principal moments. These subtleties should be considered in simulation.

## 4 Numerical integrations of $\gamma$-type and $\alpha$-type

There is no restriction on the convergence order of the numerical integrations in the error analysis, which implies that the inertia representation of $\gamma$-type can be used to improve the numerical accuracy regardless of the convergence order of the algorithm. In the following, numerical integrations with different discretization schemes are developed and compared between the $\alpha$-type and the $\gamma$-type, in order to demonstrate the widespread applicability of the inertia representation in numerical simulation.

### 4.1 Numerical integrations of order 2

Let the phase space coordinates $(\boldsymbol{q}_{k-1}, \boldsymbol{p}_{k-1})$ at $t_{k-1}$ along with the step-size $\Delta t = t_k - t_{k-1}$ be given. Define the mean values as

$$\bar{\boldsymbol{q}} = \frac{1}{2}(\boldsymbol{q}_{k-1} + \boldsymbol{q}_k), \qquad \bar{\boldsymbol{p}} = \frac{1}{2}(\boldsymbol{p}_{k-1} + \boldsymbol{p}_k) \tag{51}$$

and the increments as

$$\Delta q = q_k - q_{k-1}, \qquad \Delta p = p_k - p_{k-1} \tag{52}$$

for time interval $t \in [t_{k-1}, t_k]$.

**Implicit midpoint scheme (IMS)**  *Approximating the Hamilton's equation in* (12) *by the midpoint rule gives the following scheme:*

$$\Delta q / \Delta t = H_p(\bar{q}, \bar{p}), \qquad 0 = g(q_k),$$
$$\Delta p / \Delta t = -H_q(\bar{q}, \bar{p}) - 2\lambda \bar{q}, \tag{53}$$

*where the unknowns* $q_k$, $p_k$ *and* $\lambda$ *can be solved for the given* $q_{k-1}$ *and* $p_{k-1}$. *The schemes are abbreviated as IMS-$\alpha$ and IMS-$\gamma$ for the integrations of $\alpha$-type and $\gamma$-type, respectively.*

The second integration we are interested in is the energy–momentum preserving integration developed in Refs. [11, 16, 26]. To this end, we consider the increment of the Hamiltonian:

$$\Delta H = \Delta T + \Delta V, \tag{54}$$

where the increments $\Delta T = T(q_k, p_k) - T(q_{k-1}, p_{k-1})$, $\Delta V = V(q_k) - T(q_{k-1})$ and $\Delta g = g(q_k) - g(q_{k-1}) = 2\Delta q^T \bar{q}$. After that we can define the increment of the augmented Hamiltonian as follows:

$$\Delta H_\lambda = \Delta H + \lambda \Delta g. \tag{55}$$

The energy–momentum preserving integration is constructed by setting $\Delta H_\lambda = 0$. To this end, we first express $\Delta V$ in terms of its finite derivative $\partial V_* / \partial q^T$, defined by [27]

$$\Delta V = \Delta q^T \frac{\partial V_*}{\partial q^T}. \tag{56}$$

Substituting (22) and (56) into (54) yields

$$\Delta H = \Delta p^T \left( L(\bar{q}) J^{-1} \overline{L(q)^T p} + \alpha^{-1} \bar{q} \overline{q^T p} \right)$$
$$+ \Delta q^T \left( L(\bar{p}) J^{-1} \overline{L(p)^T q} + \alpha^{-1} \bar{p} \overline{p^T q} + \partial V_* / \partial q^T \right), \tag{57}$$

where the overbar denotes the arithmetic mean, $\overline{(\#)} = \frac{1}{2}[(\#)_{k-1} + (\#)_k]$. According to (57), we can define the finite derivatives of $H$ as

$$\bar{\nabla}_q H = L(\bar{q}) J^{-1} \overline{L(q)^T p} + \alpha^{-1} \bar{q} \overline{q^T p},$$
$$\bar{\nabla}_p H = L(\bar{p}) J^{-1} \overline{L(p)^T q} + \alpha^{-1} \bar{p} \overline{p^T q} + \partial V_* / \partial q \tag{58}$$

for the inertia representation of $\alpha$-type. Similarly, according to (26), we can define the finite derivatives of $H$ as

$$\bar{\nabla}_q H = \gamma^{-1} \bar{p} + L(\bar{q}) J_\gamma^* \overline{L(q)^T p},$$
$$\bar{\nabla}_p H = L(\bar{p}) J_\gamma^* \overline{L(p)^T q} + \partial V_* / \partial q, \tag{59}$$

for the inertia representation of $\gamma$-type. Based on (58) or (59), Eq. (55) can be expressed as

$$\Delta H_\lambda = \Delta p^T \bar{\nabla}_q H + \Delta q^T (\bar{\nabla}_p H + 2\lambda \bar{q}). \tag{60}$$

**Energy–momentum preserving scheme (EMS)**   *Let the increment $\Delta H_\lambda = 0$, the discretized scheme is constructed as follows*:

$$\Delta q / \Delta t = \bar{\nabla}_p H(q_{k-1}, p_{k-1}, q_k, p_k), \qquad 0 = g(q_k),$$
$$\Delta p / \Delta t = -\bar{\nabla}_q H(q_{k-1}, p_{k-1}, q_k, p_k) - 2\lambda \bar{q}, \tag{61}$$

*where the unknowns $q_k$, $p_k$ and $\lambda$ can be solved for the given $q_{k-1}$ and $p_{k-1}$. The schemes are abbreviated as EMS-$\alpha$ and EMS-$\gamma$ for the integrations of $\alpha$ and $\gamma$-type, respectively.*

EMS-$\alpha$ is just the same as the energy–momentum conserving integration [11] for rigid body dynamics in terms of unit quaternions. The specific algorithms of IMS and EMS are summarized in pseudocode format in Appendix B.3.

## 4.2 Numerical integrations of higher order

Runge–Kutta methods form an important class of methods for the integration of differential equations. Recently, Jay [24] proposed the specialized partitioned additive Runge–Kutta (SPARK) methods for differential-algebraic equations (DAEs), which provide a straightforward way to construct conserving algorithms with an optimal order of convergence. Small [28] further presented the unified definition of the discontinuous collocation methods (Ref. [28], p. 64), which covers all SPARK methods of interest in this paper.

### 4.2.1 Discontinuous collocation type methods

Let $c_1, \ldots, c_s$ be distinct real numbers, $\bar{c}_1, \ldots, \bar{c}_q$ be distinct real numbers, and $\tilde{c}_0, \ldots, \tilde{c}_p$ also be distinct real numbers with $\tilde{c}_0 = 0$ and $\tilde{c}_p = 1$. Consider the Hamilton system in (12) with the consistent initial values $(q_k, p_k)$ at $t_k$ and the step-size $\Delta t$, and we search for polynomials $Q(t)$ of degree $s$, $\Lambda(t)$ of degree $p$, $P^f(t)$ of degree $q$, and $P^r(t)$ of degree $p - 1$ such that

$$Q(t_k) = q_k, \qquad P(t_k) = P^f(t_k) + P^r(t_k),$$
$$P^f(t_k) = p_k - \Delta t \tilde{b}_0 \beta \hat{\delta}(t_k), \qquad P^r(t_k) = -\Delta t \tilde{b}_0 \tilde{\delta}(t_k) \tag{62}$$

with the defects

$$\hat{\delta}(t) = \dot{P}^f(t) + H_q(Q(t), P(t)), \qquad \tilde{\delta}(t) = \dot{P}^r(t) + g_q(Q(t))\Lambda(t) \tag{63}$$

and satisfying the following conditions:

$$\dot{Q}(t_k + c_i \Delta t) = H_p(Q(t_k + c_i \Delta t), P(t_k + c_i \Delta t)), \quad i = 1, \ldots, s, \tag{64}$$

$$\dot{P}^f(t_k + \bar{c}_i \Delta t) = -H_q(Q(t_k + \bar{c}_i \Delta t), P(t_k + \bar{c}_i \Delta t)), \quad i = 1, \ldots, q, \tag{65}$$

$$\dot{P}^r(t_k + \tilde{c}_i \Delta t) = -g_q(Q(t_k + \tilde{c}_i \Delta t))\Lambda(t_k + \tilde{c}_i \Delta t), \quad i = 1, \ldots, p - 1, \tag{66}$$

$$P(t) = P^f(t) + P^r(t), \tag{67}$$

$$0 = g(Q(t_k + \tilde{c}_i \Delta t)), \quad i = 0, \ldots, p. \tag{68}$$

If these polynomials exist, the numerical solution is defined by

$$q_{k+1} = Q(t_k + \Delta t), \qquad p_{k+1} = P(t_k + \Delta t) - \Delta t \tilde{b}_p (\beta \hat{\delta}(t_{k+1}) + \tilde{\delta}(t_{k+1})),$$
$$0 = g_q(q_{k+1})H_p(q_{k+1}, p_{k+1}). \tag{69}$$

**Table 1** Comparison of s-stage Gauss–Lobatto and Lobatto IIIA-B methods

| Discontinuous collocation | $s$-stage Gauss–Lobatto | $s$-stage Lobatto IIIA-B |
|---|---|---|
| $s$ for $\boldsymbol{Q}(t)$ | $s$ | $s$ |
| $c_1, c_2, \ldots, c_s$ | Gauss points | Lobatto points |
| $q$ for $\boldsymbol{P}^f(t)$ | $s$ | $s - 2$ |
| $\bar{c}_1, \bar{c}_2, \ldots, \bar{c}_q$ | $\bar{c}_i = c_i, i = 1, \ldots, s$ | $\bar{c}_{i-1} = c_i, i = 2, \ldots, s-1$ |
| $p$ for $\Lambda(t)$ | $s$ | $s - 1$ |
| $\tilde{c}_0, \tilde{c}_1, \ldots, \tilde{c}_p$ | Lobatto points | $\tilde{c}_{i-1} = c_i, i = 1, \ldots, s$ |
| $p - 1$ or $\boldsymbol{P}^r(t)$ | $s - 1$ | $s - 2$ |
| $\tilde{c}_1, \tilde{c}_2, \ldots, \tilde{c}_{p-1}$ | $\tilde{c}_i, i = 1, \ldots, s-1$ | $\tilde{c}_{i-1} = c_i, i = 2, \ldots, s-1$ |
| $\beta$ | 0 | 1 |
| Convergence order | $2s$ | $2s - 2$ |

We take into consideration two types of Gauss SPARK methods, which can be constructed in the form of discontinuous collocation type methods with the following conditions.

**s-stage Gauss–Lobatto SPARK methods (s-G-L)** *Let $\{c_i\}, i = 1, \ldots, s$, be the s nodes of the Gauss quadrature, $p = s$, $q = s$, $\beta = 0$, $\bar{c}_i = c_i$ for $i = 1, \ldots, s$, and $\{\tilde{c}_i\}, i = 0, 1, \ldots, s$ be the $s + 1$ nodes of the Lobatto quadrature and then the discontinuous collocation methods (62)–(69) lead to the s-stage Gauss–Lobatto SPARK methods of order $2s$. The schemes are abbreviated as G-L-α and G-L-γ for the integrations of α-type and γ-type, respectively.*

**s-stage Lobatto IIIA-B SPARK methods (s-LIIIA-B)** *Let $\{c_i\}, i = 1, \ldots, s$, be the s nodes of the Lobatto quadrature, $q = s - 2$, $p = s - 1$, $\beta = 1$, $\tilde{c}_{i-1} = c_i$ and $\bar{c}_{i-1} = c_i$ for $i = 2, \ldots, s - 1$, and then the discontinuous collocation methods (62)–(69) lead to the s-stage Lobatto IIIA-B SPARK methods of order $2s - 2$. The schemes are abbreviated as LIIIA-B-α and LIIIA-B-γ for the integrations of α-type and γ-type, respectively.*

Table 1 lists the comparison between s-stage Gauss Lobatto method (s-G-L) and s-stage Lobatto IIIA-B method (s-LIIIA-B). It can be observed that s-G-L is of two orders higher convergence rate than s-LIIIA-B. These two Gauss SPARK methods are symmetric and symplectic methods proposed by Jay [24, 29]. In Appendix A, we present the unified formulations of SPARK methods and the corresponding Butcher style tableaux of SPARK coefficients for LIIIA-B ($s = 2, 3$) and G-L ($s = 1, 2$). The specific algorithms of the Gauss SPARK methods are summarized in pseudocode format in Appendix B.4.

### 4.2.2 A geometric interpretation for Gauss SPARK methods

As illustrated in Fig. 1, the Gauss SPARK methods presented by (62)–(69) can be considered as a two-step projection process, where the abbreviations

$$\boldsymbol{Q}_i = \boldsymbol{Q}(t_k + c_i \Delta t), \qquad \bar{\boldsymbol{Q}}_i = \boldsymbol{Q}(t_k + \bar{c}_i \Delta t), \qquad \tilde{\boldsymbol{Q}}_i = \boldsymbol{Q}(t_k + \tilde{c}_i \Delta t),$$
$$\boldsymbol{P}_i = \boldsymbol{P}(t_k + c_i \Delta t), \qquad \bar{\boldsymbol{Q}}_i = \boldsymbol{P}(t_k + \bar{c}_i \Delta t), \qquad \tilde{\boldsymbol{Q}}_i = \boldsymbol{P}(t_k + \tilde{c}_i \Delta t) \tag{70}$$

are used to describe the physical quantities of the discrete points. In the first step, the mapping $\boldsymbol{T}_1$ defined by the implicit functions (64), (65) and (68), gives the solutions of the intermediate variables $\boldsymbol{Q}_1, \ldots, \boldsymbol{Q}_s, \boldsymbol{P}_1, \ldots, \boldsymbol{P}_s$, and $1, \ldots, \lambda_{p-1}$. With these variables, $\boldsymbol{q}_{k+1}$,

**Fig. 1** Gauss SPARK methods overview



**Fig. 2** The discrete scheme: (**a**) 2-stage-G-L, (**b**) 3-stage-LIIIA-B

$P_{k+1}$, and $\dot{P}_{k+1}$ can be solved, serving as the input parameters in the next step. The mappings are denoted by $T_1(\alpha)$ and $T_1(\gamma)$ for the Hamilton's equations of $\alpha$-type and $\gamma$-type, respectively. In the process of $T_1$ mapping, G-L and LIIIA-B methods use two different distributions of discrete points. More specifically, the discretized points $Q_1, \ldots, Q_s$ are mismatched with the constraint points $\tilde{Q}_0, \ldots, \tilde{Q}_s$ for s-G-L, as shown in Fig. 2(a). Hence $Q_i^T Q_i \neq 1, i = 1, \ldots, s$, during the simulation. According to (44), this means

$$L(Q_i)^T L(Q_i) + Q_i Q_i^T = I_4 + O(\delta Q), \quad i = 1, \ldots, s, \tag{71}$$

where $\delta Q$ denotes the numerical error in a time step. Consequently, the mapping $T_1(\alpha)$ is not mathematically equivalent to $T_1(\gamma)$, and the discretization schemes of G-L methods are highly consistent with the assumption of the error analysis in Sect. 3.2. Contrary to s-G-L, s-LIIIA-B presents a collocated distribution between the discrete points $Q_1, \ldots, Q_s$ and

constraint points $\tilde{\boldsymbol{Q}}_0, \ldots, \tilde{\boldsymbol{Q}}_s$, as illustrated in Fig. 2(b). Hence $\boldsymbol{Q}_i^T \boldsymbol{Q}_i = 1$ for $i = 1, \ldots, s$ is satisfied exactly for s-LIIIA-B during the simulation. This may render Condition (44) unsatisfied.

In the second step, the mapping $\boldsymbol{T}_2$ defined by the implicit function (69), gives the solutions of the quantities $\boldsymbol{p}_{k+1}$ and $\lambda_p$. For s-G-L, the mapping $\boldsymbol{T}_2$ which projects the value $\boldsymbol{P}_{k+1}$ onto the manifold $\mathcal{M} = \{\boldsymbol{p}^T \boldsymbol{q} = 0\}$ to obtain $\boldsymbol{p}_{k+1}$, has nothing to do with the parameter $\alpha$ or $\gamma$. Hence the parameters $\alpha$ and $\gamma$ influence the numerical performance of the integrators solely by the first step presented in Fig. 1. In contrast, s-LIIIA-B gives a mapping $\boldsymbol{T}_2$ including the correction term

$$\boldsymbol{\delta}(t_{k+1}) = \dot{\boldsymbol{P}}(t_{k+1}) + \partial T(\boldsymbol{Q}_{k+1})/\partial \boldsymbol{q}^T + \partial V(\boldsymbol{Q}_{k+1})/\partial \boldsymbol{q} + 2\lambda \boldsymbol{Q}_{k+1}. \tag{72}$$

$\boldsymbol{T}_2$ not only preserves the path of the solution on the manifold $\mathcal{M}$ but also recovers the convergence order of $\boldsymbol{P}_{k+1}$ from $2s - 4$ to $2s - 2$. Figure 3 gives a geometric interpretation of the correction term, and one can refer to Ref. [25] (pp. 35–38, 247–252) for the discussion in detail. In general, the term $\partial T(\boldsymbol{Q}_{k+1})/\partial \boldsymbol{q}^T$ in (72) can be expanded as

$$\partial T(\boldsymbol{Q}_{k+1})/\partial \boldsymbol{q}^T = \delta v(\boldsymbol{Q}_{k+1}, \boldsymbol{P}_{k+1}) + \frac{1}{4}\boldsymbol{L}(\boldsymbol{P}_{k+1})^T \boldsymbol{J} \boldsymbol{L}(\boldsymbol{P}_{k+1}) \boldsymbol{Q}_{k+1}, \tag{73}$$

where

$$\delta v(\boldsymbol{Q}_{k+1}, \boldsymbol{P}_{k+1}) = -\frac{1}{4}\gamma^{-1}\boldsymbol{L}(\boldsymbol{P}_{k+1})^T \boldsymbol{L}(\boldsymbol{P}_{k+1}) \boldsymbol{Q}_{k+1} \tag{74}$$

for the Hamilton's equations of $\gamma$-type, and

$$\delta v(\boldsymbol{Q}_{k+1}, \boldsymbol{P}_{k+1}) = \frac{1}{4}\alpha^{-1}\boldsymbol{P}_{k+1}\boldsymbol{P}_{k+1}^T \boldsymbol{Q}_{k+1} \tag{75}$$

for the Hamilton equations of $\alpha$-type. Equation (74) can be reformulated as

$$\delta v(\boldsymbol{Q}_{k+1}, \boldsymbol{P}_{k+1}) = \frac{1}{8}\gamma^{-1}\partial \big[ \boldsymbol{P}_{k+1}^T \big( \boldsymbol{I}_4 - \boldsymbol{L}(\boldsymbol{Q}_{k+1})^T \boldsymbol{L}(\boldsymbol{Q}_{k+1}) \big) \boldsymbol{P}_{k+1} \big] / \partial \boldsymbol{Q}_{k+1}^T. \tag{76}$$

Noting that $\boldsymbol{Q}_{k+1}^T \boldsymbol{Q}_{k+1} = 1$ has been satisfied in the first step, we have

$$\boldsymbol{I}_4 - \boldsymbol{L}(\boldsymbol{Q}_{k+1})^T \boldsymbol{L}(\boldsymbol{Q}_{k+1}) = \boldsymbol{Q}_{k+1}\boldsymbol{Q}_{k+1}^T. \tag{77}$$

Substituting (77) into (76) and setting $\gamma = \alpha$ lead to

$$\delta v(\boldsymbol{Q}_{k+1}, \boldsymbol{P}_{k+1}) = \frac{1}{8}\gamma^{-1}\partial \big[ \boldsymbol{P}_{k+1}^T \boldsymbol{Q}_{k+1}\boldsymbol{Q}_{k+1}^T \boldsymbol{P}_{k+1} \big] / \partial \boldsymbol{Q}_{k+1}^T = \frac{1}{4}\alpha^{-1}\boldsymbol{P}_{k+1}\boldsymbol{P}_{k+1}^T \boldsymbol{Q}_{k+1} \bigg|_{\gamma=\alpha}. \tag{78}$$

Recall that the mappings are denoted by $T_2(\alpha)$ and $T_2(\gamma)$ respectively for the Hamilton's equations of $\alpha$-type and $\gamma$-type. Equation (78) means the mapping $T_2(\gamma)$ is mathematical equivalent to $T_2(\alpha)$, and consequently $p_{k+1}$ will be solved with the same numerical accuracy for s-LIIIA-B-$\gamma$ and s-LIIIA-B-$\alpha$ if $q_{k+1}$, $P_{k+1}$ and $\dot{P}_{k+1}$ are assigned with the same inputs.

Go back to the first step to reconsider the mapping $T_1$ for s-LIIIA-B. Although $\partial T/\partial Q_i$, $i = 2, \ldots, s - 1$ and $\partial T/\partial P_i$, $i = 1, \ldots, s$ are also involved in the mapping $T_1$ in this case, we must note that $Q_i^T Q_i = 1$ for $i = 1, \ldots, s$ is only satisfied implicitly, which means $Q_i^T Q_i = 1$ is an iteration solution of the mapping $T_1$. Hence, we cannot present a straightforward and analytical way which is just like those presented by (76)–(78), to prove $T_1(\alpha) = T_1(\gamma)$. Nevertheless, the numerical testing in the next section would show that s-LIIIA-B-$\alpha$ and s-LIIIA-B-$\gamma$ are of the same numerical accuracy if $\alpha$ and $\gamma$ are assigned with the same value and this implies the equivalence between $T_1(\alpha)$ and $T_1(\gamma)$.

These differences presented above make G-L and LIIIA-B essentially two different methods in the numerical simulation, and we will further investigate the numerical performance of these methods in the following.

## 5 Numerical examples

In this section, numerical examples of a top spinning in gravitational field are considered. As shown in Fig. 4, the top fixed at $O$ is rotating in a uniform gravitational field with the acceleration $g = 9.81$ in the negative $z$-direction. $l$ represents the distance between the mass center $A$ and the fixed point $O$, and the rotational motion is expressed by the angle of nutation $\theta$, the angle of precession $\phi$, and the spin angle $\psi$. The gravitational potential energy is written as [11, 16]

$$V(q) = mgl \cos \theta = mgl q^T K q, \tag{79}$$

where $K = \text{diag}(1, -1, -1, 1)$, and then we have $\partial V(q)/\partial q^T = 2mgl K q$.

Two numerical examples are considered for comparing the numerical performance of different types of integrators in simulation, and the discussion would concentrate on the numerical influence of the parameters $\alpha$ and $\gamma$ in different inertia representations. The configuration parameters and initial conditions for the two examples are presented as follows.

**Fast spinning top**  The configuration parameters are corresponding to $m = 1$, $l = 0.04$ and the principal moments of inertia tensor with respect to the fixed point $[I_1 = I_2, I_3] = [0.002, 0.0008]$. The following initial conditions are considered:

$$[\phi_0, \theta_0, \psi_0] = [0, \pi/6, 0], \qquad \Omega_0^T = [0, 0, 40\pi]. \tag{80}$$

**Fig. 4** A heavy top

Correspondingly, the initial value $\boldsymbol{q}_0$ and $\boldsymbol{p}_0$ can be obtained by

$$
\boldsymbol{q}_0^T = \left[ \cos\left(\frac{\phi_0 + \psi_0}{2}\right) \cos\left(\frac{\theta_0}{2}\right), \cos\left(\frac{\phi_0 - \psi_0}{2}\right) \sin\left(\frac{\theta_0}{2}\right), \right.
$$
$$
\left. \sin\left(\frac{\phi_0 - \psi_0}{2}\right) \sin\left(\frac{\theta_0}{2}\right), \sin\left(\frac{\phi_0 + \psi_0}{2}\right) \cos\left(\frac{\theta_0}{2}\right) \right] \tag{81}
$$

and

$$
\boldsymbol{p}_0 = 2L(\boldsymbol{q}_0)^T J \boldsymbol{\Omega}_0. \tag{82}
$$

**Regular precession** The second example is the regular precession in gravitational field. The top is represented as a cone with the dimensions equivalent to those used in [17, 23]. The configuration parameters include the height $h = 0.1$, the radius $r = \frac{1}{2}h$, the length $l = \frac{3}{4}h$, and the mass $m = \frac{1}{3}\rho\pi^2 h$ with the density $\rho = 2700$. The three principal moments of the inertia matrix with respect to fixed point $O$ are given by

$$
I_1 = I_2 = \frac{3}{5}m\left(r^2/4 + h^2\right), \qquad I_3 = \frac{3}{10}mr^2. \tag{83}
$$

The initial conditions are as follows:

$$
[\phi_0, \theta_0, \psi_0] = [0, \pi/3, 0] \tag{84}
$$

with a precession rate $\dot{\phi} = 10$, and the velocity components $\dot{\phi}$ and $\dot{\psi}$ satisfying the relation

$$
\dot{\psi} = mgl/I_3\dot{\phi} + \frac{(I_1 - I_3)}{I_3}\dot{\phi}\cos(\theta). \tag{85}
$$

Correspondingly, the initial angular velocity vector can be obtained by

$$
\boldsymbol{\Omega}_0^T = \left[0, \dot{\phi}\sin(\theta_0), \dot{\psi} + \dot{\phi}\cos(\theta_0)\right]. \tag{86}
$$

To precisely evaluate the numerical results among these integrators, we define the mean values of errors as follows:

$$
E\big[|\Delta x_c|\big] = \frac{1}{N_t + 1}\sum_{k=0}^{N_t}|\Delta x_{c,k}|, \qquad E\big[|\Delta z_c|\big] = \frac{1}{N_t + 1}\sum_{k=0}^{N_t}|\Delta z_{c,k}|, \tag{87}
$$

$$
E\big[|\Delta H|\big] = \frac{1}{N_t + 1}\sum_{k=0}^{N_t}|\Delta H_k|, \tag{88}
$$

$$
E\big[|\Delta l_3|\big] = \frac{1}{N_t + 1}\sum_{k=0}^{N_t}|\Delta l_{3,k}|, \tag{89}
$$

$$
E\big[|g(\boldsymbol{q})|\big] = \frac{1}{N_t + 1}\sum_{k=0}^{N_t}\big|\boldsymbol{q}_k^T\boldsymbol{q}_k - 1\big|, \tag{90}
$$

$$
E\big[|\boldsymbol{q}^T\boldsymbol{p}|\big] = \frac{1}{N_t + 1}\sum_{k=0}^{N_t}\big|\boldsymbol{q}_k^T\boldsymbol{p}_k\big|. \tag{91}
$$

In the above, $N_t$ denotes the number of the total time steps; $x_c$ and $z_c$ denote the coordinates of mass center in $x$ and $z$-directions, whose errors are denoted by $\Delta x_{c,k}$ and $\Delta z_{c,k}$ at the grid point $t_k = k\Delta t$ ($k = 0, 1, \ldots, N_t$); $\Delta H_k$ denotes the pointwise absolute energy error; $\Delta l_{3,k}$ denotes the pointwise local angular momentum error of the $l_3$-component.

## 5.1 Conservation of invariants

Conservation of invariants (or termed conserved quantities) has important effects on the numerical performance of integrations. It has been suggested that numerical methods, who conserve invariants (especially quadratic invariants) automatically, tend to present good performance on long-time behavior in simulation [25]. In this paper, two different inertia representations and the corresponding parameters $\gamma$ and $\alpha$ are introduced when constructing the numerical integrations. It is interesting whether the introduction of the parameters in different inertia representations would destroy the conservation of invariants which has been held by numerical methods.
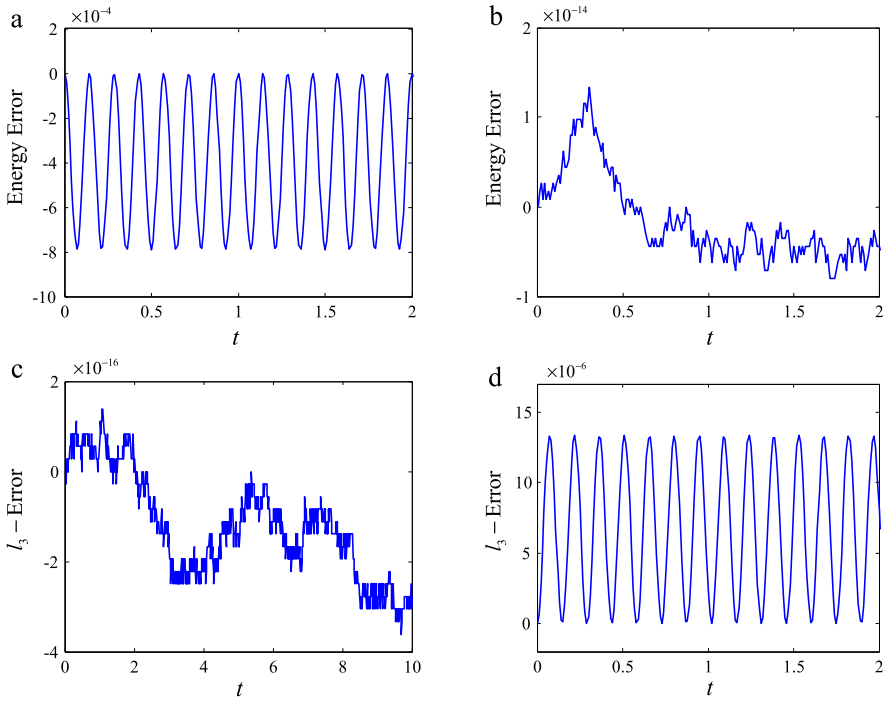
Tables 2 and 3 list the numerical results for the two examples, where the error terms $E[|\Delta H|]$, $E[|\Delta l_3|]$, $E[|g(\boldsymbol{q})|]$ and $E[|\boldsymbol{q}^T \boldsymbol{p}|]$ are defined by (88)–(91). Numerical integrations of $\gamma$-type and $\alpha$-type are both considered in the simulation. The errors are bolded if the invariants are conserved exactly by the integrations. We further show the numerical errors of the invariants (i.e., the total energy and the local angular momentum $l_3$) in Figs. 5 and 6 for IMS-$\gamma$, EMS-$\gamma$ and 3-LIIIA-B-$\gamma$ for the first example. Figure 7 shows that $\boldsymbol{q}^T \boldsymbol{p} = 0$ is not preserved exactly for IMS and EMS but presents a periodic error in simulation. We should emphasize that all these integrations present periodic errors if the invariants are not conserved exactly in the simulation. These numerical results suggest that the introduction of the parameter $\gamma$ or $\alpha$ does not destroy the conservation of an invariant as well as the good long-time behavior, which has been held automatically by the numerical methods.

**Table 2** Conservation of invariants for fast spinning top, $\gamma = \gamma_m$, $\alpha = 100\gamma_m$, $\Delta t = 0.001$, $N_t = 4000$

| Algorithms | $E[|\Delta H|]$ | $E[|\Delta l_3|]$ | $E[|g(\boldsymbol{q})|]$ | $E[|\boldsymbol{q}^T \boldsymbol{p}|]$ |
|---|---|---|---|---|
| IMS-$\gamma$ | 2.54E-06 | **1.27E-15** | **2.97E-17** | 7.09E-08 |
| IMS-$\alpha$ | 1.02E-05 | **7.66E-16** | **5.95E-17** | 7.10E-08 |
| EMS-$\gamma$ | **6.72E-14** | 8.09E-08 | **3.24E-17** | 7.09E-08 |
| EMS-$\alpha$ | **9.68E-14** | 3.25E-07 | **2.88E-17** | 1.30E-10 |
| 3-LIIIA-B-$\gamma$ | 5.90E-10 | **1.62E-14** | **4.30E-17** | **1.01E-17** |
| 3-LIIIA-B-$\alpha$ | 3.97E-10 | **3.28E-13** | **5.00E-17** | **1.02E-17** |
| 2-G-L-$\gamma$ | 5.34E-10 | **1.97E-15** | **3.09E-17** | **7.69E-18** |
| 2-G-L-$\alpha$ | 1.75E-09 | **1.36E-13** | **4.96E-17** | **7.56E-18** |

**Table 3** Conservation of invariants for regular precession, $\gamma = \gamma_m$, $\Delta t = 0.005$, $N_t = 800$

| Algorithms | $E[|\Delta H|]$ | $E[|\Delta l_3|]$ | $E[|g(\boldsymbol{q})|]$ | $E[|\boldsymbol{q}^T \boldsymbol{p}|]$ |
|---|---|---|---|---|
| IMS-$\gamma$ | 0.0139 | **4.48E-14** | **6.05E-18** | 4.82E-05 |
| IMS-$\alpha$ | 0.0806 | **6.85E-14** | **6.27E-18** | 1.16E-04 |
| EMS-$\gamma$ | **2.00E-12** | 8.23E-04 | **5.47E-18** | 9.62E-05 |
| EMS-$\alpha$ | **5.52E-12** | 0.0051 | **6.16E-18** | **2.26E-14** |
| 3-L IIIA-B-$\gamma$ | 2.34E-07 | **6.04E-13** | **1.04E-17** | **2.21E-15** |
| 3-L IIIA-B-$\alpha$ | 2.34E-07 | **8.02E-13** | **1.08E-17** | **1.02E-17** |
| 2-G-L-$\gamma$ | 6.09E-07 | **4.70E-12** | **1.19E-17** | **1.43E-15** |
| 2-G-L-$\alpha$ | 2.57E-06 | **2.78E-13** | **9.66E-18** | **1.39E-15** |

**Fig. 5** Conservation of invariants for fast spinning top: (**a**) energy error of IMS-$\gamma$, (**b**) energy error of EMS-$\gamma$, (**c**) $l_3$-error of IMS-$\gamma$, (**d**) $l_3$-error of EMS-$\gamma$; $\Delta t = 0.01$, $\gamma = \gamma_{\rm m}$



**Fig. 6** Conservation of invariants for fast spinning top with 3-LIIIA-B-$\gamma$: (**a**) energy error, (**b**) $l_3$-error, $\Delta t = 0.01$, $\gamma = \gamma_{\rm m}$

## 5.2 Comparison of the numerical accuracy

Tables 4, 5 and 6 list the numerical errors of the integrations for the two examples, where the error terms $E[|\Delta H|]$, $E[|\Delta z|]$ and $E[|\Delta x|]$ are defined by (87) and (88). Numerical integrations of $\gamma$-type and $\alpha$-type are both considered. Figures 8 and 9 further show the numerical results of IMS-$\alpha$ and IMS-$\gamma$ for the two examples. It can be observed that all the integrations of $\alpha$-type are of the same numerical accuracy no matter what values are assigned

**Fig. 7** The orthogonal condition for fast spinning top: (**a**) $q^T p$-error for IMS-$\gamma$, (**b**) $q^T p$-error for EMS-$\gamma$, $\Delta t = 0.01$, $\gamma = \gamma_m$

**Table 4** Fast spinning top: numerical accuracy for the integrations of $\gamma$-type, $I_3 = 0.0008$, $\Delta t = 0.001$, $N_t = 320$

| $\gamma$-type | | IMS | EMS | 3-LIIIA-B | 2-G-L |
|---|---|---|---|---|---|
| $E[|\Delta H|]$ | $\gamma = \gamma_m$ | 2.35E-06 | 1.12E-14 | 5.46E-10 | 4.94E-10 |
| | $\gamma = \gamma_h$ | 4.70E-06 | 8.49E-15 | 5.46E-10 | 2.72E-10 |
| | $\gamma = 100\gamma_m$ | 9.32E-06 | 1.49E-14 | 3.79E-10 | 1.61E-09 |
| $E[|\Delta z|]$ | $\gamma = \gamma_m$ | 1.48E-05 | 1.55E-05 | 7.41E-09 | 4.86E-09 |
| | $\gamma = \gamma_h$ | 2.76E-05 | 5.09E-05 | 7.41E-09 | 1.96E-09 |
| | $\gamma = 100\gamma_m$ | 1.77E-04 | 2.18E-04 | 5.41E-09 | 1.96E-08 |
| $E[|\Delta x|]$ | $\gamma = \gamma_m$ | 2.07E-04 | 2.16E-04 | 1.21E-08 | 9.83E-09 |
| | $\gamma = \gamma_h$ | 2.20E-04 | 2.20E-04 | 1.22E-08 | 8.16E-09 |
| | $\gamma = 100\gamma_m$ | 3.23E-04 | 3.23E-04 | 8.87E-09 | 3.25E-08 |

**Table 5** Fast spinning top: numerical accuracy for the integrations of $\alpha$-type, $I_3 = 0.0008$, $\Delta t = 0.001$, $N_t = 320$

| $\alpha$-type | | IMS | EMS | 3-LIIIA-B | 2-G-L |
|---|---|---|---|---|---|
| $E[|\Delta H|]$ | $\alpha = \gamma_m$ | 9.43E-06 | 1.08E-14 | 5.46E-10 | 1.60E-09 |
| | $\alpha = \gamma_h$ | 9.43E-06 | 1.39E-14 | 5.46E-10 | 1.60E-09 |
| | $\alpha = 100\gamma_m$ | 9.43E-06 | 1.03E-14 | 3.77E-10 | 1.62E-09 |
| $E[|\Delta z|]$ | $\alpha = \gamma_m$ | 1.78E-04 | 2.20E-04 | 7.41E-09 | 1.96E-08 |
| | $\alpha = \gamma_h$ | 1.78E-04 | 2.20E-04 | 7.41E-09 | 1.96E-08 |
| | $\alpha = 100\gamma_m$ | 1.78E-04 | 2.20E-04 | 5.41E-09 | 1.98E-08 |
| $E[|\Delta x|]$ | $\alpha = \gamma_m$ | 3.26E-04 | 3.80E-04 | 1.21E-08 | 3.25E-08 |
| | $\alpha = \gamma_h$ | 3.26E-04 | 3.80E-04 | 1.22E-08 | 3.25E-08 |
| | $\alpha = 100\gamma_m$ | 3.26E-04 | 3.80E-04 | 8.87E-09 | 3.27E-08 |

to the parameter $\alpha$; IMS-$\gamma$, EMS-$\gamma$ and 2-G-L-$\gamma$ with $\gamma = \gamma_m$ are of much higher numerical accuracy than those of $\alpha$-type. These numerical results suggest that the parameter $\gamma$ can be used to improve the numerical accuracy of numerical methods, whereas the parameter $\alpha$ has no influence on the numerical accuracy of the integrations. Tables 4–6 further present that

**Table 6** Regular precession: numerical accuracy for integrations, $r = \frac{1}{2}h$, $\Delta t = 0.005$, $N_t = 80$

|  |  |  | IMS | EMS | 3-LIIIA-B | 2-G-L |
|---|---|---|---|---|---|---|
| $E[\|\Delta H\|]$ | $\gamma$-type | $\gamma = \gamma_m$ | 1.69E-02 | 9.03E-12 | 2.15E-06 | 1.34E-06 |
|  |  | $\gamma = \gamma_h$ | 5.09E-01 | 3.91E-12 | 3.06E-06 | 1.08E-06 |
|  |  | $\gamma = \tilde{\gamma}_h$ | 5.57E-12 | 4.17E-12 | 1.69E-06 | 2.56E-06 |
|  | $\alpha$-type | $\alpha = \gamma_m$ | 1.07E-01 | 2.94E-12 | 2.15E-06 | 7.00E-06 |
| $E[\|\Delta z\|]$ | $\gamma$-type | $\gamma = \gamma_m$ | 1.74E-03 | 3.24E-03 | 1.96E-05 | 1.55E-05 |
|  |  | $\gamma = \gamma_h$ | 9.20E-03 | 1.44E-02 | 2.33E-05 | 1.38E-05 |
|  |  | $\gamma = \tilde{\gamma}_h$ | 5.61E-16 | 7.81E-16 | 1.73E-05 | 2.14E-05 |
|  | $\alpha$-type | $\alpha = \gamma_m$ | 4.39E-03 | 8.88E-03 | 1.96E-05 | 3.53E-05 |
| $E[\|\Delta x\|]$ | $\gamma$-type | $\gamma = \gamma_m$ | 4.35E-02 | 8.19E-03 | 5.80E-04 | 5.42E-04 |
|  |  | $\gamma = \gamma_h$ | 2.11E-01 | 4.51E-01 | 6.82E-04 | 2.91E-04 |
|  |  | $\gamma = \tilde{\gamma}_h$ | 9.05E-02 | 1.01E-01 | 5.18E-04 | 7.09E-04 |
|  | $\alpha$-type | $\alpha = \gamma_m$ | 2.03E-01 | 3.18E-01 | 5.80E-04 | 1.11E-03 |



**Fig. 8** Trajectory error of mass center for fast spinning top: (**a**) $x$-component error, (**b**) $z$-component error; $\Delta t = 0.01$; IMS-$\gamma$ ($\times$), $\gamma = \gamma_m$; IMS-$\alpha$ (- - -), $\alpha = \gamma_m$; analytical (—)

3-LIIIA-B-$\gamma$ is of the same numerical accuracy with 3-LIIIA-B-$\alpha$ if $\alpha$ and $\gamma$ are assigned with the same values. According to the discussion in Sect. 4.2.2, numerical results suggest that $T_1(\alpha) = T_1(\gamma)$ for LIIIA-B, and there is no difference in accuracy between the two discretization schemes s-LIIIA-B-$\alpha$ and s-LIIIA-B-$\gamma$.

In Table 6, besides $\gamma_m$ and $\gamma_h$, $\gamma$ is also assigned with $\tilde{\gamma}_h$ to account for the large difference between $I_1$ and $I_3$ (i.e., $I_1 = I_2 = 8.5I_3$). It can be observed in Table 6 that IMS-$\gamma$ and EMS-$\gamma$ with $\gamma = \gamma_h$ have larger numerical errors than the others and the numerical accuracy is improved evidently if we use the corrected optimal value $\tilde{\gamma}_h$ presented in (50). We further change the radius by $r = 1.5h$ for the example of regular precession which makes $I_1 = I_2 = 1.39I_3$. It can be observed in Table 7 that the numerical accuracy of IMS-$\gamma$ and EMS-$\gamma$ assigned with $\gamma = \gamma_h$ is greatly improved, especially compared with those presented in Table 6. Recall the discussions in the last paragraph of Sect. 3.3, and these numerical results are highly consistent with those discussions.

Figure 10 shows the relative error with time step increasing, where the error is obtained by evaluating the maximum of the relative error. These numerical results can be summarized as follows: IMS and EMS have the convergence of order 2; 2-G-L and 3-LIII-A-B have the

**Fig. 9** Trajectory error of mass center for regular precession top: (**a**) $x$-component error, (**b**) $z$-component error; $\Delta t = 0.007$; IMS-$\gamma$ ($\times$), $\gamma = \gamma_{\mathrm{m}}$; IMS-$\alpha$(- - -), $\alpha = \gamma_{\mathrm{m}}$; analytical (—)
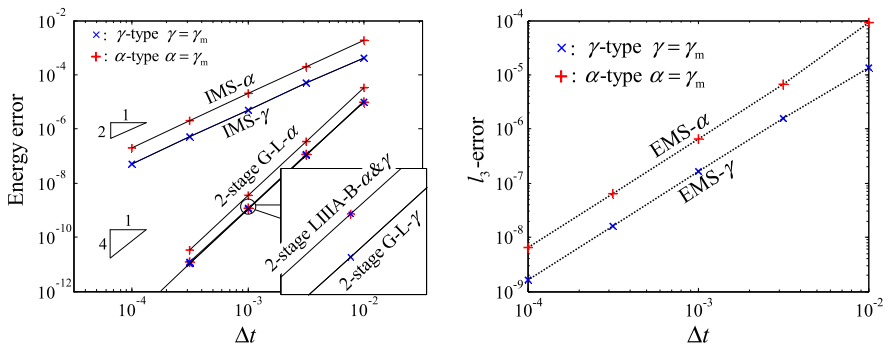
**Table 7** Regular precession: numerical accuracy for integrations, $r = 1.5h$, $\Delta t = 0.01$, $N_t = 80$

|  |  |  | IMS | EMS | 3-LIIIA-B | 2-G-L |
|---|---|---|---|---|---|---|
| $E[|\Delta H|]$ | $\gamma$-type | $\gamma = \gamma_{\mathrm{m}}$ | 3.15E-05 | 3.03E-12 | 7.32E-11 | 3.48E-12 |
|  |  | $\gamma = \gamma_{\mathrm{h}}$ | 4.99E-05 | 2.58E-12 | 7.49E-11 | 2.47E-12 |
|  | $\alpha$-type | $\alpha = \gamma_{\mathrm{m}}$ | 2.92E-03 | 1.46E-12 | 7.10E-11 | 7.72E-09 |
| $E[|\Delta z|]$ | $\gamma$-type | $\gamma = \gamma_{\mathrm{m}}$ | 5.94E-05 | 2.16E-05 | 9.03E-08 | 2.15E-08 |
|  |  | $\gamma = \gamma_{\mathrm{h}}$ | 7.49E-05 | 2.72E-05 | 9.04E-08 | 5.65E-10 |
|  | $\alpha$-type | $\alpha = \gamma_{\mathrm{m}}$ | 5.72E-04 | 2.14E-04 | 9.03E-08 | 9.32E-07 |
| $E[|\Delta x|]$ | $\gamma$-type | $\gamma = \gamma_{\mathrm{m}}$ | 1.01E-02 | 1.36E-02 | 1.91E-05 | 6.84E-06 |
|  |  | $\gamma = \gamma_{\mathrm{h}}$ | 8.74E-03 | 1.31E-02 | 1.91E-05 | 8.64E-06 |
|  | $\alpha$-type | $\alpha = \gamma_{\mathrm{m}}$ | 6.70E-02 | 3.51E-02 | 1.91E-05 | 6.74E-05 |



**Fig. 10** Fast spinning top: the relative periodic error with time step increasing

convergence of order 4. It suggests that the convergence property keeps consistent regardless of whether the integrations are derived in $\gamma$-type or $\alpha$-type, and all the integrations of $\gamma$-type present higher accuracy than those of $\alpha$-type, except for 3-LIII-A-B.

The discussions and conclusions presented in Sects. 5.1 and 5.2 are also suitable to those SPARK methods of order 2, 6 or 8, although the numerical results are only listed for the cases

of order 4. According to these numerical results, the inertia representation of $\gamma$-type can be used to improve the numerical accuracy of integrations without destroying the conservation of invariants. These results further suggest that the improvement of numerical accuracy can be achieved simultaneously for both the energy and trajectory by choosing a proper value of the parameter $\gamma$. In addition, the integrations of $\gamma$-type with the proposed optimal values $\gamma_{\mathrm{m}}$ and $\gamma_{\mathrm{h}}$ can give impressively good numerical accuracy, especially compared with those of $\alpha$-type, and $\gamma_{\mathrm{m}}$ is more robust for IMS and EMS if there are large differences among the three principal moments.

### 5.3 Comparisons between the integrations of $\gamma$ and $\alpha$-types

The conservation of invariant and the numerical accuracy of the integrations of $\alpha$-type and $\gamma$-type have been investigated in the above. Many other differences of numerical performance will be found if we consider the convergence orders, error-parameter relations, the discretization schemes and their combinations, as discussed in the following.

#### 5.3.1 The error-parameter relations for IMS, EMS and G-L SPARK methods

Figures 11, 12 and 13 shows the linear relation between the periodic error and $1/\gamma$ for IMS-$\gamma$, EMS-$\gamma$ and G-L-$\gamma$ in the example of fast spinning top. Figure 13 further shows that



**Fig. 11** Energy error of IMS-$\gamma$ for fast spinning top: (**a**) the periodic error-time curve, $\Delta t = 0.001$, (**b**) the periodic error-$1/\gamma$ curve

**Fig. 12** Fast spinning top: the periodic error of $l_3$-angular momentum for EMS-$\gamma$, $\Delta t = 0.001$

**Fig. 13** Periodic energy error with $1/\gamma$ increasing for fast spinning top: (**a**) 1-stage G-L-$\gamma$, (**b**) 2-stage G-L-$\gamma$, (**c**) 3-stage G-L-$\gamma$, (**d**) 4-stage G-L-$\gamma$

the linear relation remains unchanged for G-L SPARK methods of $\gamma$-type when the order of the integration is changed. This implies that the error estimation discussed in Sect. 3.2 is probably independent of the order conditions. Figure 14 shows the relation between the periodic error and $1/\alpha$ for the example of the fast spinning top. It can be observed that the parameter $\alpha$ has nearly no influence on the energy error for the integrations of $\alpha$-type. As referred to the discussion in Sect. 4.2.2, all three integrations IMS, EMS and G-L show the mismatch of the discretized points between the differential part and the algebraic part. We believe that this mismatch makes these results highly consistent with the error estimation in Sect. 3.2.

### 5.3.2 The error analysis for LIIIA-B SPARK methods

Figure 15 shows the periodic error of LIIIA-B-$\gamma$ changing with $1/\gamma$ for the example of the fast spinning top. The figures of LIIIA-B-$\alpha$ are not shown in this paper because s-LIIIA-B-$\alpha$ presents exactly the same results as s-LIIIA-B-$\gamma$. An interesting numerical phenomenon can be observed: the numerical errors are linear functions with $1/\gamma$ for 2-LIIIA-B and 4-LIIIA-B, whereas 3-LIIIA-B and 5-LIIIA-B present a logarithmic increasing (or decreasing) error-curve. Figure 16 further shows the convergence rate for the LIIIA-B SPARK methods. It presents that the parameter $\gamma$ influences the numerical accuracy of 2-LIIIA-B and 4-LIIIA-B regardless of the size of the time step, and in contrast the influence of parameter $\gamma$ declines rapidly for 3-LIIIA-B and 5-LIIIA-B with the time step decreasing.

**Fig. 14** Energy error changes with $1/\alpha$ for fast spinning top: (**a**) IMS-$\alpha$, (**b**) EMS-$\alpha$ (**c**) 2-stage G-L-$\alpha$, (**d**) 3-stage G-L-$\alpha$, (**e**) 4-stage G-L-$\alpha$

Although we predicted the mathematical equivalence of LIIIA-B-$\alpha$ and LIIIA-B-$\gamma$ in Sect. 4.2.2 and confirmed it by numerical results, the parameter $\alpha$ or $\gamma$ is not entirely useless for improving the numerical accuracy. Rather, it may somehow affect the integrations whose numbers of the stage are even, and we can still apply $\gamma_m$ and $\gamma_h$ to reduce the integration error, at least for even-stage schemes.

Though the results of 2-LIIIA-B and 4-LIIIA-B is preferable, the authors are not able to give an explanation to this phenomenon. Further work is needed to precisely study the role of parameters $\alpha$ and $\gamma$ for the LIIIA-B SPARK methods.

**Fig. 15** Energy error changes with $1/\gamma$ increasing for fast spinning top: (**a**) $(2, 1)$-LIIIA-B-$\gamma$, (**b**) $(3, 2)$-LIIIA-B-$\gamma$, (**c**) $(4, 3)$-LIIIA-B-$\gamma$, (**d**) $(5, 4)$-LIIIA-B-$\gamma$

**Fig. 16** Fast spinning top: the periodic relative error with time step increasing for LIIIA-B



## 5.4 Robustness of the integrations with different parameters

Robustness of the parameters $\gamma$ and $\alpha$ is also discussed in this paper. The introduction of undetermined parameters to differential equations generally involves potential numerical risks if inappropriate values of the parameters are applied in the algorithm. Firstly, this may destroy the good long-time behavior, which has been discussed for the parameters $\gamma$ and $\alpha$ in Sect. 5.1. Secondly, an inappropriate value of the introduced parameter may cause an ill-conditioned problem of the iteration matrix, which means serious restriction of the time step to obtain a convergent result of Newton iteration.

**Table 8** Fast spinning top: maximum of the time step for convergent results

| Convergence order | | Parameter | IMS | EMS | LIII A-B | G-L |
|---|---|---|---|---|---|---|
| 2nd order | $\gamma$-type | $\gamma = \gamma_m$ | 0.022 | 0.019 | 0.016 | 0.020 |
| | | $\gamma = \gamma_h$ | 0.035 | 0.022 | 0.017 | 0.022 |
| | $\alpha$-type | $\alpha = \gamma_m$ | 0.084 | 0.025 | 0.027 | divergent |
| | | $\alpha = \gamma_h$ | 0.085 | 0.024 | 0.020 | divergent |
| | | $\alpha = 100\gamma_m$ | 0.007 | 0.019 | 0.0011 | divergent |
| | | $\alpha = 0.01\gamma_m$ | 0.021 | 0.012 | 0.0014 | divergent |
| 4th order | $\gamma$-type | $\gamma = \gamma_m$ | – | – | 0.024 | 0.025 |
| | | $\gamma = \gamma_h$ | – | – | 0.028 | 0.028 |
| | $\alpha$-type | $\alpha = \gamma_m$ | – | – | 0.014 | 0.016 |
| | | $\alpha = \gamma_h$ | – | – | 0.014 | 0.017 |
| | | $\alpha = 100\gamma_m$ | – | – | 0.0016 | 0.0016 |
| | | $\alpha = 0.01\gamma_m$ | – | – | 0.0048 | 0.014 |
| 6th order | $\gamma$-type | $\gamma = \gamma_m$ | – | – | 0.029 | 0.030 |
| | | $\gamma = \gamma_h$ | – | – | 0.033 | 0.032 |
| | $\alpha$-type | $\alpha = \gamma_m$ | – | – | 0.018 | 0.021 |
| | | $\alpha = \gamma_h$ | – | – | 0.018 | 0.021 |
| | | $\alpha = 100\gamma_m$ | – | – | 0.0022 | 0.0019 |
| | | $\alpha = 0.01\gamma_m$ | – | – | 0.018 | 0.015 |
| 8th order | $\gamma$-type | $\gamma = \gamma_m$ | – | – | 0.030 | 0.029 |
| | | $\gamma = \gamma_h$ | – | – | 0.033 | 0.033 |
| | $\alpha$-type | $\alpha = \gamma_m$ | – | – | 0.022 | 0.022 |
| | | $\alpha = \gamma_h$ | – | – | 0.022 | 0.024 |
| | | $\alpha = 100\gamma_m$ | – | – | 0.0025 | 0.0025 |
| | | $\alpha = 0.01\gamma_m$ | – | – | 0.019 | 0.016 |

Annotation: $\varepsilon_r = 10^{-20}$, $iter\_max = 50$, $N_t = 160$

Table 8 lists the maximum of the size of time-step (denoted as $\Delta t_{max}$) to receive a convergent result of integrations for the example of the fast spinning top, where $iter\_max$ denotes the maximum number of iterations, $\varepsilon_r$ denotes the maximum of the iteration error and $N_t$ denotes the number of the total time steps. The Newton iteration stops and goes to the next time step if the iteration number is greater than the maximum $iter\_max$ or the iteration error is less than $\varepsilon_r$. It can be observed that: IMS-$\alpha$ and EMS-$\alpha$ allow larger sizes of $\Delta t_{max}$ than IMS-$\gamma$ and EMS-$\gamma$; G-L-$\gamma$ assigned with $\gamma = \gamma_m$ or $\gamma_h$ allows a larger size of $\Delta t_{max}$ than G-L-$\alpha$; LIII A-B-$\gamma$ assigned with $\gamma = \gamma_m$ or $\gamma_h$ allows a larger size of $\Delta t_{max}$ for $s = 3, 4$ and 5. Numerical results also suggest that 1-stage G-L-$\alpha$ is divergent in simulation. It seems that $(1, 1)$-G-L-$\alpha$ suffers from serious morbidity of the iteration matrix, and plenty of numerical tests suggest that the energy error will increase linearly for 1-G-L-$\alpha$ even if it is assigned with a very small size of time step.

Tables 9, 10, 11, 12, 13 and 14 list the mean values of iterations at which the iteration error becomes less than $\varepsilon_r$ for IMS, EMS and LIII A-B. We select four different sizes of the time step to compare the convergent speeds of different integrations. It can be observed that the integrations of $\gamma$-type assigned with $\gamma = \gamma_m$ generally have quicker convergence than those of $\alpha$-type. Table 9 shows that the iteration increases with decreasing size of the time step for IMS-$\alpha$, which means much more computational cost compared with IMS-$\gamma$.

**Table 9** Fast spinning top: iteration number for IMS

| IMS | Parameter | $\Delta t = 0.02$ | $\Delta t = 0.01$ | $\Delta t = 0.005$ | $\Delta t = 0.001$ |
|---|---|---|---|---|---|
| $\gamma$-type | $\gamma = \gamma_m$ | 9 | 6.9 | 6 | 5.1 |
| | $\gamma = \gamma_h$ | 8 | 6 | 6 | 5.0 |
| $\alpha$-type | $\alpha = \gamma_m$ | 10.9 | 14.3 | 22.7 | 54.3 |
| | $\alpha = \gamma_h$ | 10.6 | 15.1 | 29.1 | 62.8 |
| | $\alpha = 100\gamma_m$ | divergent | divergent | divergent | 92.6 |
| | $\alpha = 0.01\gamma_m$ | 12 | 8.0 | 15.65 | 41.9 |

Annotation: $\varepsilon_r = 10^{-15}$, $iter\_max = 100$, $N_t = 80$

**Table 10** Fast spinning top: iteration number for EMS

| EMS | Parameter | $\Delta t = 0.02$ | $\Delta t = 0.01$ | $\Delta t = 0.005$ | $\Delta t = 0.001$ |
|---|---|---|---|---|---|
| $\gamma$-type | $\gamma = \gamma_m$ | divergent | 7 | 6 | 5 |
| | $\gamma = \gamma_h$ | 8 | 6 | 6 | 5 |
| $\alpha$-type | $\alpha = \gamma_m$ | 9 | 7 | 6 | 5 |
| | $\alpha = \gamma_h$ | 9 | 7 | 6 | 5 |
| | $\alpha = 100\gamma_m$ | 13 | 10.0 | 14.3 | 24.6 |
| | $\alpha = 0.01\gamma_m$ | divergent | 8 | 6 | 5 |

Annotation: $\varepsilon_r = 10^{-12}$, $iter\_max = 100$, $N_t = 80$

**Table 11** Fast spinning top: iteration number for 2-LIII A-B

| 2-stage | Parameter | $\Delta t = 0.02$ | $\Delta t = 0.01$ | $\Delta t = 0.005$ | $\Delta t = 0.001$ |
|---|---|---|---|---|---|
| $\gamma$-type | $\gamma = \gamma_m$ | 73 | 30.0 | 16.0 | 9.0 |
| | $\gamma = \gamma_h$ | 46.9 | 23.9 | 14.0 | 8.0 |
| $\alpha$-type | $\alpha = \gamma_m$ | 67.8 | divergent | 25.0 | 10.0 |
| | $\alpha = \gamma_h$ | 74.9 | divergent | 24.0 | 10.0 |
| | $\alpha = 100\gamma_m$ | divergent | divergent | divergent | 61.4 |
| | $\alpha = 0.01\gamma_m$ | divergent | divergent | divergent | 10.1 |

Annotation: $\varepsilon_r = 10^{-12}$, $iter\_max = 100$, $N_t = 80$

Tables 11–14 shows that LIII A-B-$\gamma$ assigned with $\gamma = \gamma_m$ or $\gamma_h$ has quicker convergence speed than LIII A-B-$\alpha$. This means less computational cost for the integrations of $\gamma$-type. Although we have not listed the numerical results, the G-L methods present the same numerical features as the LIIIA-B methods.

These numerical results suggest that the numerical performance for the integrations of $\gamma$-type are harmonious between the numerical accuracy and stability, which means we can obtain higher numerical accuracy as well as better convergence speed by choosing the proper parameter ($\gamma_m$ or $\gamma_h$). Although these tables only list the numerical results for the example of the fast spinning top, a large amount of numerical testing demonstrates that nearly the same numerical phenomenon can be observed for other configuration parameters or initial conditions.

**Table 12** Fast spinning top: iteration number for 3-LIII A-B

| 3-stage | Parameter | $\Delta t = 0.02$ | $\Delta t = 0.01$ | $\Delta t = 0.005$ | $\Delta t = 0.001$ |
|---------|-----------|-------------------|-------------------|--------------------|--------------------|
| $\gamma$-type | $\gamma = \gamma_m$ | 27.7 | 19 | 13.0 | 8.0 |
| | $\gamma = \gamma_h$ | 21.9 | 17.0 | 12.0 | 8.0 |
| $\alpha$-type | $\alpha = \gamma_m$ | 67.8 | divergent | 25.0 | 10.0 |
| | $\alpha = \gamma_h$ | 74.9 | divergent | 24.0 | 10.0 |
| | $\alpha = 100\gamma_m$ | divergent | divergent | divergent | 42.7 |
| | $\alpha = 0.01\gamma_m$ | divergent | 30 | divergent | 9 |

Annotation: $\varepsilon_r = 10^{-12}$, $iter\_max = 100$, $N_t = 80$

**Table 13** Fast spinning top: iteration number for 4-LIII A-B

| 4-stage | Parameter | $\Delta t = 0.02$ | $\Delta t = 0.01$ | $\Delta t = 0.005$ | $\Delta t = 0.001$ |
|---------|-----------|-------------------|-------------------|--------------------|--------------------|
| $\gamma$-type | $\gamma = \gamma_m$ | 21 | 17 | 13.0 | 9.2 |
| | $\gamma = \gamma_h$ | 18 | 15 | 11.7 | 8.0 |
| $\alpha$-type | $\alpha = \gamma_m$ | 43.0 | 26.0 | 16 | 9.0 |
| | $\alpha = \gamma_h$ | 40 | 24 | 15 | 9.0 |
| | $\alpha = 100\gamma_m$ | divergent | divergent | divergent | 54.6 |
| | $\alpha = 0.01\gamma_m$ | 27.8 | 20.0 | 14.1 | 8.7 |

Annotation: $\varepsilon_r = 10^{-12}$, $iter\_max = 100$, $N_t = 80$

**Table 14** Fast spinning top: iteration number for 5-LIII A-B

| 5-stage | Parameter | $\Delta t = 0.02$ | $\Delta t = 0.01$ | $\Delta t = 0.005$ | $\Delta t = 0.001$ |
|---------|-----------|-------------------|-------------------|--------------------|--------------------|
| $\gamma$-type | $\gamma = \gamma_m$ | 19 | 15.9 | 12 | 11.7 |
| | $\gamma = \gamma_h$ | 17.0 | 14.0 | 11.0 | 14.1 |
| $\alpha$-type | $\alpha = \gamma_m$ | 29.0 | 21 | 15 | 12.5 |
| | $\alpha = \gamma_h$ | 28.0 | 20.0 | 14 | 14.9 |
| | $\alpha = 100\gamma_m$ | divergent | divergent | divergent | 55.9 |
| | $\alpha = 0.01\gamma_m$ | 23.8 | 18.1 | 12.8 | 10.2 |

Annotation: $\varepsilon_r = 10^{-12}$, $iter\_max = 100$, $N_t = 80$

## 6 Conclusion

The inertia representations of $\gamma$-type and $\alpha$-type have been developed for the simulation of the quaternion-based rigid body dynamics. The two inertia representations formally lead to different formulations of Hamilton's equations, which are theoretically equivalent if the constraint $q^T q = 1$ is satisfied exactly. However, error estimation demonstrates that the two kinds of inertia representations are different due to the discretization and suggests that the parameter $\gamma$ can be used to optimize the numerical performance of the integrations in simulation.

The implicit midpoint scheme (IMS), the energy–momentum conserving scheme (EMS), and two types of Gauss SPARK methods (G-L and LIIIA-B) are derived to investigate the two inertia representations in simulations. The numerical results show that the parameters

$\alpha$ and $\gamma$ can greatly influence the numerical performance of the integrations in simulation and the numerical influences are the result of the comprehensive effect of the discretization scheme, the inertia representations and their combinations. To be specific, IMS-$\gamma$, EMS-$\gamma$ and G-L-$\gamma$ present a linear relation between numerical errors and the parameter $\gamma$ regardless of the convergence order of the integration, whereas the parameter $\alpha$ has no influence on the numerical accuracy for these integrations of $\alpha$-type; s-LIIIA-B-$\alpha$ and s-LIIIA-B-$\gamma$ are of the same numerical accuracy in simulation, if the parameters $\alpha$ and $\gamma$ are assigned with the same value; all the integrations of $\gamma$-type can present quicker convergence speed and better stability than those of $\alpha$-type. A large amount of numerical testing demonstrates that the two values $\gamma_m$ and $\gamma_h$, referring to three principal moments, can be considered as two reasonable values of $\gamma$, with which the integrations of $\gamma$-type can present better numerical accuracy, convergence speed and stability for these integrations.

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

## Appendix A:  The Gauss SPARK methods in unified form

### A.1  A unified formulation of SPARK methods

In Sect. 4.2.1, we present two types of Gauss SPARK methods based on the discontinuous collocation type methods. It is more efficient in computation to derive a unified formulation of SPARK methods in which the discretization schemes are formulated with Butcher style tableaux of SPARK coefficients.

One step of an $(s, p)$-SPARK method applied to the Hamiltonian system of (12) with consistent initial values $(\boldsymbol{q}_k, \boldsymbol{p}_k)$ at $t_k$ and step-size $\Delta t$ is given as follows [24, 28]:

$$\boldsymbol{k}_i = H_{\boldsymbol{p}}(\boldsymbol{Q}_i, \boldsymbol{P}_i), \quad i = 1, \ldots, s,$$

$$\hat{\boldsymbol{k}}_i = -H_{\boldsymbol{q}}(\boldsymbol{Q}_i, \boldsymbol{P}_i), \quad i = 1, \ldots, s, \tag{A.1}$$

$$\tilde{\boldsymbol{k}}_i = -g_{\boldsymbol{q}}(\tilde{\boldsymbol{Q}}_i)\Lambda_i, \quad i = 0, \ldots, p,$$

$$g(\tilde{\boldsymbol{Q}}_i) = 0, \quad i = 0, 1, \ldots, p, \tag{A.2}$$

$$\boldsymbol{q}_{k+1} = \boldsymbol{q}_k + \Delta t \sum_{j=1}^{s} \boldsymbol{k}_j b_j, \tag{A.3}$$

$$g(\boldsymbol{q}_{k+1}) = 0, \tag{A.4}$$

$$\boldsymbol{p}_{k+1} = \boldsymbol{p}_k + \Delta t \sum_{j=1}^{s} b_j \hat{\boldsymbol{k}}_j + \Delta t \sum_{j=0}^{p} \tilde{b}_j \tilde{\boldsymbol{k}}_j, \tag{A.5}$$

$$g_{\boldsymbol{q}}(\boldsymbol{q}_{k+1}) H_{\boldsymbol{p}}(\boldsymbol{q}_{k+1}, \boldsymbol{p}_{k+1}) = 0, \tag{A.6}$$

where $\Lambda_i$ for $i = 0, 1, \ldots, p$ is the Lagrange multiplier. $\boldsymbol{Q}_i$, $\tilde{\boldsymbol{Q}}_i$ and $\boldsymbol{P}_i$ denote quantities discretized as the inner points, defined by

$$\boldsymbol{Q}_i = \boldsymbol{q}_k + \Delta t \sum_{j=1}^{s} \boldsymbol{k}_j a_{ij}, \quad i = 1, \ldots, s,$$

$$\tilde{\boldsymbol{Q}}_i = \boldsymbol{q}_k + \Delta t \sum_{j=1}^{s} \boldsymbol{k}_j \bar{a}_{ij}, \quad i = 0, \ldots, p, \tag{A.7}$$

$$\boldsymbol{P}_i = \boldsymbol{p}_k + \Delta t \sum_{j=1}^{s} \hat{\boldsymbol{k}}_j \hat{a}_{ij} + \Delta t \sum_{j=0}^{p} \tilde{\boldsymbol{k}}_j \tilde{a}_{ij}, \quad i = 1, \ldots, s.$$

The coefficients $(b_j, c_j)_{j=1}^{s}$ and $(\tilde{b}_j, \tilde{c}_j)_{j=0}^{p}$ are generally two distinct quadrature formulas and the SPARK coefficients can be expressed in the form of Butcher style tableaux:

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

To ensure the existence and uniqueness of the SPARK solution [24, 28], we should assume $\bar{a}_{0j} = 0$ and $\bar{a}_{pj} = b_j$ for $j = 1, \ldots, s$, which imply that $\tilde{\boldsymbol{Q}}_0 = \boldsymbol{q}_k$ and $\tilde{\boldsymbol{Q}}_p = \boldsymbol{q}_{k+1}$. Hence Eqs. (A.2), (A.4) and (A.6) give $p + 1$ independent constraints to solve $p + 1$ Lagrange's multipliers (i.e., $\Lambda_i$, $i = 0, \ldots, p$).

We are especially interested in the SPARK methods whose coefficients are the weights and nodes of Gauss quadrature, to have an optimal order of convergence. The Gauss SPARK coefficients satisfy the following conditions [24, 28]:

$$B(s): \quad \sum_{i=1}^{s} b_i c_i^{k-1} = \frac{1}{k}, \quad k = 1, \ldots, 2s,$$

$$C(s): \quad \sum_{j=1}^{s} a_{ij} c_j^{k-1} = \frac{c_i^k}{k}, \quad k = 1, \ldots, s, \ i = 1, \ldots, s, \tag{A.8}$$

$$\hat{C}(s): \quad \sum_{j=1}^{s} \hat{a}_{ij} \hat{c}_j^{k-1} = \frac{\hat{c}_i^k}{k}, \quad k = 1, \ldots, s, \ i = 1, \ldots, s, \tag{A.9}$$

$$\tilde{B}(p): \quad \sum_{i=0}^{p} \tilde{b}_i \tilde{c}_i^{k-1} = \frac{1}{k}, \quad k = 1, \ldots, 2p, \tag{A.10}$$

$$\tilde{C}(p): \quad \sum_{j=0}^{p} \tilde{a}_{ij} \tilde{c}_j^{k-1} = \frac{c_i^k}{k}, \quad k = 1, \ldots, p+1, i = 1, \ldots, s, \tag{A.11}$$

and

$$\bar{C}(s): \quad \sum_{j=1}^{s} \bar{a}_{ij} c_j^{k-1} = \frac{\tilde{c}_i^k}{k}, \quad k = 1, \ldots, s, \ i = 0, \ldots, p. \tag{A.12}$$

The two types of Gauss SPARK methods concerned in this paper can be obtained by the following conditions.

**s-stage Lobatto IIIA-B method**  *Let* $p = s - 1$, $\bar{a}_{(i-1)j} = a_{ij}$, $\tilde{a}_{i(j-1)} = \hat{a}_{ij}$, $b_j = \hat{b}_j = \tilde{b}_{j-1}$, $\tilde{c}_{i-1} = c_i$, *for* $i = 1, \ldots, s$, $j = 1, \ldots, s$, *and* $c_1, \ldots, c_s$, *be the $s$ nodes of the Lobatto quadrature. The Lobatto IIIA-B SPARK method is defined with the coefficients* $a_{ij}$, $b_j$ *and* $\hat{a}_{ij}$ *determined by*

$$\hat{a}_{i1} = \hat{b}_1, \qquad \hat{a}_{is} = 0 \quad \text{for } i = 1, \ldots, s \tag{A.13}$$

*and*

$$B(s), \quad C(s) \quad \text{and} \quad \hat{C}(s-2) \tag{A.14}$$

*with the conditions $B(s)$, $C(s)$ and $\hat{C}(s)$ defined by* (A.8) *and* (A.9).

**s-stage Gauss Lobatto method**  *Let* $p = s$, $\hat{a}_{ij} = a_{ij}$, $b_j = \hat{b}_j$, *for* $i = 1, \ldots, s$, *and* $c_1, \ldots, c_s$ *be the $s$ nodes of the Gauss quadrature and* $\tilde{c}_1, \ldots, \tilde{c}_s$ *be the $s + 1$ nodes of the Lobatto quadrature. The Gauss–Lobatto SPARK method is defined with the coefficients* $a_{ij}$, $b_j$, $\tilde{a}_{ij}$, $\tilde{b}_i$ *and* $\bar{a}_{ij}$, *determined by*

$$\tilde{a}_{i0} = \tilde{b}_0, \qquad \tilde{a}_{is} = 0 \quad \text{for } i = 1, \ldots, s \tag{A.15}$$

*and*

$$B(s), \quad C(s), \quad \tilde{B}(s), \quad \tilde{C}(s-2) \quad \text{and} \quad \bar{C}(s) \tag{A.16}$$

*with the conditions $B(s)$, $C(s)$, $\tilde{B}(p)$, $\tilde{C}(p)$ and $\bar{C}(s)$ defined by* (A.8)–(A.12).

## A.2 Butcher style tableaux of Gauss SPARK coefficients

SPARK coefficients can be calculated by substituting nodes of quadrature and the solution conditions presented by Lobatto IIIA-B or Gauss–Lobatto SPARK methods into (A.8)–(A.12). The Lobatto nodes of quadrature $c_1, c_2, \ldots, c_s$ are the zeros of

$$\frac{d^{s-2}}{dx^{s-2}}\left(x^{s-1}(x-1)^{s-1}\right). \tag{A.17}$$

The Gauss nodes of quadrature $c_1, c_2, \ldots, c_s$ are the zeros of the $s$th shifted Legendre polynomial

$$\frac{d^s}{dx^s}\left(x^s(x-1)^s\right). \tag{A.18}$$

### A.2.1 SPARK coefficients for s-stage -Lobatto IIIA-B SPARK methods

The 2-stage and 3-stage Lobatto IIIA-B SPARK methods correspond to the following Butcher style tableaux of SPARK coefficients:

$$
\begin{array}{c|cc}
0 & 0 & 0 \\
1 & 1/2 & 1/2 \\
\hline
A & 1/2 & 1/2
\end{array}
\qquad
\begin{array}{|cc}
1/2 & 0 \\
1/2 & 0 \\
\hline
\hat{A} & 1/2 & 1/2
\end{array}
\qquad
\begin{array}{|cc}
1/2 & 0 \\
1/2 & 0 \\
\hline
\tilde{A} & 1/2 & 1/2
\end{array}
\qquad
\begin{array}{c|cc}
0 & 0 & 0 \\
1 & 1/2 & 1/2 \\
\hline
\bar{A} & &
\end{array}
$$

and

$$
\begin{array}{c|ccc}
0 & 0 & 0 & 0 \\
1/2 & 5/24 & 1/3 & -1/24 \\
1 & 1/6 & 2/3 & 1/6 \\
\hline
A & 1/6 & 2/3 & 1/6
\end{array}
\qquad
\begin{array}{|ccc}
1/6 & -1/6 & 0 \\
1/6 & 1/3 & 0 \\
1/6 & 5/6 & 0 \\
\hline
\hat{A} & 1/6 & 2/3 & 1/6
\end{array}
$$

$$
\begin{array}{ccc|c}
1/6 & -1/6 & 0 & \\
1/6 & 1/3 & 0 & \\
1/6 & 5/6 & 0 & \\
\hline
\tilde{A} \; 1/6 & 2/3 & 1/6 &
\end{array}
\qquad
\begin{array}{c|ccc}
0 & 0 & 0 & 0 \\
1/2 & 5/24 & 1/3 & -1/24 \\
1 & 1/6 & 2/3 & 1/6 \\
\hline
\bar{A} &
\end{array}
$$

### A.2.2 SPARK coefficients for s-stage Gauss–Lobatto SPARK methods

The 1-stage and 2-stage Gauss–Lobatto SPARK methods correspond to the following Butcher style tableaux of SPARK coefficients:

$$
\begin{array}{c|c}
1/2 & 1/2 \\
\hline
A & 1
\end{array}
\qquad
\begin{array}{c|c}
 & 1/2 \\
\hline
\hat{A} & 1
\end{array}
\qquad
\begin{array}{c|cc}
 & 1/2 & 0 \\
\hline
\tilde{A} & 1/2 & 1/2
\end{array}
\qquad
\begin{array}{c|c}
0 & 0 \\
1 & 1 \\
\hline
\bar{A} &
\end{array}
$$

and

$$
\begin{array}{c|cc}
1/2 - \sqrt{3}/6 & 1/4 & 1/4 - \sqrt{3}/6 \\
1/2 + \sqrt{3}/6 & 1/4 + \sqrt{3}/6 & 1/4 \\
\hline
A & 1/2 & 1/2
\end{array}
\qquad
\begin{array}{c|cc}
 & 1/4 & 1/4 - \sqrt{3}/6 \\
 & 1/4 + \sqrt{3}/6 & 1/4 \\
\hline
\hat{A} & 1/2 & 1/2
\end{array}
$$

$$
\begin{array}{c|ccc}
 & 1/6 & 1/3 - \sqrt{3}/6 & 0 \\
 & 1/6 & 1/3 + \sqrt{3}/6 & 0 \\
\hline
\tilde{A} \; 1/6 & 2/3 & 1/6 &
\end{array}
\qquad
\begin{array}{c|cc}
0 & 0 & 0 \\
1/2 & 1/4 + \sqrt{3}/8 & 1/4 - \sqrt{3}/8 \\
1 & 1/2 & 1/2 \\
\hline
\bar{A} &
\end{array}
$$

## Appendix B: Algorithms in pseudocode formats

### B.3 Implicit midpoint scheme and energy–momentum conserving scheme

The implicit midpoint scheme and the energy–momentum conserving scheme presented by (53) and (61) can be summarized in pseudocode format in Table 15.

**Table 15** State-space integration algorithm of order 2

| | |
|---|---|
| 1) | Initial condition: $\boldsymbol{q}_0$, $\boldsymbol{p}_0$, $\lambda = 0$, $k = 1$ |
| 2) | Prediction step: $\boldsymbol{g}_k = \boldsymbol{g}_{k-1}$, $\boldsymbol{g}^T = [\boldsymbol{q}^T, \boldsymbol{p}^T, \lambda]$. |
| 3) | Residual calculation: $\boldsymbol{r}^T = [\boldsymbol{r}_{\boldsymbol{q}}^T, \boldsymbol{r}_{\boldsymbol{p}}^T, \boldsymbol{q}_k^T \boldsymbol{q}_k - 1]$ |
| | Implicit midpoint scheme (IMS): |
| | $\begin{cases} \boldsymbol{r}_{\boldsymbol{q}} = \Delta \boldsymbol{q}/\Delta t - H_{\boldsymbol{p}}(\bar{\boldsymbol{q}}, \bar{\boldsymbol{p}}) \\ \boldsymbol{r}_{\boldsymbol{p}} = \Delta \boldsymbol{p}/\Delta t + H_{\boldsymbol{q}}(\bar{\boldsymbol{q}}, \bar{\boldsymbol{p}}) + 2\lambda \bar{\boldsymbol{q}} \end{cases}$ |
| | Energy–momentum conserving scheme (EMS): |
| | $\begin{cases} \boldsymbol{r}_{\boldsymbol{q}} = \Delta \boldsymbol{q}/\Delta t - \bar{\nabla}_{\boldsymbol{p}} H(\boldsymbol{q}_{k-1}, \boldsymbol{p}_{k-1}, \boldsymbol{q}_k, \boldsymbol{p}_k) \\ \boldsymbol{r}_{\boldsymbol{p}} = \Delta \boldsymbol{p}/\Delta t + \bar{\nabla}_{\boldsymbol{q}} H(\boldsymbol{q}_{k-1}, \boldsymbol{p}_{k-1}, \boldsymbol{q}_k, \boldsymbol{p}_k) + 2\lambda \bar{\boldsymbol{q}} \end{cases}$ |
| 4) | Update incremental rotation parameters: |
| | $\delta \boldsymbol{g} = -\boldsymbol{K}_J^{-1} \boldsymbol{r}$, $\boldsymbol{g}_k = \boldsymbol{g}_k + \delta \boldsymbol{g}$, |
| | where $\boldsymbol{K}_J = \partial \boldsymbol{r}/\partial \boldsymbol{g}_k$ and if $\|\boldsymbol{r}\| > \varepsilon_r$, repeat from 3). |
| 5) | $k = k + 1$, return to 2), or stop. |

| **Table 16** Algorithm for Gauss SPARK methods | |
|---|---|
| 1) | Initial condition: $q_0$, $p_0$, $x_{-1} = 0$, $k = 0$ |
| 2) | Prediction step: $x_k = x_{k-1}$ where $x$ is defined by (B.2) |
| 3) | Residual calculation-1: $r^T = [r_q^T, r_p^T, r_c^T]$ |
| 4) | Update incremental parameters: |
| | $\delta x = -K_J^{-1} r$, $x_k = x_k + \delta x$, |
| | where $K_J = \partial r/\partial x_k$ and if $\|r\| > \varepsilon_r$, repeat from 3). |
| 5) | Update the configurations: $q_{k+1} = q_k + \Delta t \sum_{j=1}^s k_j b_j$, $\Lambda_{\tilde{s}} = 0$ |
| 6) | Residual calculation-2: |
| | $p_{k+1} = p_k + \Delta t \sum_{j=1}^s b_j \hat{k}_j + \Delta t \sum_{j=0}^p \tilde{b}_j \tilde{k}_j$, $\tilde{r} = p_{k+1}^T q_{k+1}$ |
| 7) | Update incremental parameter: $\delta \Lambda = -\tilde{r}/K_\Lambda$, $\Lambda_p = \Lambda_p + \delta \Lambda$ |
| | where $K_\Lambda = \partial \tilde{r}/\partial \Lambda_p$ and if $\|r\| > \tilde{\varepsilon}_r$, repeat from 6). |
| 8) | $k = k + 1$, return to 2), or stop |

## B.4 Gauss SPARK methods

Consider the SPARK method presented by (A.1)–(A.6). We can define the residual vectors in the following form:

$$
r = \begin{bmatrix} r_q \\ r_p \\ r_c \end{bmatrix} \qquad
\begin{aligned}
r_q &= [r_{q,1}^T, \ldots, r_{q,s}^T]^T & r_{q,i} &= k_i - H_p(Q_i, P_i), & i &= 1, \ldots, s, \\
r_q &= [r_{q,1}^T, \ldots, r_{q,s}^T]^T & r_{p,i} &= \hat{k}_i + H_q(Q_i, P_i), & i &= 1, \ldots, s, \\
r_c &= [r_{c,1}^T, \ldots, r_{c,p}^T]^T & r_{c,i} &= g(\tilde{Q}_i), & i &= 1, \ldots, p,
\end{aligned}
$$
(B.1)

where $Q_i$, $P_i$ and $\tilde{Q}_i$ are defined by (A.7) and $\tilde{k}_i = -g_q(\tilde{Q}_i)\Lambda_i$. Define the unknowns as

$$
x = \begin{bmatrix} k_1^T & k_2^T & \cdots & k_s^T & \hat{k}_1^T & \hat{k}_2^T & \cdots & \hat{k}_s^T & \Lambda_0 & \Lambda_1 & \cdots & \Lambda_{p-1} \end{bmatrix}^T.
$$
(B.2)

Then the Gauss SPARK methods can be summarized uniformly in pseudocode format as in Table 16.

## References

1. Goldstein, H., Poole, C.P., Safko, J.L.: Classical Mechanics, 3rd edn. Addison–Wesley, New York (2001)
2. Nikravesh, P.E., Chung, I.S.: Application of Euler parameters to the dynamic analysis of three-dimensional constrained mechanical systems. J. Mech. Des. **104**, 785–791 (1982)
3. Nikravesh, P.E., Kwon, O.K., Wehage, R.A.: Euler parameters in computational kinematics and dynamics, part 2. J. Mech. Transm. Autom. Des. **107**, 366–369 (1985)
4. Nikravesh, P.E.: Computer-Aided Analysis of Mechanical Systems. Prentice-Hall, New York (1988)
5. Haug, E.J.: Computer Aided Kinematics and Dynamics of Mechanical Systems, vol. 1: Basic Methods. Allyn and Bacon, Boston (1989)
6. Vadali, S.R.: On the Euler parameter constraint. J. Astronaut. Sci. **36**, 259–265 (1988)
7. Chou, J.C.K.: Quaternion kinematic and dynamic differential equations. IEEE Trans. Robot. Autom. **8**, 53–64 (1992)
8. Morton, H.S. Jr.: Hamiltonian and Lagrangian formulations of rigid-body rotational dynamics based on the Euler parameters. J. Astronaut. Sci. **41**, 569–591 (1993)
9. Shivarama, R., Fahrenthold, E.P.: Hamilton's equations with Euler parameters for rigid body dynamics modeling. ASME J. Dyn. Syst. Meas. Control-Trans. **126**, 124–130 (2004)
10. Sherif, K., Nachbagauer, K., Steiner, W.: On the rotational equations of motion in rigid body dynamics when using Euler parameters. Nonlinear Dyn. **81**, 343–352 (2015)
11. Betsch, P., Siebert, R.: Rigid body dynamics in terms of quaternions: Hamiltonian formulation and conserving numerical integration. Int. J. Numer. Methods Eng. **79**, 444–473 (2009)

12. Moller, M., Glocker, C.: Rigid body dynamics with a scalable body, quaternions and perfect constraints. Multibody Syst. Dyn. **27**, 437–454 (2012)
13. O'Reilly, O.M., Varadi, P.C.: Hoberman's sphere, Euler parameters and Lagrange's equations. J. Elast. **56**, 171–180 (1999)
14. Udwadia, F.E., Schutte, A.D.: An alternative derivation of the quaternion equations of motion for rigid-body rotational dynamics. ASME J. Appl. Mech.-Trans. **77**, 4 (2010)
15. Miller, T.F., Eleftheriou, M., Pattnaik, P., Ndirango, A., Newns, D., Martyna, G.J.: Symplectic quaternion scheme for biophysical molecular dynamics. J. Chem. Phys. **116**, 8649–8659 (2002)
16. Nielsen, M.B., Krenk, S.: Conservative integration of rigid body motion by quaternion parameters with implicit constraints. Int. J. Numer. Methods Eng. **92**, 734–752 (2012)
17. Simo, J.C., Wong, K.K.: Unconditionally stable algorithms for rigid body dynamics that exactly preserve energy and momentum. Int. J. Numer. Methods Eng. **31**, 19–52 (1991)
18. Wendlandt, J.M., Marsden, J.E.: Mechanical integrators derived from a discrete variational principle. Physica D **106**, 223–246 (1997)
19. Manchester, Z.R., Peck, M.A.: Quaternion variational integrators for spacecraft dynamics. J. Guid. Control Dyn. **39**, 69–76 (2016)
20. Leitz, T., Leyendecker, S.: Galerkin Lie-group variational integrators based on unit quaternion interpolation. Comput. Methods Appl. Mech. Eng. **338**, 333–361 (2018)
21. Shabana, A.A.: Euler parameters kinetic singularity. Proc. Inst. Mech. Eng., Proc., Part K, J. Multi-Body Dyn. **228**, 307–313 (2014)
22. Celledoni, E., Safstrom, N.: A Hamiltonian and multi-Hamiltonian formulation of a rod model using quaternions. Comput. Methods Appl. Mech. Eng. **199**, 2813–2819 (2010)
23. Xu, X.M., Zhong, W.X.: On the numerical influences of inertia representation for rigid body dynamics in terms of unit quaternion. ASME J. Appl. Mech.-Trans. **83**, 11 (2016)
24. Jay, L.O.: Specialized partitioned additive Runge–Kutta methods for systems of overdetermined DAEs with holonomic constraints. SIAM J. Numer. Anal. **45**, 1814–1842 (2007)
25. Hairer, E., Lubich, C., Wanner, G.: Geometric Numerical Integration. Structure-Preserving Algorithms for Ordinary Differential Equations, 2nd edn. (2006)
26. Krenk, S., Nielsen, M.B.: Conservative rigid body dynamics by convected base vectors with implicit constraints. Comput. Methods Appl. Mech. Eng. **269**, 437–453 (2014)
27. Gonzalez, O.: Exact energy and momentum conserving algorithms for general models in nonlinear elasticity. Comput. Methods Appl. Mech. Eng. **190**, 1763–1783 (2000)
28. Small, S.J.: Runge–Kutta Type Methods for Differential-Algebraic Equations in Mechanics. Dissertations & Theses – Gradworks (2011)
29. Jay, L.: Symplectic partitioned Runge–Kutta methods for constrained Hamiltonian systems. SIAM J. Numer. Anal. **33**, 368–387 (1996)