

# Conservative single-step time-integration schemes with higher-order accuracy for multi-particle dynamics with local two-point potentials

E. Graham<sup>a</sup>, G. Jelenić<sup>b,\*</sup>

<sup>a</sup> UGS, Parker's House, 46 Regent Street, Cambridge CB2 1DP, United Kingdom

<sup>b</sup> Department of Civil Engineering, University of Rijeka, Viktora Cara Emina 5, 51000 Rijeka-Fiume, Croatia

Received 4 March 2005; received in revised form 7 July 2005; accepted 8 July 2005

## Abstract

A general framework for algorithms that conserve linear and angular momenta for problems of multi-particle mechanics is presented. Conditions for energy conservation are derived, and different manners in which this may be achieved are discussed. A detailed examination of the relative equilibrium states is carried out, and conditions under which algorithms preserve these states are given; in particular, algorithms can be designed to capture the *exact* solutions of relative equilibrium problems, although these algorithms are unlikely to be energy-conserving. Following on from the approach proposed by Argyris et al. [J.H. Argyris, P.C. Dunne, T. Angelopoulos, Dynamic response by large step integration, Earthquake Engrg. Struct. Dynam. 2 (1973) 185–203], the local accuracy characteristics of algorithms are investigated thoroughly, and it is shown that there is no limit to the order of accuracy that can be achieved by algorithms in this framework, even for problems with time-dependent forces. No extra stages of calculation or additional degrees of freedom are required to be present, although the sparsity of the resulting system of equations is compromised. A few examples of new algorithms are given, and their properties verified on some model problems.

© 2005 Elsevier B.V. All rights reserved.

**Keywords:** Higher-order accuracy; Conservation; Relative equilibria; Time reversibility

## 1. Introduction

In this work, we explore the possibilities for higher-order accuracy in the design of algorithms to solve non-linear, dynamic problems involving large displacements. The exact solutions to these problems are generally unavailable in closed form, existing only as power series in  $\Delta t$ . Our sphere of interest concerns undamped systems for which the forces are derived from a single scalar function, known as *monogenic* systems [2]. We are particularly interested in *Hamiltonian* systems, for which the external forces are *conservative*; in these cases, the Hamiltonian function  $H$  defining the total energy of the system is *constant* throughout the motion.

Fundamental to the accuracy of an algorithm is the concept of *stability*. There are many different definitions of algorithmic stability that exist; e.g., [3–6] and many more besides. Loosely speaking, each relates to whether or not an algorithm produces a solution that is bounded as the total time  $T \rightarrow \infty$  (for some fixed  $\Delta t$ ), assuming the exact solution to be bounded also. In any quest for accuracy, stability is essential; therefore we seek to ensure stability first, and subsequently concentrate on improving the accuracy of an algorithm. As a consequence, we will deal exclusively with *implicit* schemes,

\* Corresponding author. Tel.: +385 51 352 114; fax: +385 51 332 816.  
E-mail address: [gordan@gradri.hr](mailto:gordan@gradri.hr) (G. Jelenić).

since explicit algorithms are known to be only conditionally stable [3,7]. For monogenic systems, we can formulate the definition of stability in terms of the change in energy within a single time-step, as was first done by Belytschko and Schoeberle [5]. Specifically, the internal energy of the system must not be allowed to increase beyond the amount of work done by the external forces during the time-step. It is possible to design algorithms that satisfy this condition for all time-step sizes; such algorithms are thus *unconditionally stable* for general non-linear problems.

For Hamiltonian systems, it is also possible to design algorithms that preserve the symplectic nature of the system: a detailed discussion of the construction of symplectic algorithms is given by Simo and Tarnow in Appendix II of [8], and also by Marsden and West [9]. It has been shown by Zhong and Marsden, however, that an algorithm cannot be symplectic and simultaneously conserve energy for a general non-integrable system, assuming the time-step size to be constant [10]. Also, algorithms which are spectrally stable and dissipate energy for linear problems cannot be symplectic, as demonstrated by Simo and Tarnow [8]. A choice must therefore be made, between the properties of symplecticity and energy conservation or dissipation, from the outset. Given the fact that a definition of stability is readily available in terms of energy growth, we elect to design algorithms based on energy criteria. In support of this decision, Simo and co-workers have shown that for stiff problems of non-linear elasticity, energy-conserving algorithms tend to perform better than symplectic algorithms [8,11,12].

Unconditional stability alone is not enough to ensure the overall accuracy of an algorithm. An extreme example of this was given by Ortiz [13], where a convergent, energy-conserving algorithm is shown to give completely inaccurate results after only a small number of time-steps. In the case of *systems with symmetries*, however, we have more information. These systems furnish two other constants of motion, namely the total linear momentum  $\mathcal{L}$  and the total angular momentum  $\mathcal{J}$ , each of which gives information about the qualitative (as well as quantitative) nature of the solution. Algorithms can also be designed to conserve these momenta for such systems, in conjunction with energy conservation, with a view to achieving better accuracy. The importance of angular momentum conservation in regard to accuracy was noted by Betsch and Steinmann [14], and classic examples of such algorithms are given by Simo and co-workers [8,15].

As a consequence of having constant linear and angular momenta, an additional property of systems with symmetries is the existence of families of fully integrable solutions, each induced by a particular combination of initial conditions. These are known as *relative equilibrium states* (or *steady states*), and give further information about the stability of the system: see [16,17] for a detailed account. Algorithms that conserve momenta can be designed to preserve these relative equilibrium states (when the initial conditions arise): those that do give solutions to a steady-state problem that physically resemble the exact solution, and thus may have enhanced stability and accuracy properties for problems of *approximately* steady-state motion. An analysis of two popular time-integration schemes in this regard is given by Gonzalez and Simo [12] and further discussion on the importance of preservation of relative equilibrium orbits is provided by Armero and Romero [18]. Examples of algorithms designed to preserve relative equilibria include the energy–momentum algorithm of Simo and Tarnow [8], and subsequent algorithms that dissipate energy by Armero and co-workers [18–20].

One further property of dynamical systems in physics that we touch upon briefly is that of time-reversibility [21], which is closely related to the uniqueness of a dynamic response. In the discrete case, however, it is not certain that, at any given point on the solution, negating the time-step would recover the solution given at the previous point in time. Algorithms that guarantee this are described as *time-reversible*. An early citation of the importance of this property in the engineering context is due to Argyris et al. [1], and further examples of such algorithms relevant to our work are the energy-conserving algorithm proposed independently by Simo and Gonzalez [11] and Reich [22], as well as the symplectic mid-point rule [8,15,23,24] and the so-called assumed distance method [25].

It is widely accepted that an algorithm should be at least second-order accurate (e.g. [4]), and the energy–momentum algorithms of Simo et al. mentioned earlier all satisfy this requirement. Various ways to increase the order of accuracy have been proposed. These include *composition methods*, whereby greater accuracy is achieved by computing intermediate results at additional points within a single time-step; example algorithms include those given independently by Yoshida [21], Forest [26] and Tarnow and Simo [27]. A disadvantage is that the procedure involves stepping backwards in time, using a larger time-step size than the original algorithm; this makes the principle less attractive for algorithms that are not time-reversible, and increases the risk of instability or divergence during the non-linear iteration process. Another approach to enhancing accuracy can be taken by discretising the equations of motion using *finite elements in time*, where the accuracy can be prescribed by the degree of the polynomial basis functions chosen; example algorithms include those of Betsch and Steinmann [28] for non-linear dynamics. These schemes bear close resemblance to Gauss Runge–Kutta methods, as described in the Appendix of [28].

Both of these strategies to improve accuracy entail additional computation cost, due to the calculation of intermediate results or the presence of extra degrees of freedom in the temporal domain of the problem. A third approach, aimed at avoiding this additional cost, is based on *Taylor series expansions* of the state variables; for linear dynamics, this method is equivalent to using Padé approximations to the exact solution. Early work was done along these lines for non-linear analysis by Argyris et al. [1,29], who presented *arbitrarily* accurate algorithms that are time-reversible, although not conservative. They were followed by LaBudde and Greenspan, who produced arbitrarily accurate schemes that also conserve

energy and angular momentum for a central-force problem [30], and similar energy-conserving schemes for the  $N$ -body problem [31]. These algorithms are not time-reversible, however, and do not preserve the orbits of relative equilibria when higher than second-order accurate. In both cases, small time-step sizes were necessary to ensure convergence of the non-linear solution procedure.

To achieve our goal of designing algorithms with desirable accuracy characteristics, we will follow this approach of using Taylor series expansions. Firstly, however, we aim to equip the algorithms with the properties of energy and momentum conservation, preservation of relative equilibrium states and time reversibility, when the relevant physical principles apply. Subject to retaining these properties of the system, we seek to maximise the order of accuracy of the algorithms we design, and present criteria by which it may be achieved.

In [32], we developed a framework for designing higher-order accurate, conservative algorithms for the central-force problem that did not entail additional computational effort in the manner described above. This work extends [32] to the realm of multi-body problems, and discusses the extra complexity involved. Specifically, we deal with non-linear elasticity of multi-element truss structures, with the view that any progress made in this area will highlight avenues of research for time integration involving problems of two- and three-dimensional continua. Also, specific to elastic structures are certain types of strain-energy potential, which in principle may be substituted with a different potential and thus make the present framework entirely applicable to the problems of multi-particle dynamics. The background theory is presented in Section 2, and in Section 3 we establish a general framework in which higher-order algorithms with the required properties exist, that does not entail intermediate calculations or additional degrees of freedom. The conservation and accuracy criteria are derived in Sections 4 and 5, and in Section 6 we take a brief look at some new and existing algorithms that fit into the framework, which are tested on a couple of model problems in Section 7.

## 2. Equations of motion

Let  $\mathbf{r}(\mathbf{X}, t), \dot{\mathbf{r}}(\mathbf{X}, t) \in \mathbb{R}^3$  be the position and velocity at time  $t$  of a point  $x \in \mathcal{B}$ , where  $\mathcal{B}_0 \subset \mathbb{R}^3$  is the initial volume of the continuum; let  $\rho(\mathbf{X}, t)$  represent the current density of the material, with  $\rho_0(\mathbf{X}) \equiv \rho(\mathbf{X}, 0)$  the initial density; and let  $\mathbf{X}$  be the position vector of  $x$  in the reference configuration. Using the conservation of mass, we have the following definitions (e.g. [33]):

$$M := \int_{\mathcal{B}_0} \rho_0 dV, \quad (2.1)$$

$$\mathcal{L} := \int_{\mathcal{B}_0} \rho_0 \dot{\mathbf{r}} dV, \quad (2.2)$$

$$\mathcal{J} := \int_{\mathcal{B}_0} \mathbf{r} \times \rho_0 \dot{\mathbf{r}} dV \quad \text{and} \quad (2.3)$$

$$T := \frac{1}{2} \int_{\mathcal{B}_0} \rho_0 \dot{\mathbf{r}} \cdot \dot{\mathbf{r}} dV, \quad (2.4)$$

for the mass  $M$ , linear momentum  $\mathcal{L}$ , angular momentum  $\mathcal{J}$  and the kinetic energy  $T$  of the continuum. We also define  $\Phi$  to be the total potential energy function from which the forces are derived, which is independent of velocities. Furthermore, we will require knowledge of the *centre of mass* of the continuum, which is defined as

$$\mathbf{r}^c := \frac{1}{M} \int_{\mathcal{B}_0} \rho_0 \mathbf{r} dV \quad (2.5)$$

(using the conservation of mass).

Introducing *spatial* discretisation of the position vector, we have

$$\mathbf{r}(\mathbf{X}, t) = \mathbf{N}(\mathbf{X})\mathbf{R}(t), \quad (2.6)$$

where  $\mathbf{N}(\mathbf{X}) \in \mathbb{R}^{3 \times 3N}$  is a matrix of shape functions and

$$\mathbf{R}(t) = \langle \mathbf{r}^1(t) \quad \dots \quad \mathbf{r}^N(t) \rangle \in \mathbb{R}^{3N}$$

a vector of nodal positions, with  $N$  the number of spatial nodes used in the discretisation. (Here and throughout the paper we use the notation  $\langle \cdot \rangle$  to describe a column vector thereby enabling such expressions to be contained within a line of text.)

Thus  $\ddot{\mathbf{r}}(\mathbf{X}, t) = \mathbf{N}(\mathbf{X})\ddot{\mathbf{R}}(t)$  and, by defining the mass matrix  $\mathbf{M} \in \mathbb{R}^{3N \times 3N}$  to be

$$\mathbf{M} := \int_{\mathcal{B}_0} \rho_0 \mathbf{N}^T \mathbf{N} dV, \quad (2.7)$$

we can express the general semi-discrete equation for undamped motion of a continuum as

$$\mathbf{M}\ddot{\mathbf{R}} + \nabla_{\mathbf{R}}\Phi = \mathbf{0}. \quad (2.8)$$

For the problems we consider in this work, the potential function  $\Phi(\mathbf{R}, t)$  can be written as

$$\Phi(\mathbf{R}, t) = \phi(\mathbf{R}) - U(\mathbf{R}, t), \tag{2.9}$$

where  $\phi$  denotes the strain energy of the system, and  $U$  is the function from which the externally applied force is derived. Thus we have

$$\nabla_{\mathbf{R}} U(\mathbf{R}, t) = \mathbf{F}(\mathbf{R}, t)$$

for an external force  $\mathbf{F}$ , and hence we may express (2.8) as

$$\mathbf{M}\ddot{\mathbf{R}} + \nabla\phi = \mathbf{F}, \tag{2.10}$$

which describes the dynamic equilibrium of the continuum.

For a single linear (two-noded) bar element of initial length  $l_0$  and cross-sectional area  $A_0$  (assumed constant along the length of the bar), we have the position vector  $\mathbf{R} = \langle \mathbf{r}^1 \ \mathbf{r}^2 \rangle \in \mathbb{R}^6$  and shape function matrix

$$\mathbf{N}(s) = [N_1(s)\mathbf{I}_3 \quad N_2(s)\mathbf{I}_3] \in \mathbb{R}^{3 \times 6}; \quad N_1(s) = \frac{l_0 - s}{l_0}, \quad N_2(s) = \frac{s}{l_0} \tag{2.11}$$

for  $s \in [0, l_0]$ , with  $\mathbf{I}_3$  the three-dimensional identity matrix. We also have the mass and the stiffness of the bar as

$$m := A_0\rho_0 \int_0^{l_0} ds \equiv A_0\rho_0 l_0 \quad \text{and} \quad k := \frac{E\bar{A}}{\bar{l}}$$

for constant density  $\rho_0$ , where  $\bar{l}$  and  $\bar{A}$  represent the natural (undeformed) length and cross-sectional area, respectively, and  $E$  is Young’s modulus of the material. For an elastic, homogeneous bar with assumed linear interpolation given by (2.6) and (2.11), the strain energy function  $\phi(\cdot)$  is dependent only on the length  $l := \|\mathbf{r}^2 - \mathbf{r}^1\|$  of the bar, as it is in the equivalent central-force problem [32]. Thus we have

$$\nabla\phi(l) = \phi'(l)\nabla\left\{\sqrt{(\mathbf{r}^2 - \mathbf{r}^1) \cdot (\mathbf{r}^2 - \mathbf{r}^1)}\right\} = \frac{\phi'(l)}{l}\tilde{\mathbf{I}}\mathbf{R}, \quad \text{where } \tilde{\mathbf{I}} := \begin{pmatrix} \mathbf{I}_3 & -\mathbf{I}_3 \\ -\mathbf{I}_3 & \mathbf{I}_3 \end{pmatrix}.$$

Using the abbreviation

$$f(l) := \frac{\phi'(l)}{l}$$

the equation of motion becomes, for a single bar element,

$$\mathbf{M}\ddot{\mathbf{R}} + f(l)\tilde{\mathbf{I}}\mathbf{R} = \mathbf{F}. \tag{2.12}$$

In general, we consider systems comprising many bar elements. Applying Hamilton’s principle to the whole structure, the equilibrium equation (2.10) becomes an assembly of elemental force contributions in the form of (2.12) (e.g. [4], Section 7.2). Therefore we have

$$\mathbf{M}\ddot{\mathbf{R}} := \sum_{\substack{i=1, \\ j>i}}^N \mathbf{M}_{ij}\ddot{\mathbf{R}}_{ij}, \quad \nabla\phi := \sum_{\substack{i=1, \\ j>i}}^N \nabla\phi_{ij} = \sum_{\substack{i=1, \\ j>i}}^N f_{ij}(l_{ij})\tilde{\mathbf{I}}\mathbf{R}_{ij} \quad \text{and} \quad \mathbf{F} := \sum_{\substack{i=1, \\ j>i}}^N \mathbf{F}_{ij}$$

for  $\mathbf{M} \in \mathbb{R}^{3N \times 3N}$ ;  $\mathbf{R}, \ddot{\mathbf{R}}, \mathbf{F}, \nabla\phi \in \mathbb{R}^{3N}$ ;  $\mathbf{M}_{ij} \in \mathbb{R}^{6 \times 6}$ ; and  $\mathbf{R}_{ij}, \ddot{\mathbf{R}}_{ij}, \mathbf{F}_{ij}, \nabla\phi_{ij} \in \mathbb{R}^6$ , where  $(\cdot)_{ij}$  denotes a quantity pertaining to the element connecting nodes  $i$  and  $j$ , with

$$\mathbf{R}_{ij} := \langle \mathbf{r}^i \ \mathbf{r}^j \rangle. \tag{2.13}$$

(Note that the summation symbol  $\sum$  here conveys an assembly of elemental contributions.) Thus

$$m_{ij} := \begin{cases} A_{ij,0}\rho_{ij,0}l_{ij,0} & : \exists \text{ an element connecting nodes } i \text{ and } j, \text{ and} \\ 0 & : \text{ otherwise} \end{cases}$$

denotes an elemental mass; similarly,  $k_{ij}$  and  $\phi_{ij}$  denote the stiffness and strain energy of the element connecting nodes  $i$  and  $j$  if it exists, respectively, with  $k_{ij} := 0$  and  $\phi_{ij} := 0$  otherwise. We also have the total mass of the whole structure as

$$M = \sum_{\substack{i=1, \\ j>i}}^N m_{ij}. \tag{2.14}$$

In addition, we have the centre of mass for the whole system as

$$\mathbf{r}^c = \frac{1}{M} \sum_{\substack{i=1, \\ j>i}}^N \frac{m_{ij}}{2} (\mathbf{r}^i + \mathbf{r}^j) = \frac{1}{M} \mathcal{J}_N \mathbf{M} \mathbf{R}, \tag{2.15}$$

where the matrix  $\mathcal{J}_N \in \mathbb{R}^{3 \times 3N}$  is defined as

$$\mathcal{J}_N := \underbrace{[\mathbf{I}_3 \ \cdots \ \mathbf{I}_3]}_{N \text{ times}}. \tag{2.16}$$

We define the symmetric global matrix  $\mathcal{F}$  with submatrices  $\mathcal{F}^{ij} \in \mathbb{R}^{3 \times 3}$  such that

$$\mathcal{F}^{ij} := \begin{cases} \left( \sum_{k=1}^N f_{ik}(l_{ik}) \right) \mathbf{I}_3 & : \quad i = j, \\ -f_{ij}(l_{ij}) \mathbf{I}_3 & : \quad i \neq j, \end{cases} \tag{2.17}$$

where  $f_{ij} = \frac{\phi^{ij}(l_{ij})}{l_{ij}}$ . (Recall the notational difference between an *elemental matrix*  $\mathbf{M}_{ij} \in \mathbb{R}^{6 \times 6}$  and a *submatrix*  $\mathbf{M}^{ij} \in \mathbb{R}^{3 \times 3}$ .) Thus we have

$$\nabla \phi = \mathcal{F} \mathbf{R}, \tag{2.18}$$

and so we can now express (2.10) globally as

$$\mathbf{M} \ddot{\mathbf{R}} + \mathcal{F} \mathbf{R} = \mathbf{F}, \tag{2.19}$$

where the global mass matrix is symmetric and positive-definite. For linear systems,  $\mathcal{F}$  is constant, and we have  $\mathcal{F} \mathbf{R} = \mathbf{K} \mathbf{U}$ , where  $\mathbf{K}$  is the *stiffness matrix* and  $\mathbf{U}$  is a vector of nodal displacements. Here, we emphasise that  $\mathcal{F}$  is a function of  $\mathbf{R}$ . In view of our forthcoming time-integration schemes, we split (2.19) into a coupled first-order system with momenta

$$\mathbf{P} = \mathbf{M} \dot{\mathbf{R}} = \langle \mathbf{p}^1 \ \cdots \ \mathbf{p}^N \rangle \tag{2.20}$$

and positions  $\mathbf{R}$  as the primary variables. Thus we arrive at

$$\begin{aligned} \dot{\mathbf{P}} + \mathcal{F} \mathbf{R} &= \mathbf{F}, \\ \dot{\mathbf{R}} &= \mathbf{M}^{-1} \mathbf{P} \end{aligned} \tag{2.21}$$

which gives the equations of motion for an assembly of bar elements.

Turning our attention to the conservation properties of the semi-discrete system, for the linear momentum we have, from (2.2), (2.5) and (2.15),

$$\mathcal{L} = M \dot{\mathbf{r}}^c = \sum_{\substack{i=1, \\ j>i}}^N \frac{m_{ij}}{2} (\dot{\mathbf{r}}^i + \dot{\mathbf{r}}^j) = \sum_{i=1}^N \dot{\mathbf{p}}^i = \mathcal{J}_N \dot{\mathbf{P}} \tag{2.22}$$

and from (2.17) we see that  $\sum_{i=1}^N \mathcal{F}^{ij} = \mathbf{0}_3 \ \forall 1 \leq j \leq N$ , where  $\mathbf{0}_3$  is the three-dimensional zero matrix. Hence

$$\dot{\mathcal{L}} = \sum_{i=1}^N \mathbf{F}^i \tag{2.23}$$

and for systems where the external forces and reactions (due to supports) sum to zero, the total linear momentum is conserved. With regard to the angular momentum, from (2.3), (2.6) and (2.11), it can be shown that

$$\mathcal{J} = \sum_{i=1}^N \mathbf{r}^i \times \mathbf{p}^i \tag{2.24}$$

and hence

$$\dot{\mathcal{J}} = \sum_{i=1}^N (\dot{\mathbf{r}}^i \times \mathbf{p}^i + \mathbf{r}^i \times \dot{\mathbf{p}}^i) = \sum_{i=1}^N \mathbf{r}^i \times \mathbf{F}^i, \tag{2.25}$$

since both  $\mathbf{M}$  and  $\mathcal{F}$  are symmetric. Hence for systems where the moments of the external forces and reactions sum to zero, the total angular momentum is conserved.

The total energy  $H$  of the (undamped) system is the sum of the kinetic and potential energies, i.e.

$$H := T + \Phi. \tag{2.26}$$

We can write the kinetic energy  $T$  as

$$T = \frac{1}{2} \mathbf{P} \cdot \mathbf{M}^{-1} \mathbf{P} \quad (2.27)$$

and from (2.9) we have

$$\Phi(\mathbf{R}, t) = \phi(\mathbf{R}) - U(\mathbf{R}, t).$$

The potential function  $U$  can be split into conservative and non-conservative components, i.e.

$$U(\mathbf{R}, t) = U^C(\mathbf{R}) + U^{NC}(\mathbf{R}, t). \quad (2.28)$$

Differentiating  $H$  with respect to time then gives, from (2.8) and (2.9),

$$\dot{H} = -\frac{\partial U^{NC}}{\partial t}. \quad (2.29)$$

Thus for (monogenic) systems where  $U$  is not explicitly dependent on time, the total energy is conserved.

The conservation of momenta gives rise to a set of *relative equilibrium states* as solutions to system (2.21) when  $\mathbf{F} = \mathbf{0}$ . These are described in detail in [12,16–18] and references therein; here, we describe those characteristics of relative equilibria that are relevant to truss structures. We also derive the initial conditions under which these states exist. Corresponding results for 3D continua are given in Chapter 3 of [17].

For each type of momentum conserved by the system, there exists an associated *group motion*: thus from linear momentum conservation we have the group of translations, and from angular momentum conservation the group of rotations. Physically speaking, relative equilibria are solutions to system (2.21) when  $\mathbf{F} = \mathbf{0}$  that are *group orbits*; that is, solutions that differ only by a group motion. Therefore relative equilibrium motion can consist of either uniform translation, uniform rotation or, more generally, a combination of the two. This motion is a consequence of the invariance of the states of stresses and strains with respect to rigid body motion.

Mathematically speaking, relative equilibria are solutions to system (2.21) (with  $\mathbf{F} = \mathbf{0}$ ) that make the total energy  $H$  stationary for prescribed values of the linear momentum  $\mathcal{L}$  and angular momentum  $\mathcal{J}$  [17]. In other words, they are solutions  $\langle \tilde{\mathbf{P}}(t), \tilde{\mathbf{R}}(t) \rangle$  such that

$$\nabla H(\tilde{\mathbf{P}}, \tilde{\mathbf{R}}) = \mathbf{0} \quad \text{with constraints} \quad \mathcal{L}(\tilde{\mathbf{P}}) = \tilde{\mathcal{L}} \quad \text{and} \quad \mathcal{J}(\tilde{\mathbf{P}}, \tilde{\mathbf{R}}) = \tilde{\mathcal{J}}. \quad (2.30)$$

Accordingly, we need to solve

$$\begin{aligned} \nabla_{\mathbf{R}} H - (\nabla_{\mathbf{R}} \otimes \mathcal{J}) \boldsymbol{\omega} &= \mathbf{0}, \\ \nabla_{\mathbf{P}} H - (\nabla_{\mathbf{P}} \otimes \mathcal{L}) \boldsymbol{\eta} - (\nabla_{\mathbf{P}} \otimes \mathcal{J}) \boldsymbol{\omega} &= \mathbf{0}, \end{aligned} \quad (2.31)$$

where  $\boldsymbol{\eta}, \boldsymbol{\omega} \in \mathbb{R}^3$  are the Lagrange multipliers (e.g. [34, Section 4.9]), as shown in [17]. We define the skew-symmetric operator  $\hat{(\cdot)} : \mathbb{R}^3 \rightarrow \mathbb{R}^{3 \times 3}$  such that

$$\hat{\mathbf{u}} := \begin{pmatrix} 0 & -u_3 & u_2 \\ u_3 & 0 & -u_1 \\ -u_2 & u_1 & 0 \end{pmatrix} \iff \hat{\mathbf{u}} \mathbf{v} \equiv \mathbf{u} \times \mathbf{v} \quad \forall \mathbf{u}, \mathbf{v} \in \mathbb{R}^3, \quad (2.32)$$

and also the global matrix

$$\hat{\mathbf{\Omega}} := \begin{pmatrix} \hat{\boldsymbol{\omega}} & \boldsymbol{\theta}_3 & \cdots & \boldsymbol{\theta}_3 \\ \boldsymbol{\theta}_3 & \hat{\boldsymbol{\omega}} & \ddots & \vdots \\ \vdots & \ddots & \ddots & \boldsymbol{\theta}_3 \\ \boldsymbol{\theta}_3 & \cdots & \boldsymbol{\theta}_3 & \hat{\boldsymbol{\omega}} \end{pmatrix} \in \mathbb{R}^{3N \times 3N}. \quad (2.33)$$

From (2.22) and (2.24), (2.31) becomes, after some algebra,

$$\begin{aligned} \mathcal{F} \mathbf{R} + \hat{\mathbf{\Omega}} \mathbf{P} &= \mathbf{0}, \\ \mathbf{M}^{-1} \mathbf{P} - \mathcal{J}_N^T \boldsymbol{\eta} - \hat{\mathbf{\Omega}} \mathbf{R} &= \mathbf{0}. \end{aligned} \quad (2.34)$$

System (2.34), along with the constraints  $\mathcal{L}(\mathbf{P}) = \tilde{\mathcal{L}}$  and  $\mathcal{J}(\mathbf{P}, \mathbf{R}) = \tilde{\mathcal{J}}$ , can now be solved for  $\mathbf{R}(t)$ ,  $\mathbf{P}(t)$ ,  $\boldsymbol{\eta}(t)$  and  $\boldsymbol{\omega}(t)$ , thus furnishing the relative equilibrium solution  $\langle \tilde{\mathbf{P}}(t), \tilde{\mathbf{R}}(t) \rangle$ .

More information concerning the physical nature of relative equilibrium solutions is yet available, however, when we consider that  $\mathbf{R}(t)$  and  $\mathbf{P}(t)$  also satisfy (2.21) for  $\mathbf{F} = \mathbf{0}$ . Summarising the relevant findings of [17], we have the following important facts regarding the Lagrange multipliers  $\boldsymbol{\eta}(t)$  and  $\boldsymbol{\omega}(t)$ :

$$\begin{aligned} \text{(i)} \quad & \boldsymbol{\omega} \times \dot{\mathbf{r}}^c = \mathbf{0}, \\ \text{(ii)} \quad & \boldsymbol{\eta} = \dot{\mathbf{r}}^c - \boldsymbol{\omega} \times \mathbf{r}^c, \text{ and} \\ \text{(iii)} \quad & \boldsymbol{\omega} = \boldsymbol{\omega}_0 \quad \forall t \text{ i.e. } \boldsymbol{\omega} \text{ is constant.} \end{aligned} \tag{2.35}$$

From (2.2), (2.5), (2.23) and (2.35) we see that  $\boldsymbol{\eta}$  is also constant. We now define the *relative position vectors*

$$\bar{\mathbf{r}}^i := \mathbf{r}^i - \mathbf{r}^c; \quad 1 \leq i \leq N \quad \text{and} \quad \bar{\mathbf{R}} := \langle \bar{\mathbf{r}}^1 \quad \dots \quad \bar{\mathbf{r}}^N \rangle.$$

By substituting (2.35)<sub>2,3</sub> into (2.34)<sub>2</sub>, we get

$$\dot{\bar{\mathbf{R}}} - \widehat{\boldsymbol{\Omega}}_0 \bar{\mathbf{R}} = \mathbf{0}.$$

Since  $\sum_{j=1}^N \mathcal{F}^{ij} = \mathbf{0}_3 \quad \forall 1 \leq i \leq N$ , we have  $\mathcal{F} \bar{\mathbf{R}} = \mathcal{F} \bar{\mathbf{R}}$ ; furthermore, it can be shown that  $\widehat{\boldsymbol{\Omega}}_0$  commutes with  $\mathbf{M}$ , and so by virtue of (2.35)<sub>1,3</sub> we can write  $\widehat{\boldsymbol{\Omega}} \mathbf{P} = \widehat{\boldsymbol{\Omega}}_0 \mathbf{M} \bar{\mathbf{R}}$ . Therefore (2.34) becomes

$$\begin{aligned} \mathcal{F} \bar{\mathbf{R}} + \widehat{\boldsymbol{\Omega}}_0 \mathbf{M} \bar{\mathbf{R}} &= \mathbf{0}, \\ \dot{\bar{\mathbf{R}}} - \widehat{\boldsymbol{\Omega}}_0 \bar{\mathbf{R}} &= \mathbf{0}. \end{aligned} \tag{2.36}$$

System (2.36) thus gives necessary and sufficient conditions for the relative equilibrium solution  $\mathbf{P}(t) = \tilde{\mathbf{P}}(t)$  and  $\mathbf{R}(t) = \tilde{\mathbf{R}}(t)$  of (2.21) to exist, in terms of the constant vector  $\boldsymbol{\omega}$  and the position and velocity of the centre of mass of the structure (which are related). Thus the initial conditions for relative equilibrium states are given by

$$\begin{aligned} \mathcal{F}_0 \bar{\mathbf{R}}_0 + \widehat{\boldsymbol{\Omega}}_0 \mathbf{M} \bar{\mathbf{V}}_0 &= \mathbf{0}, \\ \bar{\mathbf{V}}_0 &= \widehat{\boldsymbol{\Omega}}_0 \bar{\mathbf{R}}_0, \end{aligned} \tag{2.37}$$

where  $\dot{\mathbf{R}}(0) = \mathbf{V}_0$ ,  $\dot{\mathbf{r}}^c(0) = \mathbf{v}_0^c$ ,  $\bar{\mathbf{V}}_0 := \langle \mathbf{v}_0^1 - \mathbf{v}_0^c \quad \dots \quad \mathbf{v}_0^N - \mathbf{v}_0^c \rangle$  and  $\mathcal{F}_0 \equiv \mathcal{F}(\mathbf{R}_0)$ . Finally, we can obtain the explicit form of  $\tilde{\mathbf{P}}(t)$  and  $\tilde{\mathbf{R}}(t)$  from (2.36)<sub>2</sub>. Introducing the matrix exponential  $\exp(t\mathbf{A})$  defined by

$$\exp(t\mathbf{A}) := \sum_{s=0}^{\infty} \frac{t^s}{s!} \mathbf{A}^s, \tag{2.38}$$

which is convergent for all  $\mathbf{A} \in \mathbb{R}^{3N \times 3N}$  (e.g. [35, Chapter 5]), we have

$$\begin{aligned} \dot{\bar{\mathbf{R}}} = \widehat{\boldsymbol{\Omega}}_0 \bar{\mathbf{R}} &\iff \bar{\mathbf{R}}(t) = \exp(t\widehat{\boldsymbol{\Omega}}_0) \bar{\mathbf{R}}_0 \iff \bar{\mathbf{r}}^i = \exp(t\widehat{\boldsymbol{\omega}}_0) \bar{\mathbf{r}}_0^i, \\ \text{and } \dot{\bar{\mathbf{R}}}(t) = \exp(t\widehat{\boldsymbol{\Omega}}_0) \bar{\mathbf{V}}_0 &\iff \dot{\bar{\mathbf{r}}}^i = \exp(t\widehat{\boldsymbol{\omega}}_0) \dot{\bar{\mathbf{r}}}_0^i; \quad 1 \leq i \leq N \\ \text{with } \dot{\mathbf{r}}^c(t) &= \mathbf{v}_0^c + t\mathbf{v}_0^c. \end{aligned} \tag{2.39}$$

It can be shown (e.g. [36]) that the transformation matrix  $\exp(t\widehat{\boldsymbol{\omega}}_0)$  represents a rotation of angle  $\|\boldsymbol{\omega}_0\|t$  about an axis parallel to  $\boldsymbol{\omega}_0$ . Thus all of the nodes remain fixed in relation to one another, so that *the structure rotates as a rigid body, with (constant) angular velocity  $\boldsymbol{\omega}_0$* . Recalling (2.35), we also note that *the velocity of the centre of mass is aligned with the axis of rotation*. Therefore the overall motion is *a superposition of uniform rotation and uniform translation along the axis of rotation*. The two special cases are then pure rotation (when  $\mathbf{v}_0^c = \mathbf{0}$ ) and pure translation (when  $\boldsymbol{\omega}_0 = \mathbf{0}$ ).

### 3. Algorithm derivation

Following the approach given in [32], we now describe a family of single-step time-integration schemes to solve system (2.21) approximately, that can be specialised to conserve various constants of motion. These algorithms are given in terms of global position and momentum vectors, which implies that these quantities are defined to be continuous across element boundaries. The general form for such a family is

$$\begin{aligned} \mathbf{R}_{n+1} &= \mathbf{A}\mathbf{R}_n + \mathbf{B}\mathbf{P}_n + \mathbf{R}_F, \\ \mathbf{P}_{n+1} &= \mathbf{C}\mathbf{R}_n + \mathbf{D}\mathbf{P}_n + \mathbf{P}_F, \end{aligned} \tag{3.1}$$

where  $\mathbf{R}_k, \mathbf{P}_k \in \mathbb{R}^{3N}$  are the discrete approximations to the positions  $\mathbf{R}(t_k)$  and momenta  $\mathbf{P}(t_k)$  at time  $t_k \geq 0$ ;  $\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D} \in \mathbb{R}^{3N \times 3N}$  are matrices of parameters that depend upon the configuration and the time-step  $\Delta t$ ; and  $\mathbf{R}_F, \mathbf{P}_F \in \mathbb{R}^{3N}$  are vectors that depend on the external force  $\mathbf{F}$ . Given the form of the matrices  $\mathcal{F}$  and  $\mathbf{M}$  in (2.21), we now restrict matrices  $\mathbf{A}, \mathbf{B}, \mathbf{C}$  and  $\mathbf{D}$  to have scaled unit submatrices also; therefore

$$\mathbf{A}^{ij} = a^{ij} \mathbf{I}_3, \quad \mathbf{B}^{ij} = b^{ij} \mathbf{I}_3, \quad \mathbf{C}^{ij} = c^{ij} \mathbf{I}_3, \quad \text{and} \quad \mathbf{D}^{ij} = d^{ij} \mathbf{I}_3; \quad 1 \leq i, j \leq N.$$

Thus in general there are  $4N^2$  unspecified parameters in (3.1), excluding those related to  $\mathbf{R}_F$  and  $\mathbf{P}_F$ .

We now define

$$\mathbf{Z}_n := \begin{Bmatrix} \mathbf{R}_n \\ \mathbf{P}_n \end{Bmatrix} \quad \text{and} \quad \mathbf{Z}_F(\mathbf{Z}_{n+1}, \mathbf{Z}_n, \mathbf{F}, \Delta t) := \begin{Bmatrix} \mathbf{R}_F \\ \mathbf{P}_F \end{Bmatrix}$$

and can thus express algorithm (3.1) in matrix form as

$$\mathbf{Z}_{n+1} = \mathcal{B}_{n+1} \mathbf{Z}_n + \mathbf{Z}_F, \quad \text{where} \quad \mathcal{B}_{n+1} \equiv \mathcal{B}(\mathbf{Z}_{n+1}, \mathbf{Z}_n, \Delta t) := \begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix} \tag{3.2}$$

(with  $\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}, \mathbf{P}_F$  and  $\mathbf{R}_F$  dependent on  $\mathbf{Z}_{n+1}, \mathbf{Z}_n$  and  $\Delta t$ , and  $\mathbf{P}_F$  and  $\mathbf{R}_F$  also dependent on  $\mathbf{F}$ ). We require that  $\mathcal{B}_{n+1}$  be non-singular, to prevent the possible occurrence of the solution  $\mathbf{Z}_{n+1} = \mathbf{0}$  when  $\mathbf{Z}_F = \mathbf{0}$  and  $\mathbf{Z}_n \neq \mathbf{0}$ . Thus  $\mathbf{A}, \mathbf{B}, \mathbf{C}$  and  $\mathbf{D}$  are constrained by the condition  $\det(\mathcal{B}_{n+1}) \neq 0$ . An expression for the determinant of a  $2 \times 2$  block matrix in the form of (3.2)<sub>2</sub> can be given in terms of its component matrices as

$$\det(\mathcal{B}_{n+1}) = \det(\mathbf{A})[\det(\mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B})] \tag{3.3}$$

provided that the upper left matrix  $\mathbf{A}$  is non-singular (e.g. Chapter 0 of [35]).

### 3.1. Inherent momentum conservation

We now look at the conditions under which algorithm (3.1) conserves both linear and angular momenta in the absence of external forces. From (2.22), the total discrete linear momentum  $\mathcal{L}_n$  at time-step  $n$  is computed as

$$\mathcal{L}_n = \sum_{i=1}^N \mathbf{p}_n^i = \mathcal{J}_N \mathbf{P}_n.$$

From (3.1), we have

$$\mathcal{L}_\Delta = \mathcal{J}_N \mathbf{P}_\Delta = \mathcal{J}_N [\mathbf{C}\mathbf{R}_n + (\mathbf{D} - \mathbf{I}_{3N})\mathbf{P}_n + \mathbf{P}_F],$$

where here and throughout the paper we define  $(\cdot)_\Delta := (\cdot)_{n+1} - (\cdot)_n$  for any given physical quantity  $(\cdot)$ , and we note that  $\mathbf{R}_n$  and  $\mathbf{P}_n$  are arbitrary. Therefore algorithm (3.1) will conserve linear momentum in general if and only if

$$\mathcal{J}_N \mathbf{C} = \mathcal{O}_N, \quad \mathcal{J}_N \mathbf{D} = \mathcal{J}_N \quad \text{and} \quad \mathcal{J}_N \mathbf{P}_F = \mathbf{0}, \tag{3.4}$$

where the matrix  $\mathcal{O}_N \in \mathbb{R}^{3 \times 3N}$  is defined as

$$\mathcal{O}_N := \underbrace{[\mathbf{0}_3 \quad \cdots \quad \mathbf{0}_3]}_{N \text{ times}}. \tag{3.5}$$

Thus the columns of  $\mathbf{C}$  must sum to zero (meaning that  $\mathbf{C}$  is singular), and those of  $\mathbf{D}$  to one, whilst the force-related momentum components  $\mathbf{p}_F^i$  must also sum to zero. This amounts to  $2N$  conditions involving  $\mathbf{C}$  and  $\mathbf{D}$  that algorithm (3.1) must satisfy.

From (2.24), the total discrete angular momentum  $\mathcal{J}_n$  at time-step  $n$  is given as

$$\mathcal{J}_n = \sum_{i=1}^N \mathbf{r}_n^i \times \mathbf{p}_n^i.$$

From (3.1), we have

$$\begin{aligned} \mathcal{J}_\Delta &= \sum_{i=1}^N (\mathbf{r}_{n+1}^i \times \mathbf{p}_{n+1}^i - \mathbf{r}_n^i \times \mathbf{p}_n^i) = \sum_{i=1}^N \left( \left[ \sum_{j=1}^N (a^{ij} \mathbf{r}_n^j + b^{ij} \mathbf{p}_n^j) + \mathbf{r}_F^i \right] \times \left[ \sum_{k=1}^N (c^{ik} \mathbf{r}_n^k + d^{ik} \mathbf{p}_n^k) + \mathbf{p}_F^i \right] - \mathbf{r}_n^i \times \mathbf{p}_n^i \right) \\ &= \sum_{i=1}^N \left[ \sum_{j,k=1}^N (a^{ij} c^{ik} \mathbf{r}_n^j \times \mathbf{r}_n^k + b^{ij} d^{ik} \mathbf{p}_n^j \times \mathbf{p}_n^k + a^{ij} d^{ik} \mathbf{r}_n^j \times \mathbf{p}_n^k + b^{ij} c^{ik} \mathbf{p}_n^j \times \mathbf{r}_n^k) + \mathbf{r}_{n+1}^i \times \mathbf{p}_F^i + \mathbf{r}_F^i \times \mathbf{p}_{n+1}^i - \mathbf{r}_n^i \times \mathbf{p}_F^i - \mathbf{r}_n^i \times \mathbf{p}_n^i \right] \\ &= \sum_{i=1}^N \left( \sum_{\substack{j=1, \\ k>j}}^N [(a^{ij} c^{ik} - a^{ik} c^{ij}) \mathbf{r}_n^j \times \mathbf{r}_n^k + (b^{ij} d^{ik} - b^{ik} d^{ij}) \mathbf{p}_n^j \times \mathbf{p}_n^k] \right) \\ &\quad + \sum_{i=1}^N \left( \sum_{j,k=1}^N (a^{ij} d^{ik} - b^{ik} c^{ij}) \mathbf{r}_n^j \times \mathbf{p}_n^k - \mathbf{r}_n^i \times \mathbf{p}_n^i + \mathbf{r}_{n+1}^i \times \mathbf{p}_F^i + \mathbf{r}_F^i \times \mathbf{p}_{n+1}^i - \mathbf{r}_n^i \times \mathbf{p}_F^i \right), \end{aligned}$$

where the vectors  $\mathbf{r}_n^j \times \mathbf{r}_n^k$ ,  $\mathbf{p}_n^j \times \mathbf{p}_n^k$  and  $\mathbf{r}_n^j \times \mathbf{p}_n^k$  are unrelated for all  $j, k > j; 1 \leq j, k \leq N$ . Therefore algorithm (3.1) will conserve angular momentum in general if and only if

$$\left. \begin{aligned} \sum_{i=1}^N (a^{ij} c^{ik} - a^{ik} c^{ij}) &= 0, \\ \sum_{i=1}^N (b^{ij} d^{ik} - b^{ik} d^{ij}) &= 0 \end{aligned} \right\} \forall j, k > j; \quad \sum_{i=1}^N (a^{ij} d^{ik} - b^{ik} c^{ij}) = \delta_{ij} \quad \forall j, k$$

and  $\mathbf{r}_F^i = \mathbf{p}_F^i = \mathbf{0} \quad \forall i; 1 \leq i, j, k \leq N$ ,

where  $\delta_{ij}$  is the Kronecker delta such that  $\delta_{ij} = 1$  if  $i = j$  and 0 otherwise. This amounts to a further  $2N^2 - N$  conditions on  $\mathbf{A}, \mathbf{B}, \mathbf{C}$  and  $\mathbf{D}$  which, when expressed in matrix form, become

$$\mathbf{A}^T \mathbf{C} = \mathbf{C}^T \mathbf{A}, \quad \mathbf{B}^T \mathbf{D} = \mathbf{D}^T \mathbf{B} \quad \text{and} \quad \mathbf{A}^T \mathbf{D} - \mathbf{C}^T \mathbf{B} = \mathbf{I}_{3N}, \tag{3.6}$$

along with  $\mathbf{R}_F = \mathbf{P}_F = \mathbf{0}$ . Now, from (3.3) we see that

$$\det(\mathcal{B}_{n+1}) = \det(\mathbf{A}^T) [\det(\mathbf{D} - \mathbf{C} \mathbf{A}^{-1} \mathbf{B})] = \det(\mathbf{A}^T \mathbf{D} - \mathbf{A}^T \mathbf{C} \mathbf{A}^{-1} \mathbf{B})$$

and when incorporating conditions (3.6), this becomes

$$\det(\mathcal{B}_{n+1}) = \det(\mathbf{I}_{3N} + \mathbf{C}^T \mathbf{B} - \mathbf{C}^T \mathbf{A} \mathbf{A}^{-1} \mathbf{B}) = \det(\mathbf{I}_{3N}) = 1. \tag{3.7}$$

Hence angular momentum conservation automatically assures that  $\mathcal{B}_{n+1}$  is non-singular, as required, with unit determinant. Thus any family of single-step algorithms that conserves both linear and angular momenta can have at most  $2N^2 - N$  free parameters, excluding those related to the external force. Note that in linear dynamics,  $\mathcal{B}_{n+1}$  is a constant amplification matrix; thus (3.7) implies preservation of the phase space (Liouville’s theorem) for any angular momentum-conserving algorithm in linear dynamics.

### 3.2. Choice of parameters

We wish to express our algorithms in a form that will relate to our previous work on the central-force problem [32]; thus we choose

$$\begin{aligned} \frac{1}{\Delta t} (\mathbf{P}_\Delta + \mathcal{G} \mathbf{P}_{1/2}) &= -\mathcal{X} \mathbf{R}_{1/2} + \mathbf{F}_a, \\ \frac{1}{\Delta t} (\mathbf{R}_\Delta - \mathcal{G}^T \mathbf{R}_{1/2}) &= \mathcal{M}^{-1} \mathbf{P}_{1/2} - \mathbf{V}_a, \end{aligned} \tag{3.8}$$

where  $\mathcal{G}, \mathcal{X}$  and  $\mathcal{M}$  are the parameter matrices (with  $\mathcal{M}$  necessarily non-singular),  $\mathbf{F}_a$  and  $\mathbf{V}_a$  are algorithmic force and velocity vectors that pertain to the external force; here and throughout the paper we define  $(\cdot)_{1/2} := \frac{1}{2}[(\cdot)_n + (\cdot)_{n+1}]$ . As before, each matrix will consist of unit submatrices:

$$\mathcal{G}^{ij} = g^{ij} \mathbf{I}_3, \quad \mathcal{X}^{ij} = x^{ij} \mathbf{I}_3 \quad \text{and} \quad \mathcal{M}^{ij} = \mu^{ij} \mathbf{I}_3. \tag{3.9}$$

We will refer to this family of momentum-conserving algorithms collectively as Algorithm MC. Note that for the specific case where  $\mathcal{G} := \mathbf{0}_{3N}$ ,  $\mathcal{M} := \mathbf{M}$ ,  $\mathbf{F}_a := \mathbf{F}$  and  $\mathbf{V}_a := \mathbf{0}$ , we have

$$\begin{aligned} \frac{1}{\Delta t} \mathbf{P}_\Delta &= -\mathcal{X} \mathbf{R}_{1/2} + \mathbf{F}, \\ \frac{1}{\Delta t} \mathbf{R}_\Delta &= \mathbf{M}^{-1} \mathbf{P}_{1/2}, \end{aligned} \quad (3.10)$$

which is the form of several familiar time-integration schemes (e.g. [8,15,37]), each distinguished by its definition of  $\mathcal{X}$ .

We can relate this set of parameters to the canonical set  $\{\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}, \mathbf{P}_F, \mathbf{R}_F\}$  by forming the equation

$$\mathbf{Z}_\Delta = \mathcal{C}_{n+1} \mathbf{Z}_{1/2} + \mathbf{Z}_a$$

from (3.8), where

$$\mathcal{C}_{n+1} := \begin{pmatrix} \mathcal{G}^\top & \Delta t \mathcal{M}^{-1} \\ -\Delta t \mathcal{X} & -\mathcal{G} \end{pmatrix} \quad \text{and} \quad \mathbf{Z}_a := \Delta t \begin{Bmatrix} -\mathbf{V}_a \\ \mathbf{F}_a \end{Bmatrix}. \quad (3.11)$$

This in turn gives us

$$\mathbf{Z}_{n+1} = \left( \mathbf{I}_{6N} - \frac{1}{2} \mathcal{C}_{n+1} \right)^{-1} \left[ \left( \mathbf{I}_{6N} + \frac{1}{2} \mathcal{C}_{n+1} \right) \mathbf{Z}_n + \mathbf{Z}_a \right],$$

and thus from (3.2)<sub>1</sub> we have

$$\mathcal{B}_{n+1} = \left[ \mathbf{I}_{6N} - \frac{1}{2} \mathcal{C}_{n+1} \right]^{-1} \left[ \mathbf{I}_{6N} + \frac{1}{2} \mathcal{C}_{n+1} \right] \quad \text{and} \quad \mathbf{Z}_F = \left[ \mathbf{I}_{6N} - \frac{1}{2} \mathcal{C}_{n+1} \right]^{-1} \mathbf{Z}_a. \quad (3.12)$$

Combining (3.11) and (3.12) leads to  $\{\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}, \mathbf{P}_F, \mathbf{R}_F\}$  in terms of  $\{\mathcal{G}, \mathcal{X}, \mathcal{M}, \mathbf{F}_a, \mathbf{V}_a\}$ , provided all the relevant inverses exist.

#### 4. Properties of the algorithm

We will now derive conditions for parameters  $\mathcal{G}$ ,  $\mathcal{X}$ ,  $\mathcal{M}$ ,  $\mathbf{F}_a$  and  $\mathbf{V}_a$  under which the algorithm will conserve linear and angular momenta, conserve the total energy, preserve relative equilibria and be time-reversible.

##### 4.1. Conservation of linear momentum

**Proposition 1.** *Algorithm MC gives the discrete linear momentum derivative as*

$$\frac{1}{\Delta t} \mathcal{L}_\Delta = \sum_{i=1}^N \mathbf{F}_a^i \quad (4.1)$$

provided that

$$\sum_{i=1}^N \mathcal{G}^{ij} = \sum_{i=1}^N \mathcal{X}^{ij} = \mathbf{0}_3 \quad \forall 1 \leq j \leq N. \quad (4.2)$$

Thus linear momentum is conserved whenever  $\sum_{i=1}^N \mathbf{F}_a^i = \mathbf{0}$ .

The proof is given in Appendix A.1; note the similarity between (4.1) and (2.23). Eq. (4.2) therefore imposes  $2N$  conditions on matrices  $\mathcal{G}$  and  $\mathcal{X}$ .

##### 4.2. Conservation of angular momentum

**Proposition 2.** *Algorithm MC gives the discrete angular momentum derivative as*

$$\frac{1}{\Delta t} \mathcal{J}_\Delta = \sum_{i=1}^N \left( \mathbf{r}_{1/2}^i \times \mathbf{F}_a^i + \mathbf{p}_{1/2}^i \times \mathbf{v}_a^i \right) \quad (4.3)$$

provided that

$$\mathcal{X} = \mathcal{X}^\top \quad \text{and} \quad \mathcal{M} = \mathcal{M}^\top. \quad (4.4)$$

Thus angular momentum is conserved whenever  $\mathbf{F}_a = \mathbf{V}_a = \mathbf{0}$ .

The proof is given in [Appendix A.2](#); note the similarity between (4.3) and (2.25). Eq. (4.4) imposes a further  $N^2 - N$  conditions on matrices  $\mathcal{X}$  and  $\mathcal{M}$ . The total number of free parameters remaining in  $\mathcal{G}$ ,  $\mathcal{X}$  and  $\mathcal{M}$  after conservation of momenta has been secured is  $2N^2 - N$ . This tallies with the amount given in Section 3 for  $\mathbf{A}$ ,  $\mathbf{B}$ ,  $\mathbf{C}$  and  $\mathbf{D}$ ; hence *Algorithm MC is a fully general form for single-step momentum-conserving algorithms*, provided that (4.2) and (4.4) are satisfied. Note that this means  $\mathcal{G}$  and  $\mathcal{X}$  are singular, and that  $\mathcal{X}$  and  $\mathcal{M}$  are symmetric.

#### 4.3. Conservation of energy

From (2.9), (2.26) and (2.27), the discrete total energy at time-step  $n$  is given as

$$H_n = \frac{1}{2} \mathbf{P}_n \cdot \mathbf{M}^{-1} \mathbf{P}_n + \phi_n - U_n. \quad (4.5)$$

For Algorithm MC to be energy-conserving, we therefore require

$$H_{n+1} = H_n \quad (4.6)$$

to hold for all  $n$  under the appropriate conditions.

##### 4.3.1. A single, global condition for energy conservation

**Proposition 3.** *Algorithm MC gives the discrete energy derivative as*

$$\frac{1}{\Delta t} H_\Delta = -\frac{U_\Delta^{\text{NC}}}{\Delta t} \quad (4.7)$$

if  $\mathcal{X} := \kappa \overline{\mathcal{X}}$  for an arbitrary matrix  $\overline{\mathcal{X}}$ , where  $\kappa$  is defined by

$$\begin{aligned} A\kappa + B &= 0 \quad \text{for} \\ A &= -\Delta t \overline{\mathcal{X}} \mathbf{R}_{1/2} \cdot \mathbf{M}^{-1} \mathbf{P}_{1/2} \quad \text{and} \quad B = \phi_\Delta - U_\Delta^{\text{C}} + (\Delta t \mathbf{F}_a - \mathcal{G} \mathbf{P}_{1/2}) \cdot \mathbf{M}^{-1} \mathbf{P}_{1/2}, \end{aligned} \quad (4.8)$$

with  $U^{\text{C}}$  and  $U^{\text{NC}}$  denoting the conservative and non-conservative parts of the potential function  $U$  as defined in (2.28). Thus energy is conserved whenever  $U^{\text{NC}} = 0$ .

The proof is given in [Appendix A.3](#), and we see the similarity between (4.7) and (2.29). Eq. (4.8) therefore fixes one more of the free parameters; thus *families of single-step energy-momentum algorithms can have up to  $2N^2 - N - 1$  free parameters, excluding those related to the external force*.

Algorithms that do not conserve energy can be converted to energy-conserving algorithms by transforming  $\mathcal{X} \mapsto \kappa \overline{\mathcal{X}}$ , with  $\kappa$  as defined in (4.8). This leaves any momentum-conserving attributes of the algorithm intact, since the structure of  $\overline{\mathcal{X}}$  is unaltered. For example, the *symplectic-momentum mid-point algorithm* (e.g. [15,23,24,38]), which is an instance of (3.10) with  $\overline{\mathcal{X}}$  defined such that  $\overline{\mathcal{X}} \mathbf{R}_{1/2} := \nabla_{\mathbf{R}_{1/2}} \phi(\mathbf{R}_{1/2})$ , can be made to conserve energy with  $\kappa$  defined by (4.8); this becomes

$$[\nabla_{\mathbf{R}_{1/2}} \phi(\mathbf{R}_{1/2}) \cdot \mathbf{R}_\Delta] \kappa = \phi_\Delta \quad (4.9)$$

for constant  $\mathbf{F}$ , which matches Eq. (2.27) of [15], where this algorithm was first proposed. The property of symplecticity is not present in the new algorithm, however. A similar conversion can be performed on another instance of (3.10) known as the *assumed distance method* [37,25], which preserves relative equilibria but does not conserve energy in general; defining  $\kappa$  according to (4.8) generates an energy-conserving algorithm, but the property of preserving relative equilibria is lost, as will be explained in Section 4.5.

There are serious flaws in using a single scalar condition to establish energy conservation, however. Firstly, (4.8) admits a solution for  $\kappa$  only when  $A \neq 0$  or  $\lim_{A \rightarrow 0} \{-B/A\}$  exists. Except for the particular case of a single bar element fixed at one end with no external force (thereby equivalent to a central-force problem), this cannot be guaranteed; hence a drawback with this method of energy conservation is the fact that, for typical  $\mathcal{G}$ ,  $\mathcal{M}$  and  $\overline{\mathcal{X}}$ , it does not guarantee an energy-conserving solution for all possible configurations and time-step sizes, since  $\kappa$  cannot always be defined by (4.8). (This issue was also raised in [15].) Another objection is that it introduces a global coupling between all of the degrees of freedom; this cannot be motivated physically, as mentioned in [22], and it destroys any sparsity in the tangent stiffness matrix used during the non-linear solution procedure. Furthermore, adding energy conservation to a non-conserving algorithm does not necessarily enhance its performance: results have shown that in cases where a non-conserving scheme fails due to a large increase in the total energy, its conserving counterpart invariably fails to provide a solution (i.e., the non-linear solution procedure does not converge) [39,40]. In the case of the assumed distance method, adding energy conservation in this way can actually degrade its performance [40].

In summary, it does not seem as though algorithms should be designed to conserve energy in this way, since there are many undesirable side-effects. An alternative approach to energy conservation [8,11,22,41] will now be presented.

#### 4.3.2. Elemental conditions for energy conservation

From Appendix A.3 we see that

$$H_{\Delta} = -U_{\Delta}^{\text{NC}} \iff (\Delta t \mathbf{F}_a - \mathcal{G} \mathbf{P}_{1/2}) \cdot \mathbf{M}^{-1} \mathbf{P}_{1/2} + \phi_{\Delta} - U_{\Delta}^{\text{C}} = \Delta t \mathcal{X} \mathbf{R}_{1/2} \cdot \mathbf{M}^{-1} \mathbf{P}_{1/2},$$

and thus using (3.8)<sub>2</sub> to substitute for  $\mathbf{P}_{1/2}$  we arrive at the condition

$$\begin{aligned} \mathcal{X} \mathbf{R}_{1/2} \cdot \mathbf{M}^{-1} \mathcal{M} (\mathbf{R}_{\Delta} - \mathcal{G}^{\text{T}} \mathbf{R}_{1/2}) &= \phi_{\Delta} - U_{\Delta}^{\text{C}} - \Delta t \mathcal{X} \mathbf{R}_{1/2} \cdot \mathbf{M}^{-1} \mathcal{M} \mathbf{V}_a + \left( \Delta t \mathbf{F}_a - \mathcal{G} \mathcal{M} \left[ \frac{1}{\Delta t} (\mathbf{R}_{\Delta} - \mathcal{G}^{\text{T}} \mathbf{R}_{1/2}) + \mathbf{V}_a \right] \right) \\ &\quad \cdot \mathbf{M}^{-1} \mathcal{M} \left[ \frac{1}{\Delta t} (\mathbf{R}_{\Delta} - \mathcal{G}^{\text{T}} \mathbf{R}_{1/2}) + \mathbf{V}_a \right] \end{aligned} \quad (4.10)$$

for energy conservation.

For algorithms of the form (3.10), (4.10) reduces to

$$\mathcal{X} \mathbf{R}_{1/2} \cdot \mathbf{R}_{\Delta} = \phi_{\Delta} - U_{\Delta}^{\text{C}} + \mathbf{F} \cdot \mathbf{R}_{\Delta} = \phi_{\Delta}, \quad (4.11)$$

since  $\mathbf{F} \cdot \mathbf{R}_{\Delta} = U_{\Delta}^{\text{C}}$ . Given that these algorithms conserve both momenta,  $\mathcal{X}$  must be symmetric with rows and columns that sum to zero. Here and throughout the paper we will use the shorthand notation  $\langle \mathbf{B}_{ij} \rangle$  to denote symmetric block matrices  $\mathbf{A} \in \mathbb{R}^{nm \times nm}$  of the form

$$\mathbf{A}^{ij} = \begin{cases} \sum_{k=1}^N \mathbf{B}_{ik} & : i = j, \\ -\mathbf{B}_{ij} & : i \neq j, \end{cases}$$

for some  $\mathbf{B}_{ij} \in \mathbb{R}^{m \times m}$ ; note that each row (and column) sums to zero. Hence we can write

$$\mathcal{X} \equiv \langle \xi_{ij} \mathbf{I}_3 \rangle \quad (4.12)$$

for some  $\xi_{ij}$  (to be defined). It can easily be shown that, as a consequence of the form of  $\mathcal{X}$ ,

$$\mathcal{X} \mathbf{R}_{1/2} \cdot \mathbf{R}_{\Delta} = \sum_{\substack{i=1, \\ j>i}}^N \xi_{ij} \mathbf{r}_{1/2}^{ij} \cdot \mathbf{r}_{\Delta}^{ij} \quad \forall \mathbf{R}_n, \mathbf{R}_{n+1} \text{ and } \xi_{ij},$$

where  $\mathbf{r}^{ij} := \mathbf{r}^j - \mathbf{r}^i$ . Recalling the definition of  $\phi_{\Delta}$ , (4.11) becomes

$$\sum_{\substack{i=1, \\ j>i}}^N \xi_{ij} \mathbf{r}_{1/2}^{ij} \cdot \mathbf{r}_{\Delta}^{ij} = \sum_{\substack{i=1, \\ j>i}}^N \phi_{ij\Delta}. \quad (4.13)$$

This equation can clearly be satisfied by defining

$$\xi_{ij} := \frac{\phi_{ij\Delta}}{\mathbf{r}_{1/2}^{ij} \cdot \mathbf{r}_{\Delta}^{ij}} \quad \forall i, j, \quad (4.14)$$

implying  $\xi_{ij} = \xi_{ji}$  and  $\xi_{ii} = 0$ . The resulting algorithm is known as the *energy–momentum mid-point algorithm* [8,11,22,41], and does not require a global energy condition of the form (4.8). (This equation is satisfied by  $\kappa := 1$ .) Since the  $\xi_{ij} \equiv \xi_{ij}(\mathbf{r}^{ij})$  depend only on the element connecting nodes  $i$  and  $j$  (if it exists), we call these *elemental conditions* for energy conservation. The energy–momentum algorithm suffers none of the drawbacks mentioned in conjunction with global energy conservation, being well-defined when  $\mathbf{r}_{1/2}^{ij} \cdot \mathbf{r}_{\Delta}^{ij} = 0$ , and has been shown to perform well for a wide range of problems [8,11,32,40].

We now look at the conditions for general  $\mathcal{G}$ ,  $\mathcal{X}$  and  $\mathcal{M}$  under which *elemental energy conservation*—that is, energy conservation through satisfaction of simultaneous elemental equations—can be achieved alongside conservation of momenta. Elemental energy conservation is a much more stringent requirement than global energy conservation, with a far greater number of conditions to be satisfied. These amount to balancing each individual  $\phi_{ij\Delta}$  term in (4.10) with an algorithmic quantity wholly defined by  $\mathbf{r}^{ij}$  (the relative position vector for the element joining nodes  $i$  and  $j$ ), i.e.

$$\phi_{ij\Delta} = \alpha_{ij}(\mathbf{r}^{ij}) \quad \forall i, j,$$

for some  $\alpha_{ij}$ . Note that this is tantamount to stipulating that (4.10) consist *entirely* of elemental terms (i.e. those that are functions of some  $\mathbf{r}^{ij}$  only), except for those pertaining to the external force: if other terms were to appear, a further (non-

elemental) condition would be required to ensure that they sum to zero. Given the form of  $\mathcal{X}$ , we have the left-hand side of (4.10) as

$$\mathcal{X} \mathbf{R}_{1/2} \cdot \mathbf{M}^{-1} \mathcal{M} (\mathbf{R}_\Delta - \mathcal{G}^T \mathbf{R}_{1/2}) = \sum_{\substack{i=1, \\ j>i}}^N \xi_{ij} \mathbf{r}_{1/2}^{ij} \cdot (\check{\mathbf{r}}_\Delta^{ij} - \check{\mathbf{r}}_{1/2}^{ij}),$$

where  $\check{\mathbf{R}} = \mathbf{M}^{-1} \mathcal{M} \mathbf{R}$  and  $\check{\mathbf{R}} = \mathbf{M}^{-1} \mathcal{M} \mathcal{G}^T \mathbf{R}$ . Elemental energy conservation then requires  $\check{\mathbf{r}}^{ij} = a_{ij} \mathbf{r}^{ij}$  and  $\check{\mathbf{r}}^{ij} = b_{ij} \mathbf{r}^{ij}$ , and hence we must have  $\mathcal{M} := a \mathbf{M}$  and  $\mathcal{G} := \frac{b}{a} \mathbf{I}_n$  for scalars  $a$  and  $b$  (not necessarily constant), i.e.  $a_{ij} = a$ ,  $b_{ij} = b \forall i, j$ . Given that  $\mathcal{G}$  is singular (from Proposition 1), however, we must have  $\mathcal{G} := \mathbf{0}_{3N}$ , meaning  $b = 0$ . So the whole of (4.10) now becomes

$$a \mathcal{X} \mathbf{R}_{1/2} \cdot \mathbf{R}_\Delta = \phi_\Delta - U_\Delta^C - a \Delta t \mathcal{X} \mathbf{R}_{1/2} \cdot \mathbf{V}_a + a \mathbf{F}_a \cdot \mathbf{R}_\Delta + a \Delta t \mathbf{F}_a \cdot \mathbf{V}_a. \tag{4.15}$$

For problems without external forces, we have  $\mathbf{F}_a = \mathbf{V}_a = \mathbf{0}$  from Proposition 2, which requires that  $\mathcal{X} := \frac{1}{a} \langle \xi_{ij} \mathbf{I}_3 \rangle$  for  $\xi_{ij}$  defined by (4.14). For the general conservative case, then, (4.15) becomes

$$-U_\Delta^C - \Delta t \langle \xi_{ij} \mathbf{I}_3 \rangle \mathbf{R}_{1/2} \cdot \mathbf{V}_a + a \mathbf{F}_a \cdot (\mathbf{R}_\Delta + a \mathbf{V}_a) = 0,$$

which is satisfied by  $\mathbf{F}_a := \frac{1}{a} \mathbf{F}$ ,  $\mathbf{V}_a := \mathbf{0}$ . Therefore the only instances of Algorithm MC that give rise to elemental energy conservation are scaled variations of the energy–momentum mid-point algorithm, i.e.

$$\mathcal{G} := \mathbf{0}_{3N}, \quad \mathcal{X} := \alpha \langle \xi_{ij} \mathbf{I}_3 \rangle, \quad \mathcal{M} := \frac{1}{\alpha} \mathbf{M}, \quad \mathbf{F}_a := \alpha \mathbf{F} \quad \text{and} \quad \mathbf{V}_a := \mathbf{0} \tag{4.16}$$

for some  $\alpha(\mathbf{Z}_{n+1}, \mathbf{Z}_n, \Delta t) \in \mathbb{R}$ , with  $\xi_{ij}$  defined by (4.14).

#### 4.4. Preservation of relative equilibria

From (2.39), we have the exact solution of a relative equilibrium problem as

$$\bar{\mathbf{R}}(t) = \exp(t \hat{\Omega}_0) \bar{\mathbf{R}}_0, \quad \bar{\mathbf{R}}(t) = \exp(t \hat{\Omega}_0) \bar{\mathbf{V}}_0 \quad \text{and} \quad \mathbf{r}^c(t) = \mathbf{r}_0^c + t \mathbf{v}_0^c$$

with  $\exp(t \hat{\Omega}_0) \equiv \text{diag}[\exp(t \hat{\omega}_0)]$ , where  $\bar{\mathbf{R}}_0$  and  $\bar{\mathbf{V}}_0$  are given by

$$\bar{\mathbf{R}}_0 = \mathbf{R}_0 - \mathbf{R}_0^c, \quad \mathbf{R}_0^c := \underbrace{\langle \mathbf{r}_0^c \ \cdots \ \mathbf{r}_0^c \rangle}_{N \text{ times}}, \quad \bar{\mathbf{V}}_0 = \mathbf{V}_0 - \mathbf{V}_0^c \quad \text{and} \quad \mathbf{V}_0^c := \underbrace{\langle \mathbf{v}_0^c \ \cdots \ \mathbf{v}_0^c \rangle}_{N \text{ times}}.$$

We note from Section 2 that  $\omega_0 \times \mathbf{v}_0^c = \mathbf{0}$ , and that  $\mathbf{R}_0$  and  $\mathbf{V}_0$  are defined by (2.37). We now define a *relative equilibrium path* as the discrete solution

$$\mathbf{R}_k = \mathbf{R}_k^c + \exp(\lambda_r k \Delta t \hat{\Omega}_0) \bar{\mathbf{R}}_0, \quad \mathbf{P}_k = \mathbf{M} \left( \mathbf{V}_0^c + \exp(\lambda_t k \Delta t \hat{\Omega}_0) \bar{\mathbf{V}}_0 \right) \quad \text{and} \quad \mathbf{R}_k^c = \mathbf{R}_0^c + \lambda_t k \Delta t \mathbf{V}_0^c \quad \forall k, \tag{4.17}$$

where  $\lambda_r$  and  $\lambda_t$  are constants. Solution (4.17) describes a combination of uniform rotation about an axis of inertia and uniform translation along it, as does the exact solution. However, the angular velocity of the discrete solution is  $\lambda_r \omega_0$ , and the translational velocity is  $\lambda_t \mathbf{v}_0^c$ . The exact solution is evidently captured if  $\lambda_r = \lambda_t = 1$ ; for  $\lambda_r = \lambda_t \neq 1$ , the trajectory of the exact solution is recovered (with different velocities). In general, the discrete solution need not necessarily produce points along the orbit of the exact solution, as was the case for the central-force problem in [32].

For Algorithm MC to produce paths of relative equilibria under initial conditions (2.37), we require that the solution given in (4.17) be inserted into the algorithm without conflict for  $k = n$  and  $k = n + 1$ , for certain values of  $\lambda_r$  and  $\lambda_t$ . We now introduce the notation

$$\begin{aligned} (\cdot)^{\text{RE}} &:= (\cdot) |_{\mathbf{R}_n, \mathbf{P}_n, \mathbf{R}_{n+1}, \mathbf{P}_{n+1}} \text{ defined by (4.17)} \\ &\mathbf{R}_0, \mathbf{V}_0 \text{ defined by (2.37) and } \mathbf{F} = \mathbf{0} \end{aligned}$$

for a given quantity  $(\cdot)$ , to denote the value taken when a relative equilibrium solution is in effect. Note that this is equivalent to the notation used in [32] in the case of a central-force problem.

**Proposition 4.** Under initial conditions (2.37), Algorithm MC produces paths of relative equilibria provided that

$$\begin{aligned} \left( \mathcal{X}^{\text{RE}} - \lambda_r \frac{\tan(\frac{1}{2} \theta^{\text{RE}})}{\frac{1}{2} \theta^{\text{RE}}} \mathcal{F}_0 \right) \bar{\mathbf{R}}_0 = \mathbf{0}, \quad \left( \lambda_r \frac{\tan(\frac{1}{2} \theta^{\text{RE}})}{\frac{1}{2} \theta^{\text{RE}}} \mathcal{M}^{\text{RE}} - \mathbf{M} \right) \bar{\mathbf{V}}_0 = \mathbf{0}, \\ \mathcal{G}^{\text{RE}} = \mathbf{0}_{3N}, \quad (\mathbf{M} - \lambda_t \mathcal{M}^{\text{RE}}) \mathbf{V}_0^c = \mathbf{0} \quad \text{and} \quad \mathbf{F}_a^{\text{RE}} = \mathbf{V}_a^{\text{RE}} = \mathbf{0} \end{aligned} \tag{4.18}$$

for some fixed  $\lambda_r, \lambda_t$ , where  $\theta^{\text{RE}} := \lambda_r \|\omega_0\| \Delta t$  and  $-\pi < \theta^{\text{RE}} < \pi$ , and that for all  $n$  it gives a unique solution for  $\mathbf{R}_{n+1}$  and  $\mathbf{P}_{n+1}$  given  $\mathbf{R}_n, \mathbf{P}_n$  and  $\Delta t$ .

The proof is given in Appendix A.4. Eq. (4.18) therefore gives criteria that  $\mathcal{G}, \mathcal{X}, \mathcal{M}, \mathbf{F}_a$  and  $\mathbf{V}_a$  must satisfy under initial conditions (2.37) for the algorithm to produce a relative equilibrium path. In fact, it is possible to design algorithms to capture the relative equilibrium solution *exactly*, with  $\lambda_r = \lambda_t = 1$ ; an example of this shall be given in Section 6.3.

#### 4.5. Conservation of energy and the preservation of relative equilibria

**Lemma 1.** *Any algorithm that produces paths of relative equilibria also conserves energy along those paths.*

The proof is given in Appendix A.5; thus no inherent conflict arises from having conservation of energy and preservation of relative equilibria within the same algorithm.

**Lemma 2.** *Any algorithm that conserves energy globally cannot produce paths of relative equilibria in general, in that an algorithm conserving energy via (4.8) cannot always be well-defined under relative equilibrium conditions.*

This can be argued by contradiction as follows: if an algorithm produces paths of relative equilibria, it satisfies (4.18). It can then be seen from (A.1) in Appendix A.3 that  $A$  and  $B$  from (4.8) are such that  $A^{\text{RE}} = B^{\text{RE}} = 0$ , with no obvious value for

$$\lim_{\substack{A \rightarrow A^{\text{RE}} \\ B \rightarrow B^{\text{RE}}}} \left\{ \frac{-B}{A} \right\}$$

that is uniquely defined regardless of how  $A \rightarrow A^{\text{RE}}$  and  $B \rightarrow B^{\text{RE}}$ . This is one of those instances where  $\kappa$  cannot be evaluated via (4.8), as mentioned in Section 4.3.1.

**Lemma 3.** *Any algorithm that conserves energy elementally produces paths of relative equilibria such that*

$$\lambda_r = \frac{\frac{1}{2} \alpha^{\text{RE}} \theta^{\text{RE}}}{\tan(\frac{1}{2} \theta^{\text{RE}})} \quad \text{and} \quad \lambda_t = \alpha^{\text{RE}}, \quad (4.19)$$

for  $\theta^{\text{RE}} := \lambda_r \|\omega_0\| \Delta t$  and  $\alpha$  given in (4.16).

It is straightforward to show that  $\mathcal{X}$  defined by (4.12) and (4.14) is such that  $\mathcal{X}^{\text{RE}} = \mathcal{F}_0$ ; the result can then be read directly from (4.16) and (4.18). A corollary of Lemma 3 is that *no elementally energy-conserving algorithm can capture the exact trajectory of a general relative equilibrium problem*, given that  $\tan^{-1} x \neq x$  for  $x \neq 0$ . (For the special cases of pure translation and pure rotation, the exact trajectories can be recovered.)

Taken together, these three lemmas imply that instances of Algorithm MC are unlikely to both conserve energy and produce paths of relative equilibria unless they are scaled variations of the energy–momentum mid-point algorithm, as given by (4.16).

#### 4.6. Time reversibility

An algorithm is described as time-reversible if, at any given configuration  $\mathbf{Z}_{n+1}$ , applying a negative time-step of  $-\Delta t$  recovers the previous configuration  $\mathbf{Z}_n$  [1,21]. From (3.2), an algorithm is thus time-reversible if

$$\mathbf{Z}_{n+1} = \mathcal{B}(\mathbf{Z}_{n+1}, \mathbf{Z}_n, \Delta t) \mathbf{Z}_n + \mathbf{Z}_F(\mathbf{Z}_{n+1}, \mathbf{Z}_n, \Delta t) \iff \mathbf{Z}_n = \mathcal{B}(\mathbf{Z}_n, \mathbf{Z}_{n+1}, -\Delta t) \mathbf{Z}_{n+1} + \mathbf{Z}_F(\mathbf{Z}_n, \mathbf{Z}_{n+1}, -\Delta t). \quad (4.20)$$

We now introduce for any quantity  $(\cdot)$  the notation

$$(\cdot)^{\text{TR}} := (\cdot)|_{\mathbf{Z}_{n+1} \leftrightarrow \mathbf{Z}_n, \Delta t \leftrightarrow -\Delta t}.$$

**Proposition 5.** *Algorithm MC is time-reversible if*

$$\mathcal{X}^{\text{TR}} = \mathcal{X}, \quad \mathcal{M}^{\text{TR}} = \mathcal{M}, \quad \mathcal{G}^{\text{TR}} = -\mathcal{G}, \quad \mathbf{F}_a^{\text{TR}} = \mathbf{F}_a \quad \text{and} \quad \mathbf{V}_a^{\text{TR}} = \mathbf{V}_a. \quad (4.21)$$

The proof is given in Appendix A.6.

### 5. Local accuracy analysis

We now analyse the local accuracy characteristics of Algorithm MC, and investigate its capacity for higher-order accuracy when applied to general non-linear problems. We will also derive the series form of the exact solution. This section takes a very similar approach to that of [32].

We define the *local error vector* as

$$\epsilon := \mathbf{Z}_{n+1} - \mathbf{Z}(t_{n+1}) \quad \text{when } \mathbf{Z}_n = \mathbf{Z}(t_n), \tag{5.1}$$

with  $\mathbf{Z} = \langle \mathbf{R} \ \mathbf{P} \rangle$  as before: throughout this section, we will assume the solution at time-step  $n$  to be exact, i.e.  $\mathbf{Z}_n = \mathbf{Z}(t_n)$ . We also define the *residual vector*

$$\mathbf{g}(\mathbf{X}) := \mathcal{B}(\mathbf{X}, \mathbf{Z}_n, \Delta t)\mathbf{Z}_n + \mathbf{Z}_F(\mathbf{X}, \mathbf{Z}_n, \mathbf{F}, \Delta t) - \mathbf{X}, \tag{5.2}$$

where  $\mathcal{B}$  and  $\mathbf{Z}_F$  were introduced in Section 3. Consequently we have

$$\mathbf{g}(\mathbf{Z}_{n+1}) = \mathbf{0} \quad \text{and hence } \mathbf{g}[\mathbf{Z}(t_{n+1})] = \mathbf{g}(\mathbf{Z}_{n+1} - \epsilon) = -\nabla \mathbf{g}(\mathbf{Z}_{n+1})\epsilon + \mathcal{O}(\|\epsilon\|^2),$$

where  $\nabla \mathbf{g}$  is known as the *Jacobian matrix*. Given that  $\mathbf{g}(\mathbf{X}) \in \mathcal{O}(1)$  (i.e.  $\mathcal{O}[\Delta t^0]$ ) for general  $\mathbf{X}$ , we have  $\nabla \mathbf{g}(\mathbf{Z}_{n+1}) \in \mathcal{O}(1)$  also; thus the dependence of  $\mathbf{g}[\mathbf{Z}(t_{n+1})]$  on  $\Delta t$  reveals the size of the local error  $\epsilon$ .

We begin by introducing the abbreviations

$$\zeta := \langle \mathbf{Z}(t_{n+1})\mathbf{Z}(t_n)\Delta t \rangle \quad \text{and} \quad \tilde{\zeta} := \langle \mathbf{Z}(t_{n+1})\mathbf{Z}(t_n)\mathbf{F}\Delta t \rangle,$$

and setting  $\mathbf{X} = \mathbf{Z}(t_{n+1})$  in (5.2) to get

$$\mathbf{g}[\mathbf{Z}(t_{n+1})] := \mathcal{B}(\zeta)\mathbf{Z}_n + \mathbf{Z}_F(\tilde{\zeta}) - \mathbf{Z}(t_{n+1}). \tag{5.3}$$

Assuming that  $\mathbf{Z}(t)$  is analytic in a neighbourhood of  $t_n$ , we have

$$\mathbf{Z}(t_{n+1}) = \mathbf{Z}(t_n + \Delta t) = \sum_{s=0}^{\infty} \frac{\mathbf{Z}^{(s)}(t_n)}{s!} \Delta t^s, \tag{5.4}$$

where  $(\cdot)^{(s)} \equiv \frac{d^s}{dt^s} \{(\cdot)\}$ . We now express (2.21) in matrix form as

$$\dot{\mathbf{Z}} = \Psi \mathbf{Z} + \tilde{\mathbf{F}}, \quad \text{where } \Psi(t) = \begin{pmatrix} \mathbf{0}_{3N} & \mathbf{M}^{-1} \\ -\mathcal{F} & \mathbf{0}_{3N} \end{pmatrix} \quad \text{and} \quad \tilde{\mathbf{F}}(t) = \begin{Bmatrix} \mathbf{0} \\ \mathbf{F} \end{Bmatrix}. \tag{5.5}$$

We now derive the series solution of (5.5)<sub>1</sub> for known data at time  $t_n$ . By repeated differentiation of (5.5)<sub>1</sub>, we can express the derivative  $\mathbf{Z}^{(s+1)}$  in terms of lower-order derivatives of  $\mathbf{Z}$  and  $\Psi$ , i.e.

$$\begin{aligned} \ddot{\mathbf{Z}} &= \dot{\Psi} \mathbf{Z} + \Psi \dot{\mathbf{Z}} + \dot{\tilde{\mathbf{F}}}, & \mathbf{Z}^{(3)} &= \ddot{\Psi} \mathbf{Z} + 2\dot{\Psi} \dot{\mathbf{Z}} + \Psi \ddot{\mathbf{Z}} + \dot{\tilde{\mathbf{F}}}^{(2)}, \\ \mathbf{Z}^{(4)} &= \Psi^{(3)} \mathbf{Z} + 3\ddot{\Psi} \dot{\mathbf{Z}} + 3\dot{\Psi} \ddot{\mathbf{Z}} + \Psi \mathbf{Z}^{(3)} + \tilde{\mathbf{F}}^{(3)}, \dots \end{aligned}$$

and so on. Summarising this procedure, we have

$$\mathbf{Z}^{(s+1)} = \sum_{r=0}^s \binom{s}{r} \Psi^{(s-r)} \mathbf{Z}^{(r)} + \tilde{\mathbf{F}}^{(s)}, \tag{5.6}$$

where  $\binom{s}{r}$  represents the binomial coefficient  $\frac{s!}{(s-r)!r!}$ . Inserting (5.6) into (5.4) gives

$$\begin{aligned} \mathbf{Z}(t_{n+1}) &= \mathbf{Z}_n + \sum_{s=0}^{\infty} \frac{\mathbf{Z}_n^{(s+1)}}{(s+1)!} \Delta t^{s+1} = \mathbf{Z}_n + \sum_{s=0}^{\infty} \frac{\Delta t^{s+1}}{(s+1)!} \left[ \sum_{r=0}^s \binom{s}{r} \Psi_n^{(s-r)} \mathbf{Z}_n^{(r)} + \tilde{\mathbf{F}}_n^{(s)} \right] \\ &= \mathbf{Z}_n + \sum_{s=0}^{\infty} \frac{\Delta t^{s+1}}{(s+1)!} \left( \Psi_n^{(s)} \mathbf{Z}_n + \tilde{\mathbf{F}}_n^{(s)} \right) + \sum_{s=0}^{\infty} \frac{\Delta t^{s+1}}{(s+1)!} \sum_{r=1}^s \binom{s}{r} \Psi_n^{(s-r)} \mathbf{Z}_n^{(r)} \\ &= \mathbf{Z}_n + \sum_{s=0}^{\infty} \frac{\Delta t^{s+1}}{(s+1)!} \left( \Psi_n^{(s)} \mathbf{Z}_n + \tilde{\mathbf{F}}_n^{(s)} \right) + \sum_{s=0}^{\infty} \frac{\Delta t^{s+2}}{(s+2)!} \sum_{r=0}^s \binom{s+1}{r+1} \Psi_n^{(s-r)} \mathbf{Z}_n^{(r+1)}, \end{aligned}$$

where  $\sum_{r=1}^s (\cdot) := 0$  for  $s < r$ . Continuing in the same manner, we now substitute for  $\mathbf{Z}_n^{(r+1)}$  using (5.6) to get

$$\begin{aligned} \mathbf{Z}(t_{n+1}) &= \mathbf{Z}_n + \sum_{s=0}^{\infty} \frac{\Delta t^{s+1}}{(s+1)!} \left( \Psi_n^{(s)} \mathbf{Z}_n + \tilde{\mathbf{F}}_n^{(s)} \right) + \sum_{s=0}^{\infty} \frac{\Delta t^{s+2}}{(s+2)!} \sum_{r=0}^s \binom{s+1}{r+1} \Psi_n^{(s-r)} \left( \Psi^{(r)} \mathbf{Z}_n + \tilde{\mathbf{F}}^{(r)} \right) \\ &\quad + \sum_{s=0}^{\infty} \frac{\Delta t^{s+3}}{(s+3)!} \sum_{r=0}^s \binom{s+2}{r+2} \Psi_n^{(s-r)} \sum_{q=0}^r (r+1q+1) \Psi^{(r-q)} \mathbf{Z}_n^{(q+1)}. \end{aligned}$$

The expression for  $\mathbf{Z}(t_{n+1})$  can now be seen to consist of a sum of terms of the form

$$\left[ \sum_{s_0=0}^{\infty} \frac{\Delta t^{s_0+m+1}}{(s_0+m+1)!} \sum_{s_1=0}^{s_0} \binom{s_0+m}{s_1+m} \Psi_n^{(s_0-s_1)} \sum_{s_2=0}^{s_1} \binom{s_1+m-1}{s_2+m-1} \Psi^{(s_1-s_2)} \dots \sum_{s_m=0}^{s_{m-1}} \binom{s_{m-1}+1}{s_m+1} \Psi^{(s_{m-1}-s_m)} \right] \mathcal{V} \tag{5.7}$$

for  $m = 0, 1, 2, \dots$ , where  $\mathcal{V}$  stands for either  $\Psi^{(s_m)} \mathbf{Z}_n$  or  $\tilde{\mathbf{F}}^{(s_m)}$ . We can write (5.7) more compactly as

$$\sum_{s_0=0}^{\infty} \frac{\Delta t^{s_0+m+1}}{(s_0+m+1)!} \left( \prod_{i=0}^{m-1} \left[ \sum_{s_{i+1}=0}^{s_i} \binom{s_i+m-i}{s_{i+1}+m-i} \Psi_n^{(s_i-s_{i+1})} \right] \right) \mathcal{V},$$

with  $s_{-1} := -\infty$  in the case  $m = 0$ . Thus we can express the solution to (5.5)<sub>1</sub> at time  $t_{n+1}$  in (relatively) compact form as

$$\mathbf{Z}(t_{n+1}) = \mathcal{B}^c \mathbf{Z}_n + \mathbf{Z}_F^c, \tag{5.8}$$

where

$$\begin{aligned} \mathcal{B}^c &:= \sum_{m=0}^{\infty} \left[ \sum_{s_0=0}^{\infty} \frac{\Delta t^{s_0+m+1}}{(s_0+m+1)!} \left( \prod_{i=0}^{m-1} \left[ \sum_{s_{i+1}=0}^{s_i} \binom{s_i+m-i}{s_{i+1}+m-i} \Psi_n^{(s_i-s_{i+1})} \right] \Psi_n^{(s_m)} \right) \right] + \mathbf{I}_{6N} \quad \text{and} \\ \mathbf{Z}_F^c &:= \sum_{m=0}^{\infty} \left[ \sum_{s_0=0}^{\infty} \frac{\Delta t^{s_0+m+1}}{(s_0+m+1)!} \left( \prod_{i=0}^{m-1} \left[ \sum_{s_{i+1}=0}^{s_i} \binom{s_i+m-i}{s_{i+1}+m-i} \Psi_n^{(s_i-s_{i+1})} \right] \tilde{\mathbf{F}}_n^{(s_m)} \right) \right]. \end{aligned} \tag{5.9}$$

(Note that for  $m = 0$ , we have  $\prod_{i=0}^{m-1}(\cdot)_i := \mathbf{I}_{6N}$ .) Eqs. (5.8) and (5.9) thus give the exact solution to the semi-discrete equations of motion (2.21), to which we can compare the solutions obtained from our algorithms. From (5.3), (5.8) and (5.9) we can see that the order of  $\mathbf{g}[\mathbf{Z}(t_{n+1})]$  with respect to  $\Delta t$  is governed by how closely  $\mathcal{B}(\zeta)$  and  $\mathbf{Z}_F^c(\zeta)$  match  $\mathcal{B}^c$  and  $\mathbf{Z}_F^c$ , respectively. Using (5.3), we see that since  $\mathbf{Z}_n \in \mathcal{O}(1)$ ,

$$\mathbf{g}[\mathbf{Z}(t_{n+1})] \in \mathcal{O}(\Delta t^{p+1}) \iff \mathcal{B}(\zeta) - \mathcal{B}^c \in \mathcal{O}(\Delta t^{p+p_1}), \quad \mathbf{Z}_F^c(\zeta) - \mathbf{Z}_F^c \in \mathcal{O}(\Delta t^{p+p_2}),$$

where  $p_1, p_2 \in \mathbb{Z}^+$  are such that  $\text{Min}\{p_1, p_2\} = 1$ , and thus

$$\epsilon \in \mathcal{O}(\Delta t^{p+1}) \iff \mathcal{B}(\zeta) - \mathcal{B}^c \in \mathcal{O}(\Delta t^{p+p_1}), \quad \mathbf{Z}_F^c(\zeta) - \mathbf{Z}_F^c \in \mathcal{O}(\Delta t^{p+p_2}). \tag{5.10}$$

Eq. (5.10) thus contains the criteria for Algorithm MC to be  $p$ th-order accurate. To express this in terms of parameters  $\mathcal{G}$ ,  $\mathcal{X}$ ,  $\mathcal{M}$ ,  $\mathbf{F}_a$  and  $\mathbf{V}_a$ , we note from (3.12) that

$$\left[ \mathbf{I}_{6N} - \frac{1}{2} \mathcal{C}(\zeta) \right] \mathcal{B}(\zeta) = \left[ \mathbf{I}_{6N} + \frac{1}{2} \mathcal{C}(\zeta) \right] \quad \text{and} \quad \left[ \mathbf{I}_{6N} - \frac{1}{2} \mathcal{C}(\zeta) \right] \mathbf{Z}_F^c(\zeta) = \mathbf{Z}_a(\zeta),$$

with  $\mathcal{C}$  and  $\mathbf{Z}_a$  given in (3.11) (and repeated here for convenience) as

$$\mathcal{C} = \begin{pmatrix} \mathcal{G}^T & \Delta t \mathcal{M}^{-1} \\ -\Delta t \mathcal{X} & -\mathcal{G} \end{pmatrix} \quad \text{and} \quad \mathbf{Z}_a = \Delta t \begin{Bmatrix} -\mathbf{V}_a \\ \mathbf{F}_a \end{Bmatrix}. \tag{5.11}$$

Therefore pre-multiplying (5.10) by  $-\left[ \mathbf{I}_{6N} - \frac{1}{2} \mathcal{C}(\zeta) \right]$  gives us

$$\epsilon \in \mathcal{O}(\Delta t^{p+1}) \iff \left[ \mathbf{I}_{6N} - \frac{1}{2} \mathcal{C}(\zeta) \right] \mathcal{B}^c - \left[ \mathbf{I}_{6N} + \frac{1}{2} \mathcal{C}(\zeta) \right] \in \mathcal{O}(\Delta t^{p+p_1}), \quad \left[ \mathbf{I}_{6N} - \frac{1}{2} \mathcal{C}(\zeta) \right] \mathbf{Z}_F^c - \mathbf{Z}_a(\zeta) \in \mathcal{O}(\Delta t^{p+p_2}), \tag{5.12}$$

since the left-hand side is unaffected provided that  $\mathcal{C}(\zeta)$  is bounded as  $\Delta t \rightarrow 0$ . We assume that these parameter matrices are also analytic functions of  $t$  within the same neighbourhood of  $t_n$  as  $\mathbf{Z}(t)$ , and thus we can write

$$\mathcal{C}(\zeta) = \sum_{s=0}^{\infty} \mathcal{C}_s(t_n) \Delta t^s \quad \text{and} \quad \mathbf{Z}_a(\zeta) = \sum_{s=0}^{\infty} \mathbf{Z}_{a,s}(t_n) \Delta t^s, \tag{5.13}$$

where the coefficients are wholly determined at time  $t_n$ . This implies similar series expansions for  $\mathcal{G}$ ,  $\mathcal{X}$ ,  $\mathcal{M}^{-1}$ ,  $\mathbf{F}_a$  and  $\mathbf{V}_a$ , and we can relate the expansion of  $\mathcal{M}^{-1}$  to that of  $\mathcal{M}$  via the formula

$$\begin{aligned} \mathcal{M}(\zeta) &= \sum_{s=0}^{\infty} \mathcal{M}_s \Delta t^s = \mathcal{M}_0 \left( \mathbf{I}_{3N} + \sum_{s=0}^{\infty} \mathcal{M}_0^{-1} \mathcal{M}_{s+1} \Delta t^{s+1} \right) \Rightarrow \mathcal{M}^{-1}(\zeta) = \left( \mathbf{I}_{3N} + \sum_{s=0}^{\infty} \mathcal{M}_0^{-1} \mathcal{M}_{s+1} \Delta t^{s+1} \right)^{-1} \mathcal{M}_0^{-1} \\ &= \sum_{r=0}^{\infty} \left( - \sum_{s=0}^{\infty} \mathcal{M}_0^{-1} \mathcal{M}_{s+1} \Delta t^{s+1} \right)^r \mathcal{M}_0^{-1}, \end{aligned} \tag{5.14}$$

which is derived using the binomial formula (e.g. [34]) and valid for sufficiently small  $\Delta t$ . From (5.9), we also have

$$\mathcal{B}^c = \sum_{s=0}^{\infty} \mathcal{B}_s^c(t_n) \Delta t^s \quad \text{and} \quad \mathcal{Z}_F^c = \sum_{s=0}^{\infty} \mathcal{Z}_{F,s}^c(t_n) \Delta t^s. \tag{5.15}$$

Therefore (5.12), (5.13) and (5.15) combine to give

$$\epsilon \in \mathcal{O}(\Delta t^{p+1}) \iff \mathcal{B}_s^c - \frac{1}{2} \sum_{r=0}^s \mathcal{C}_r \mathcal{B}_{s-r}^c - \frac{1}{2} \mathcal{C}_s = \begin{cases} \mathbf{I}_{6N} & : s = 0, \\ \mathbf{0}_{6N} & : s = 1, \dots, p \end{cases} \tag{5.16}$$

and  $\mathcal{Z}_{F,s}^c - \frac{1}{2} \sum_{r=0}^s \mathcal{C}_r \mathcal{Z}_{F,s-r}^c = \mathbf{Z}_{a,s}$  for  $s = 0, \dots, p$ .

We will now derive cumulative criteria for  $p$ th-order accuracy when  $p = 0, 1, 2$ .

5.1. Zeroth-order accuracy

For  $p = 0$ , we have

$$\mathcal{B}_0^c - \frac{1}{2} \mathcal{C}_0 \mathcal{B}_0^c - \frac{1}{2} \mathcal{C}_0 = \mathbf{I}_{6N} \quad \text{and} \quad \mathcal{Z}_{F,0}^c - \frac{1}{2} \mathcal{C}_0 \mathcal{Z}_{F,0}^c = \mathbf{Z}_{a,0}.$$

From (5.9) we have

$$\mathcal{B}_0^c = \mathbf{I}_{6N} \quad \text{and} \quad \mathcal{Z}_{F,0}^c = \mathbf{0}, \tag{5.17}$$

and so we have zeroth-order accuracy if and only if

$$\mathcal{C}_0 = \mathbf{0}_{6N} \quad \text{and} \quad \mathbf{Z}_{a,0} = \mathbf{0}. \tag{5.18}$$

From (5.11), this equates to

$$\mathcal{G}_0 = \mathbf{0}_{3N}. \tag{5.19}$$

Note that this property is *not* sufficient to assure convergence to the true solution as  $\Delta t \rightarrow 0$ . An algorithm that is zeroth-order accurate will in fact give  $\mathbf{Z}_n = \mathbf{Z}_0 \forall n$  using an infinitesimal time-step size, which can be considered the most basic approximation to the true solution.

5.2. First-order accuracy

For  $p = 1$ , the conditions include those for  $p = 0$  and also

$$\mathcal{B}_1^c - \frac{1}{2} \mathcal{C}_0 \mathcal{B}_1^c - \frac{1}{2} \mathcal{C}_1 \mathcal{B}_0^c - \frac{1}{2} \mathcal{C}_1 = \mathbf{0}_{6N} \quad \text{and} \quad \mathcal{Z}_{F,1}^c - \frac{1}{2} \mathcal{C}_0 \mathcal{Z}_{F,1}^c - \frac{1}{2} \mathcal{C}_1 \mathcal{Z}_{F,0}^c = \mathbf{Z}_{a,1}$$

from (5.16). After incorporating (5.17) and (5.18), these become

$$\mathcal{C}_1 = \mathcal{B}_1^c \quad \text{and} \quad \mathbf{Z}_{a,1} = \mathcal{Z}_{F,1}^c.$$

From (5.9) we see that

$$\mathcal{B}^c = \mathbf{I}_{6N} + \underbrace{\sum_{s_0=0}^{\infty} \frac{\Delta t^{s_0+1}}{(s_0+1)!} \Psi_n^{(s_0)}}_{m=0} + \mathcal{O}(\Delta t^2)$$

and similarly for  $\mathcal{Z}_F^c$ ; thus we have

$$\mathcal{B}_1^c = \Psi_n \quad \text{and} \quad \mathcal{Z}_{F,1}^c = \tilde{F}_n. \tag{5.20}$$

Therefore we have first-order accuracy if and only if

$$\mathcal{C}_1 = \Psi_n \quad \text{and} \quad \mathbf{Z}_{a,1} = \tilde{F}_n. \tag{5.21}$$

Using (5.11) and (5.5)<sub>2,3</sub>, this gives us

$$\mathcal{G}_1 = \mathbf{0}_{3N}, \quad \mathcal{M}_0 = \mathbf{M}, \quad \mathcal{X}_0 = \mathcal{F}_n, \quad \mathbf{F}_{a,0} = \mathbf{F}_n \quad \text{and} \quad \mathbf{V}_{a,0} = \mathbf{0}, \tag{5.22}$$

after using (5.14). This property is known as *consistency*, and implies convergence for stable algorithms.

5.3. Second-order accuracy

For  $p = 2$ , the conditions include those for  $p = 1$  and also

$$\mathcal{B}_2^c - \frac{1}{2}\mathcal{C}_0\mathcal{B}_2^c - \frac{1}{2}\mathcal{C}_1\mathcal{B}_1^c - \frac{1}{2}\mathcal{C}_2\mathcal{B}_0^c - \frac{1}{2}\mathcal{C}_2 = \mathbf{0}_{6N} \quad \text{and}$$

$$\mathbf{Z}_{F,2}^c - \frac{1}{2}\mathcal{C}_0\mathbf{Z}_{F,2}^c - \frac{1}{2}\mathcal{C}_1\mathbf{Z}_{F,1}^c - \frac{1}{2}\mathcal{C}_2\mathbf{Z}_{F,0}^c = \mathbf{Z}_{a,2}$$

from (5.16). After incorporating (5.17), (5.18), (5.20) and (5.21), these reduce to

$$\mathcal{C}_2 = \mathcal{B}_2^c - \frac{1}{2}(\Psi_n)^2 \quad \text{and} \quad \mathbf{Z}_{a,2} = \mathbf{Z}_{F,2}^c - \frac{1}{2}\Psi_n\tilde{\mathbf{F}}_n.$$

From (5.9) we see that

$$\mathcal{B}^c = \mathbf{I}_{6N} + \underbrace{\sum_{s_0=0}^{\infty} \frac{\Delta t^{s_0+1}}{(s_0+1)!} \Psi_n^{(s_0)}}_{m=0} + \underbrace{\sum_{s_0=0}^{\infty} \frac{\Delta t^{s_0+2}}{(s_0+2)!} \sum_{s_1=0}^{s_0} \binom{s_0+1}{s_1+1} \Psi_n^{(s_0-s_1)} \Psi_n^{(s_1)}}_{m=1} + \mathcal{O}(\Delta t^3)$$

and similarly for  $\mathbf{Z}_F^c$ ; thus we have

$$\mathcal{B}_2^c = \frac{1}{2}\dot{\Psi}_n + \frac{1}{2}(\Psi_n)^2 \quad \text{and} \quad \mathbf{Z}_{F,2}^c = \frac{1}{2}\dot{\tilde{\mathbf{F}}}_n + \frac{1}{2}\Psi_n\tilde{\mathbf{F}}_n.$$

Therefore we have second-order accuracy if and only if

$$\mathcal{C}_2 = \frac{1}{2}\dot{\Psi}_n \quad \text{and} \quad \mathbf{Z}_{a,2} = \frac{1}{2}\dot{\tilde{\mathbf{F}}}_n. \tag{5.23}$$

Using (5.11) and (5.5)<sub>2,3</sub>, this becomes

$$\mathcal{G}_2 = \mathcal{M}_1 = \mathbf{0}_{3N}, \quad \mathcal{X}_1 = \frac{1}{2}\dot{\tilde{\mathcal{F}}}_n, \quad \mathbf{F}_{a,1} = \frac{1}{2}\dot{\tilde{\mathbf{F}}}_n \quad \text{and} \quad \mathbf{V}_{a,1} = \mathbf{0}, \tag{5.24}$$

after again using (5.14).

It appears as though this process can be continued indefinitely. This is indeed the case, and can be shown inductively as follows: Suppose Algorithm MC is  $(p - 1)$ th-order accurate for  $p \geq 0$ . Since  $\mathcal{B}_0^c = \mathbf{I}_{6N}$  we have, from (5.16),

$$\mathcal{C}_s = \mathcal{B}_s^c - \frac{1}{2} \sum_{r=0}^{s-1} \mathcal{C}_r \mathcal{B}_{s-r}^c \quad \text{and} \quad \mathbf{Z}_{a,s} = \mathbf{Z}_{F,s}^c - \frac{1}{2} \sum_{r=0}^s \mathcal{C}_r \mathbf{Z}_{F,s-r}^c \tag{5.25}$$

as the requirements for  $p$ th-order accuracy. Since  $\mathcal{C}_r; 0 \leq r \leq s - 1$  and  $\mathbf{Z}_{F,s}, \mathcal{B}_s^c; 0 \leq s < \infty$  are all known, (5.25) can be immediately solved (in order) to furnish  $\mathcal{C}_s$  and  $\mathbf{Z}_{a,s}$ . Hence with appropriate choices for  $\mathcal{G}, \mathcal{X}, \mathcal{M}, \mathbf{F}_a$  and  $\mathbf{V}_a$ , Algorithm MC can be made arbitrarily accurate. The criteria for accuracy up to fifth order are given in Table 1, and are continued up to eighth order in [40]. We note the symmetries of the expressions for  $\mathcal{G}$  and  $\mathbf{V}_a$ , and also for  $\mathcal{X}$  and  $\mathbf{F}_a$ .

From Table 1, we see immediately that for time-integration schemes with constant  $\mathcal{G} := \mathcal{G}_0$  or  $\mathcal{M} := \mathcal{M}_0$ , the limit is second-order accuracy for problems with general strain energy functions  $\phi(l)$ , which is consistent with our results in [32]. Therefore schemes defined by (4.16), that conserve energy elementally, cannot be higher than second-order accurate; in turn, this means that higher-order schemes are unlikely to conserve energy and preserve relative equilibria, as discussed in Section 4.5. We also note that they will not retain the sparsity of time-integration schemes of the form (3.10), and will thus be computationally more expensive.

We now verify that the accuracy requirements given in this section do not conflict with the conservation conditions from Section 4 by encapsulating the fact that higher-order accuracy does not hinder conservation of a physical quantity in the following result:

**Proposition 6.** Let  $\mathbf{f}[\mathbf{Z}(t)]$  be a constant quantity of the motion governed by (5.5). Then

$$\epsilon \in \mathcal{O}(\Delta t^{p+1}) \Rightarrow \mathbf{f}(\mathbf{Z}_{n+1}) - \mathbf{f}(\mathbf{Z}_n) \in \mathcal{O}(\Delta t^{p+q}),$$

where  $q \geq 1$ , assuming  $\mathbf{Z}_n = \mathbf{Z}(t_n)$ . That is to say, any  $p$ th-order algorithm will conserve a constant of motion up to order  $p$  or higher.

The proof is given in Appendix A.7. Hence any algorithm that does not conserve one or more of the constants of motion must be limited in its order of accuracy.

Table 1  
Cumulative conditions for  $p$ th-order accuracy

$p$	Conditions needed
0	$\mathcal{G}_0 = \mathbf{0}_{3N}$
1	$\mathcal{G}_1 = \mathbf{0}_{3N}, \quad \mathcal{X}_0 = \mathcal{F}_n, \quad \mathcal{M}_0 = \mathbf{M}, \quad \mathbf{F}_{a,0} = \mathbf{F}_n, \quad \mathbf{V}_{a,0} = \mathbf{0}$
2	$\mathcal{G}_2 = \mathbf{0}_{3N}, \quad \mathcal{X}_1 = \frac{1}{2}\dot{\mathcal{F}}_n, \quad \mathcal{M}_1 = \mathbf{0}_{3N}, \quad \mathbf{F}_{a,1} = \frac{1}{2}\dot{\mathbf{F}}_n, \quad \mathbf{V}_{a,1} = \mathbf{0}$
3	$\mathcal{G}_3 = \frac{1}{12}\ddot{\mathcal{F}}_n\mathbf{M}^{-1}, \quad \mathcal{X}_2 = \frac{1}{6}\ddot{\mathcal{F}}_n + \frac{1}{12}\mathcal{F}_n\mathbf{M}^{-1}\dot{\mathcal{F}}_n, \quad \mathcal{M}_2 = -\frac{1}{12}\ddot{\mathcal{F}}_n,$  $\mathbf{F}_{a,2} = \frac{1}{6}\ddot{\mathbf{F}}_n + \frac{1}{12}\mathcal{F}_n\mathbf{M}^{-1}\dot{\mathbf{F}}_n, \quad \mathbf{V}_{a,2} = \frac{1}{12}\mathbf{M}^{-1}\dot{\mathbf{F}}_n$
4	$\mathcal{G}_4 = \frac{1}{24}\dddot{\mathcal{F}}_n\mathbf{M}^{-1}, \quad \mathcal{X}_3 = \frac{1}{24}(\mathcal{F}_n^{(3)} + \mathcal{F}_n\mathbf{M}^{-1}\dot{\mathcal{F}}_n + \dot{\mathcal{F}}_n\mathbf{M}^{-1}\dot{\mathcal{F}}_n),$  $\mathcal{M}_3 = -\frac{1}{24}\ddot{\mathcal{F}}_n, \quad \mathbf{F}_{a,3} = \frac{1}{24}(\mathbf{F}_n^{(3)} + \mathcal{F}_n\mathbf{M}^{-1}\dot{\mathbf{F}}_n + \dot{\mathcal{F}}_n\mathbf{M}^{-1}\dot{\mathbf{F}}_n),$  $\mathbf{V}_{a,3} = \frac{1}{24}\mathbf{M}^{-1}\ddot{\mathbf{F}}_n$
5	$\mathcal{G}_5 = \frac{1}{80}\mathcal{F}_n^{(3)}\mathbf{M}^{-1} + \frac{1}{120}\dot{\mathcal{F}}_n\mathbf{M}^{-1}\mathcal{F}_n\mathbf{M}^{-1} + \frac{1}{240}\mathcal{F}_n\mathbf{M}^{-1}\dot{\mathcal{F}}_n\mathbf{M}^{-1},$  $\mathcal{X}_4 = \frac{7}{240}\dot{\mathcal{F}}_n\mathbf{M}^{-1}\dot{\mathcal{F}}_n + \frac{1}{80}(\ddot{\mathcal{F}}_n\mathbf{M}^{-1}\dot{\mathcal{F}}_n + \dot{\mathcal{F}}_n\mathbf{M}^{-1}\ddot{\mathcal{F}}_n) + \frac{1}{120}(\mathcal{F}_n^{(4)} + \mathcal{F}_n\mathbf{M}^{-1}\ddot{\mathcal{F}}_n\mathbf{M}^{-1}\dot{\mathcal{F}}_n),$  $\mathcal{M}_4 = -\frac{1}{60}\ddot{\mathcal{F}}_n - \frac{1}{720}\mathcal{F}_n\mathbf{M}^{-1}\ddot{\mathcal{F}}_n,$  $\mathbf{F}_{a,4} = \frac{7}{240}\dot{\mathcal{F}}_n\mathbf{M}^{-1}\dot{\mathbf{F}}_n + \frac{1}{80}(\mathcal{F}_n\mathbf{M}^{-1}\ddot{\mathbf{F}}_n + \ddot{\mathcal{F}}_n\mathbf{M}^{-1}\dot{\mathbf{F}}_n) + \frac{1}{120}(\mathbf{F}_n^{(4)} + \mathcal{F}_n\mathbf{M}^{-1}\ddot{\mathcal{F}}_n\mathbf{M}^{-1}\dot{\mathbf{F}}_n),$  $\mathbf{V}_{a,4} = \frac{1}{80}\mathbf{M}^{-1}\mathcal{F}_n^{(3)} + \frac{1}{120}\mathbf{M}^{-1}\dot{\mathcal{F}}_n\mathbf{M}^{-1}\dot{\mathbf{F}}_n + \frac{1}{240}\mathbf{M}^{-1}\mathcal{F}_n\mathbf{M}^{-1}\dot{\mathcal{F}}_n\mathbf{M}^{-1}\dot{\mathbf{F}}_n$

### 6. Example momentum-conserving algorithms

We now give a few examples of momentum-conserving algorithms that fit into the general framework provided by Algorithm MC, and discuss their properties.

#### 6.1. Current second-order algorithms

The four momentum-conserving time-integration schemes mentioned briefly in Section 4 all have the form (3.10). To repeat, the symplectic-momentum mid-point algorithm [15,23,24,38] defines

$$\mathcal{X} := \left\langle \frac{\phi'_{ij}(\|\mathbf{r}_{1/2}^{ij}\|)}{\|\mathbf{r}_{1/2}^{ij}\|} \mathbf{I}_3 \right\rangle, \tag{6.1}$$

where  $\mathbf{r}^{ij} := \mathbf{r}^j - \mathbf{r}^i$ . This algorithm neither conserves energy nor preserves relative equilibria in general. We will refer to this algorithm as SM. The energy–momentum mid-point algorithm [8,11,22,41] has

$$\mathcal{X} := \left\langle \frac{\phi_{ij\Delta}}{\frac{1}{2}(l_{ijn+1}^2 - l_{ijn}^2)} \mathbf{I}_3 \right\rangle, \tag{6.2}$$

for  $l_{ij} = \|\mathbf{r}^j - \mathbf{r}^i\|$ . This algorithm conserves energy [8] and also preserves relative equilibria [12,18]. We will refer to this algorithm as EM. The assumed distance method [37,25] is given by

$$\mathcal{X} := \left\langle \frac{\phi'_{ij}(l_{ij1/2})}{l_{ij1/2}} \mathbf{I}_3 \right\rangle, \quad (6.3)$$

where  $l_{ij1/2} = \frac{1}{2}(l_{ijn} + l_{ijn+1})$ . This algorithm is not energy-conserving in general, but coincides with the energy–momentum algorithm for quadratic functions  $\phi_{ij}$  [37,42]: it also preserves relative equilibria [42]. Lastly, the energy-conserving algorithm of [15], which can be constructed by applying a global energy condition to the symplectic-momentum algorithm, has

$$\mathcal{X} := \kappa \left\langle \frac{\phi'_{ij}(\|\mathbf{r}_{1/2}^{ij}\|)}{\|\mathbf{r}_{1/2}^{ij}\|} \mathbf{I}_3 \right\rangle, \quad (6.4)$$

with  $\kappa$  defined by (4.8), which coincides with (4.9). This algorithm does not preserve relative equilibria in general. Each of these algorithms is time-reversible, as can be seen from (4.21). They are also second-order accurate. Note that the second-order accuracy of these algorithms in their given form, with  $\mathbf{F}_a := \mathbf{F}$ , tacitly assumes that the force  $\mathbf{F}$  is constant. In the case of a non-constant force, this definition would need to be altered to  $\mathbf{F}_a := \mathbf{F}_{1/2}$ , as can be seen from Table 1.

We now analyse the paths of relative equilibria produced by the energy–momentum (EM) and assumed distance algorithms; in these cases we have, from Lemma 3,  $\lambda_r = \frac{1}{2}\theta^{\text{RE}} / \tan(\frac{1}{2}\theta^{\text{RE}})$  and  $\lambda_t = 1$ . Thus these algorithms will produce paths of relative equilibria with

$$\lambda_r = \frac{\tan^{-1}(\frac{1}{2}\omega_0\Delta t)}{\frac{1}{2}\omega_0\Delta t} \quad \text{and} \quad \lambda_t = 1. \quad (6.5)$$

These paths do *not* follow the trajectory of the exact solution in general, as mentioned in Section 4.5. In the case of pure translation, however, we have  $\omega_0 = 0$  and thus the exact solution is recovered. For general steady-state motion, the translational part  $\mathbf{R}_k^c$  is computed exactly, and the overall error in the positions  $\mathbf{R}_n$  at time-step  $n$  is given by

$$\mathbf{R}_n - \mathbf{R}(t_n) = [\exp(\lambda_r n \Delta t \hat{\boldsymbol{\Omega}}_0) - \exp(n \Delta t \hat{\boldsymbol{\Omega}}_0)] \bar{\mathbf{R}}_0 \quad (6.6)$$

from (4.17)<sub>1</sub>. This error is contained entirely in the *rotational part* of the motion and, owing to its being defined as a difference in *positions* between the exact and the numerical solution, is necessarily *bounded* by the radius of the structure (as measured from the centre of mass). In other words, while the error in rotation (and consequently the length of the path traced) will accumulate as  $t \rightarrow \infty$ , the error in position will stay bounded and only change periodically in time. Given that  $\tan^{-1} x < x \forall x > 0$ , (6.5)<sub>1</sub> tells us that  $\lambda_r < 1$ , and thus the rotation of the structure lags behind the exact solution, with the error in rotation given as  $(1 - \lambda_r)n\omega_0\Delta t = 0$  at time-step  $n$ .

Note that the local error for  $\lambda_r$  defined by (6.5) turns out to be

$$\mathbf{R}_{n+1} - \mathbf{R}(t_{n+1}) = [\text{cay}(\Delta t \hat{\boldsymbol{\Omega}}_0) - \exp(\Delta t \hat{\boldsymbol{\Omega}}_0)] \bar{\mathbf{R}}_n,$$

where the Cayley transform  $\text{cay}(\omega_0\Delta t) = \mathbf{I} + \frac{1}{1 + \frac{1}{4}\|\omega_0\Delta t\|^2} \left( \Delta t \hat{\boldsymbol{\omega}}_0 + \frac{\Delta t^2}{2} \hat{\boldsymbol{\omega}}_0^2 \right)$  [43] in this equation is only a second-order approximation to the exponential mapping  $\exp(\hat{\boldsymbol{\omega}}_0\Delta t) = \mathbf{I} + \frac{\sin\|\omega_0\Delta t\|}{\|\omega_0\Delta t\|} \Delta t \hat{\boldsymbol{\omega}}_0 + \frac{1 - \cos\|\omega_0\Delta t\|}{\|\omega_0\Delta t\|^2} \Delta t^2 \hat{\boldsymbol{\omega}}_0^2$  [36] and is responsible for the position error.

## 6.2. Higher-order accurate algorithms

To illustrate the theory developed in Sections 4 and 5, we now present a sample fourth-order accurate, momentum-conserving algorithm that is time-reversible and preserves relative equilibria. We will call this algorithm M4; we acknowledge that there are many other possibilities. For fourth-order accuracy we require, from Table 1,

$$\begin{aligned} \mathcal{G} &= \frac{1}{12} \dot{\mathcal{F}}_n \mathbf{M}^{-1} \Delta t^3 + \frac{1}{24} \ddot{\mathcal{F}}_n \mathbf{M}^{-1} \Delta t^4 + \mathcal{O}(\Delta t^5), \\ \mathcal{X} &= \mathcal{F}_n + \frac{1}{2} \dot{\mathcal{F}}_n \Delta t + \left( \frac{1}{6} \ddot{\mathcal{F}}_n + \frac{1}{12} \mathcal{F}_n \mathbf{M}^{-1} \mathcal{F}_n \right) \Delta t^2 \\ &\quad + \frac{1}{24} (\mathcal{F}_n^{(3)} + \mathcal{F}_n \mathbf{M}^{-1} \dot{\mathcal{F}}_n + \dot{\mathcal{F}}_n \mathbf{M}^{-1} \mathcal{F}_n) \Delta t^3 + \mathcal{O}(\Delta t^4), \\ \mathcal{M} &= \mathbf{M} - \frac{1}{12} \mathcal{F}_n \Delta t^2 - \frac{1}{24} \dot{\mathcal{F}}_n \Delta t^3 + \mathcal{O}(\Delta t^4), \\ \mathbf{F}_a &= \mathbf{F}_n + \frac{1}{2} \dot{\mathbf{F}}_n \Delta t + \left( \frac{1}{6} \ddot{\mathbf{F}}_n + \frac{1}{12} \mathcal{F}_n \mathbf{M}^{-1} \mathbf{F}_n \right) \Delta t^2 \\ &\quad + \frac{1}{24} (\mathbf{F}_n^{(3)} + \mathcal{F}_n \mathbf{M}^{-1} \dot{\mathbf{F}}_n + \dot{\mathcal{F}}_n \mathbf{M}^{-1} \mathbf{F}_n) \Delta t^3 + \mathcal{O}(\Delta t^4) \quad \text{and} \end{aligned}$$

$$V_a = \frac{1}{12} M^{-1} \dot{F}_n \Delta t^2 + \frac{1}{24} M^{-1} \ddot{F}_n \Delta t^3 + \mathcal{O}(\Delta t^4) \tag{6.7}$$

We can therefore define

$$\begin{aligned} \mathcal{G} &:= \frac{\Delta t^2}{12} \mathcal{F}_\Delta M^{-1}, \quad \mathcal{M} := M \left( M + \frac{\Delta t^2}{12} \mathcal{F}_{1/2} \right)^{-1} M, \\ \mathcal{X} &:= \tilde{\mathcal{X}} + \frac{\Delta t^2}{12} \tilde{\mathcal{X}} M^{-1} \tilde{\mathcal{X}}, \quad \text{where } \tilde{\mathcal{X}} := \mathcal{F}_{1/2} - \frac{\Delta t}{12} \dot{\mathcal{F}}_\Delta, \\ F_a &:= \frac{1}{\Delta t} \int_{t_n}^{t_{n+1}} F(t) dt + \frac{\Delta t^2}{12} \mathcal{F}_{1/2} M^{-1} F_{1/2} \quad \text{and} \quad V_a := \frac{\Delta t}{12} M^{-1} F_\Delta \end{aligned} \tag{6.8}$$

in order to fulfil the respective accuracy criteria, which are verified by taking a Taylor series expansion of each expression. (The definition of  $F_a$  obviously presumes the function  $F(t)$  to be integrable.) The term  $\mathcal{F}_\Delta = \mathcal{F}_{n+1} - \mathcal{F}_n$  in the definition of  $\mathcal{X}$  involves dot-product calculations when evaluating the entries  $f_{ijn}$  and  $f_{ijn+1}$ , which are known to cause numerical difficulties when dealing with stiff problems [32]. Therefore we do not advocate the use of this algorithm in practice.

Regarding the properties of the algorithm defined by (6.8), we see from the structure of  $\mathcal{F}$  given in (2.17) that

$$\mathcal{I}_N \mathcal{G} = \frac{\Delta t^2}{12} \mathcal{I}_N \mathcal{F}_\Delta M^{-1} = \frac{\Delta t^2}{12} \mathcal{O}_N M^{-1} = \mathcal{O}_N,$$

and similarly  $\mathcal{I}_N \mathcal{X} = \mathcal{O}_N$ . Thus  $\mathcal{G}$  and  $\mathcal{X}$  satisfy the requirements for linear momentum conservation given in (4.2), and for a constant force  $F$ , we have

$$F_a = F + \frac{\Delta t^2}{12} \mathcal{F}_{1/2} M^{-1} F.$$

Hence  $\sum_{i=1}^N F_a^i = \sum_{i=1}^N F^i + \frac{\Delta t^2}{12} \mathcal{I}_N \mathcal{F}_{1/2} M^{-1} F = \mathbf{0}$  whenever  $\sum_{i=1}^N F^i = \mathbf{0}$ . Therefore linear momentum is conserved for  $\sum_{i=1}^N F^i = \mathbf{0}$ , as desired. Angular momentum conservation is also assured for  $F = \mathbf{0}$  due to the symmetry of  $\mathcal{M}$  and  $\mathcal{X}$ . Time reversibility follows from (4.21), since  $\mathcal{F}_\Delta^{(s) \text{ TR}} = -\mathcal{F}_\Delta^{(s)} \forall s$  and  $\mathcal{F}_{1/2}^{\text{TR}} = \mathcal{F}_{1/2}$ .

Under relative equilibrium conditions, we have from (6.8), when  $F = \mathbf{0}$ ,

$$\begin{aligned} \mathcal{G}^{\text{RE}} &= \mathbf{0}_{3N}, \quad \mathcal{M}^{\text{RE}} = M \left( M + \frac{\Delta t^2}{12} \mathcal{F}_0 \right)^{-1} M, \quad F_a^{\text{RE}} = V_a^{\text{RE}} = \mathbf{0} \quad \text{and} \\ \mathcal{X}^{\text{RE}} &= \mathcal{F}_0 + \frac{\Delta t^2}{12} \mathcal{F}_0 M^{-1} \mathcal{F}_0. \end{aligned}$$

When inserted into conditions (4.18), this leads to, after rearranging,

$$\begin{aligned} \frac{\Delta t^2}{12} \mathcal{F}_0 M^{-1} \mathcal{F}_0 \bar{R}_0 &= (\alpha - 1) \mathcal{F}_0 \bar{R}_0, \quad \alpha M \bar{V}_0 = \left( M + \frac{\Delta t^2}{12} \mathcal{F}_0 \right) \bar{V}_0 \\ \text{and} \quad \left( M + \frac{\Delta t^2}{12} \mathcal{F}_0 \right) V_0^c &= \lambda_t M V_0^c, \end{aligned} \tag{6.9}$$

where  $\alpha := \lambda_r \frac{\tan(\frac{1}{2} \theta^{\text{RE}})}{\frac{1}{2} \theta^{\text{RE}}}$  and  $\theta^{\text{RE}} := \lambda_r \omega_0 \Delta t$ .

Using (2.37) and the fact that  $\hat{\Omega}_0^3 = -\omega_0^2 \hat{\Omega}_0$  from (2.32), it can be shown that

$$(M^{-1} \mathcal{F}_0)^s \bar{R}_0 = -\omega_0^{2(s-1)} \hat{\Omega}_0^2 \bar{R}_0 \quad \text{and} \quad (M^{-1} \mathcal{F}_0)^s \bar{V}_0 = \omega_0^{2s} \bar{V}_0 \quad \forall s \in \mathbb{Z}^+.$$

It follows that M4 produces paths of relative equilibria with

$$\lambda_r = \frac{\tan^{-1}(\frac{1}{2} \alpha \omega_0 \Delta t)}{\frac{1}{2} \omega_0 \Delta t} \quad \text{for } \alpha = 1 + \frac{\omega_0^2 \Delta t^2}{12} \quad \text{and} \quad \lambda_t = 1; \tag{6.10}$$

thus the translational part of any relative equilibrium motion is exactly recovered.

### 6.3. Angle-preserving algorithms

In this section, we extend to multi-element problems the ideas introduced in Section 8.4 of [32] for central-force algorithms designed to eliminate the error in rotation (or *period error*) for steady-state problems. It is possible to remove this error entirely by modifying the definitions of  $\mathcal{X}$  and  $\mathcal{M}$ . For  $\mathcal{G}^{\text{RE}} := \mathbf{0}_{3N}$ ,  $F_a^{\text{RE}} := \mathbf{0}$  and  $V_a^{\text{RE}} := \mathbf{0}$  we have, from (4.18)

$$\left( \mathcal{X}^{\text{RE}} - \lambda_r \frac{\tan(\frac{1}{2}\theta^{\text{RE}})}{\frac{1}{2}\theta^{\text{RE}}} \mathcal{F}_0 \right) \bar{\mathbf{R}}_0 = \mathbf{0}, \quad \left( \lambda_r \frac{\tan(\frac{1}{2}\theta^{\text{RE}})}{\frac{1}{2}\theta^{\text{RE}}} \mathcal{M}^{\text{RE}} - \mathbf{M} \right) \bar{\mathbf{V}}_0 = \mathbf{0} \quad \text{and} \quad (\mathbf{M} - \lambda_t \mathcal{M}^{\text{RE}}) \mathbf{V}_0^c = \mathbf{0}. \quad (6.11)$$

To remove the error in rotation, we need to ensure that  $\lambda_r = 1$ , which leads to

$$\mathcal{X}^{\text{RE}} \bar{\mathbf{R}}_0 = \beta^{\text{RE}} \mathcal{F}_0 \bar{\mathbf{R}}_0 \quad \text{and} \quad \mathcal{M}^{\text{RE}} \bar{\mathbf{V}}_0 = \frac{1}{\beta^{\text{RE}}} \mathbf{M} \bar{\mathbf{V}}_0 \quad \text{for} \quad \beta^{\text{RE}} = \frac{\tan(\frac{1}{2}\omega_0 \Delta t)}{\frac{1}{2}\omega_0 \Delta t}. \quad (6.12)$$

Thus algorithms that have  $\mathcal{X}^{\text{RE}} \bar{\mathbf{R}}_0 = \mathcal{F}_0 \bar{\mathbf{R}}_0$  and  $\mathcal{M}^{\text{RE}} \bar{\mathbf{V}}_0 = \mathbf{M} \bar{\mathbf{V}}_0$ , such as EM and the assumed distance method, can be modified via  $\mathcal{X} \mapsto \beta \mathcal{X}$ ,  $\mathcal{M} \mapsto \frac{1}{\beta} \mathcal{M}$  to give (6.12), where  $\beta$  is defined as

$$\beta := \frac{\tan(\frac{1}{2}\theta)}{\frac{1}{2}\theta}, \quad \text{with } \theta \text{ defined so that } \theta^{\text{RE}} = \omega_0 \Delta t. \quad (6.13)$$

(Note that (6.13) is equivalent to (8.3)<sub>2</sub> in [32].) These algorithms will return the exact angle of rotation for general relative equilibrium problems.

We now define the *angle-preserving scheme* EM $\theta$ , arising from EM via the modification

$$\mathcal{G} := \mathbf{0}_{3N}, \quad \mathcal{X} := \beta \langle \xi_{ij} \mathbf{I}_3 \rangle, \quad \mathcal{M} := \frac{1}{\beta} \mathbf{M}, \quad \mathbf{F}_a := \alpha \mathbf{F} \quad \text{and} \quad \mathbf{V}_a := \mathbf{0}, \quad (6.14)$$

with  $\beta$  defined by (6.13)<sub>1</sub> and  $\xi_{ij}$  by (4.14); the definition of  $\theta$  will follow shortly. Note that EM $\theta$  retains the property of elemental energy conservation, being of the form (4.16). From Lemma 3, we have

$$\lambda_r = 1 \quad \text{and} \quad \lambda_t = \beta^{\text{RE}} = \frac{\tan(\frac{1}{2}\omega_0 \Delta t)}{\frac{1}{2}\omega_0 \Delta t}, \quad (6.15)$$

which may be compared with (6.5). Since  $\tan x > x \forall 0 < x < \frac{\pi}{2}$ , the translational movement with this scheme is *faster* than the exact solution. The position errors are given by

$$\mathbf{R}_n - \mathbf{R}(t_n) = (\lambda_t - 1)n\Delta t \mathbf{V}_0^c, \quad (6.16)$$

as seen from (4.17). This error is not bounded in time; thus for a general relative equilibrium problem, EM $\theta$  will accumulate (translational) errors in positions as  $t \rightarrow \infty$ . This is in contrast to schemes given by (3.10), where the position error is bounded in time, but the rotation error accumulates as  $t \rightarrow \infty$ . From (4.17), however, it can be seen that the error in momenta  $\mathbf{P}_n$  is zero.

We desire an algorithm that calculates both translational and rotational components of the positions (and also the momenta) exactly for general steady-state problems: such algorithms are possible, as mentioned in Section 4.4, although they are unlikely to conserve energy (for they can do so neither globally nor elementally, as seen from Lemmas 2 and 3). By way of an example, we define the angle-preserving algorithm A $\theta$  by

$$\mathcal{G} := \mathbf{0}_{3N}, \quad \mathcal{X} := \beta \mathcal{F}_{1/2}, \quad \mathcal{M} := \mathbf{M} + c\Delta t^2 \mathcal{F}_{1/2}, \quad \mathbf{F}_a := \mathbf{F} \quad \text{and} \quad \mathbf{V}_a := \mathbf{0} \quad (6.17)$$

for some  $c(\theta)$  to be determined, with  $\beta$  and  $\theta$  defined by (6.13). Under relative equilibrium conditions, this algorithm immediately satisfies conditions (4.18)<sub>3,5</sub>, and (4.18)<sub>1,4</sub> are satisfied by  $\lambda_r = \lambda_t = 1$  given that  $\mathcal{F}_{1/2}^{\text{RE}} = \mathcal{F}_0$ , and  $\mathcal{F}_0 \mathbf{V}_0^c = \mathbf{0}$  due to the form of  $\mathcal{F}_0$ . Eq. (4.18)<sub>2</sub> then becomes

$$(\lambda_r \beta^{\text{RE}} - 1) \mathbf{M} \bar{\mathbf{V}}_0 + \lambda_r \beta^{\text{RE}} c^{\text{RE}} \Delta t^2 \mathcal{F}_0 \bar{\mathbf{V}}_0 = \mathbf{0},$$

which simplifies to

$$[\lambda_r \beta^{\text{RE}} (1 + c^{\text{RE}} \omega_0^2 \Delta t^2) - 1] \mathbf{M} \bar{\mathbf{V}}_0 = \mathbf{0}.$$

(Note from (2.32) that  $\hat{\mathbf{\Omega}}_0^3 = -\omega_0^2 \hat{\mathbf{\Omega}}_0$ , and that  $\mathcal{F}_0 \bar{\mathbf{R}}_0 = -\hat{\mathbf{\Omega}}_0 \mathbf{M} \bar{\mathbf{V}}_0$  and  $\bar{\mathbf{V}}_0 = \hat{\mathbf{\Omega}}_0 \bar{\mathbf{R}}_0$  from (2.37).) This can be solved for  $c^{\text{RE}}$  to give

$$c^{\text{RE}} = \frac{1 - \lambda_r \beta^{\text{RE}}}{\lambda_r \beta^{\text{RE}} \omega_0^2 \Delta t^2}. \quad (6.18)$$

So for  $c$  defined as

$$c := \frac{\frac{1}{2}\theta - \tan(\frac{1}{2}\theta)}{\theta^2 \tan(\frac{1}{2}\theta)},$$

(6.18) is satisfied for  $\lambda_r = 1$ ; hence (4.18)<sub>2</sub> is also. Note that A $\theta$  does not conserve energy in general, although it does produce exact solutions to all steady-state problems: a fact that could have relevance when designing algorithms for solving stiff, approximately steady-state problems at large time-steps.

We now turn to the requirement on  $\theta$  that  $\theta^{\text{RE}} = \omega_0 \Delta t$ . There are several plausible definitions of  $\theta$  that satisfy this. The simplest involves taking the average of all of the angles moved through by the vectors  $\vec{r}^i, i = 1, \dots, N$  (representing the position of each node relative to the centre of mass) during one time-step. This does not take into account the actual displacement of each node, however, so a better choice would be a weighted average of the angles  $\theta_i$  that takes into account the associated lengths  $\bar{l}_n^i$  and  $\bar{l}_{n+1}^i$ , namely

$$\theta := \frac{\sum_{i=1}^N \bar{l}_{1/2}^i \theta_i}{\sum_{i=1}^N \bar{l}_{1/2}^i} \quad \text{for } \theta_i = \cos^{-1} \left( \frac{\bar{r}_n^i \cdot \bar{r}_{n+1}^i}{\bar{l}_n^i \bar{l}_{n+1}^i} \right) \text{ and } \bar{l}^i := \|\bar{r}^i\|. \tag{6.19}$$

For a relative equilibrium problem, we have  $\theta_i = \omega_0 \Delta t \forall i$  from (4.17), and thus  $\theta^{\text{RE}} = \omega_0 \Delta t$ .

In summary, both algorithms EM $\theta$  and A $\theta$  are designed to eliminate the rotational error in a general steady-state problem. Linear and angular momenta are conserved by these schemes, since the form of  $\mathcal{X}$  and  $\mathcal{M}$  are unchanged. Energy is conserved by EM $\theta$  but not by A $\theta$ . Each is time-reversible, since  $\theta^{\text{TR}} = \theta$ , and also second-order accurate, since

$$\frac{\tan(\frac{1}{2}\theta)}{\frac{1}{2}\theta} = 1 + \frac{1}{12}\theta^2 + \mathcal{O}(\theta^4)$$

with  $\theta \in \mathcal{O}(\Delta t)$ .

### 7. Numerical results

We will now test the new algorithms described in the previous chapter on a couple of model problems, to verify the properties discussed in Sections 4 and 5. We will run each experiment with a range of time-step sizes, and assess the relative errors in the positions  $\mathbf{R}$  and velocities  $\mathbf{V}$  at two different sampling times, in order to verify the results. These errors are calculated as

$$\frac{\|\mathbf{R}_n - \mathbf{R}(t_n)\|}{\|\mathbf{R}(t_n)\|} \quad \text{and} \quad \frac{\|\mathbf{V}_n - \mathbf{V}(t_n)\|}{\|\mathbf{V}(t_n)\|}$$

respectively, where  $\{\mathbf{R}_n, \mathbf{V}_n\}$  denotes the approximate solution and  $\{\mathbf{R}(t_n), \mathbf{V}(t_n)\}$  the reference solution. The reference solution for each problem was obtained by running the established energy-conserving scheme EM at a suitably low time-step size, as will be fully described in each case. In order to minimise the effects of round-off error, quadruple precision arithmetic was used when generating the reference solutions, and in some cases for the tests themselves, as will be mentioned.

#### 7.1. Spring-mass system

Our first pair of examples is based on the spring-mass system described in Section 5.1 of [15], shown on the left-hand side of Fig. 1. This consists of four separate masses, with each pair of masses connected by a spring. The strain energy function for each pair is

$$\phi_{ij}(l_{ij}) := \frac{1}{2} k_{ij} (l_{ij} - \bar{l}_{ij})^2, \tag{7.1}$$

where  $k_{ij}$  and  $\bar{l}_{ij}$  represent the spring stiffness and undeformed distance between masses  $i$  and  $j$ , respectively. This strain energy function gives rise to the *Engineering strain* measure (e.g. [44, Section 3.1.1]). The model parameters are  $k_{ij} = 1$ ,

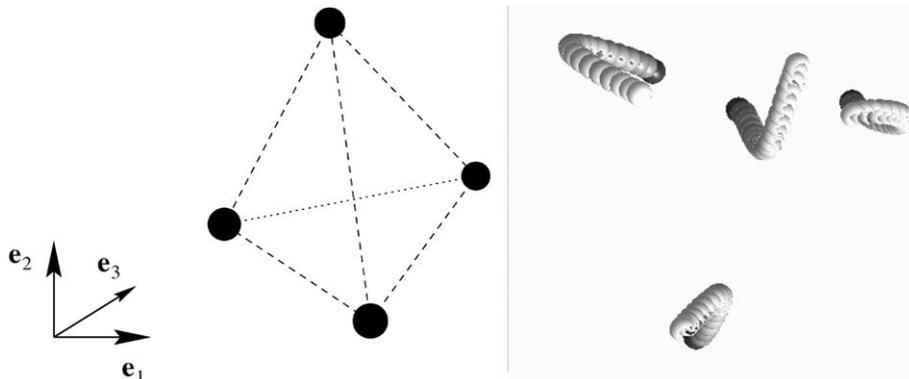


Fig. 1. Spring-mass system at rest and in motion.

$m^i = 1$  and  $\bar{l}_{ij} = 1 \forall i, j$  from [15], thus each spring has a natural length of unity. We note that the problem is *non-stiff*. The initial position and momentum vectors  $\mathbf{r}_0^i$  and  $\mathbf{p}_0^i$  were randomly generated by the authors as

$$\begin{aligned} \mathbf{r}_0^1 &= \langle 0.2340, -0.2166, -0.0109 \rangle, & \mathbf{p}_0^1 &= \langle 0.04095, -0.01483, 0.04325 \rangle, \\ \mathbf{r}_0^2 &= \langle 0.0772, 0.7605, 0.0061 \rangle, & \mathbf{p}_0^2 &= \langle -0.02980, 0.04400, -0.02959 \rangle, \\ \mathbf{r}_0^3 &= \langle 0.8054, 0.6466, -0.1059 \rangle, & \mathbf{p}_0^3 &= \langle -0.02328, -0.01432, -0.03716 \rangle, \\ \mathbf{r}_0^4 &= \langle 0.3903, 0.6187, 0.9678 \rangle & \text{and } \mathbf{p}_0^4 &= \langle 0.04152, 0.00114, 0.02621 \rangle. \end{aligned}$$

We will analyse two versions of this problem: a force-free case (giving a fully conservative system, as used in [15]), and an asymptotically conservative case with a decaying external force  $\mathbf{F}$ : note that this second version involves a *time-dependent* force. The total response time in each case is 30 s, and we measure the relative errors in the solutions obtained at  $t = 10$  and  $t = 30$  s accordingly. The reference solutions for all versions were run using algorithm EM with a time-step size of  $\Delta t = 10^{-5}$  and convergence tolerance of  $10^{-18}$  for the Newton–Raphson iteration. The image on the right-hand side of Fig. 1 depicts the trajectories of the masses during the initial seconds of the motion in the force-free case, as given by the reference solution, with the darkest colour denoting the initial configuration.

This problem is used to show the theoretical properties of the fourth-order accurate algorithm M4: we will therefore test this algorithm alongside the established algorithms SM and EM, for comparison. (Given the nature of the strain energy function in (7.1), the scheme EM and the assumed distance method are equivalent, as was proved in [42].) Each algorithm will be run using the following time-step sizes:  $\Delta t = 1, 0.25, 0.0625$  and  $0.015625$ . All tests will be carried out using double precision arithmetic, with a convergence tolerance of  $10^{-10}$ .

### 7.1.1. Force-free case

Our first version of the spring-mass system has no external force, and is thus identical to that in [15] (and subsequently [42]). Note that the total response time was 50 s in [15]: in all other respects, our results may be compared to those given there. Both energy and momenta are constant, as given in Fig. 2; the small values of the momenta reflect the fact that this problem is dominated by axial vibration. The amount of rotational and translational movement is tiny in comparison, so the group motions are very small. From Fig. 8 of [15] we can estimate the period of axial vibration for each pair of masses to be approximately 5 s: thus all of the time-steps used are small enough to resolve this mode approximately.

Fig. 3 shows the corresponding energy and momenta given by algorithm M4 at the largest time-step size. We can see that, in accordance with the properties listed in Section 6, this scheme conserves linear and angular momenta. (The established algorithms SM and EM also conserve momenta, and EM additionally conserves energy.)

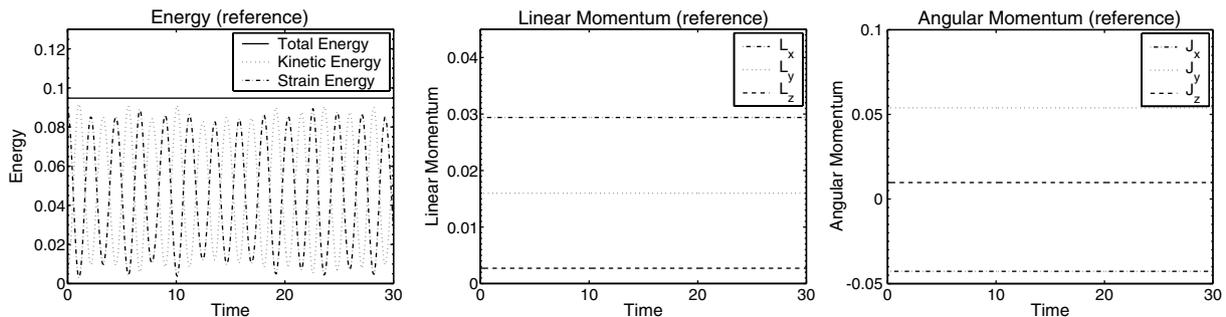


Fig. 2. Reference data for the force-free spring-mass system.

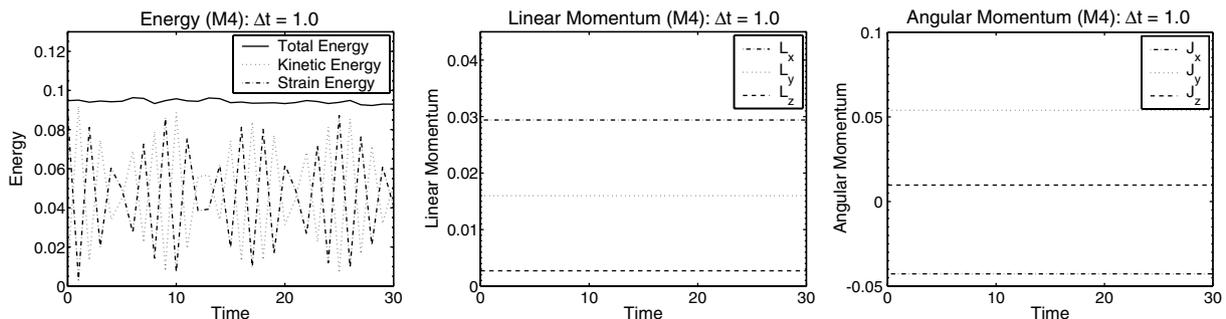


Fig. 3. Energy and momenta as given by M4.

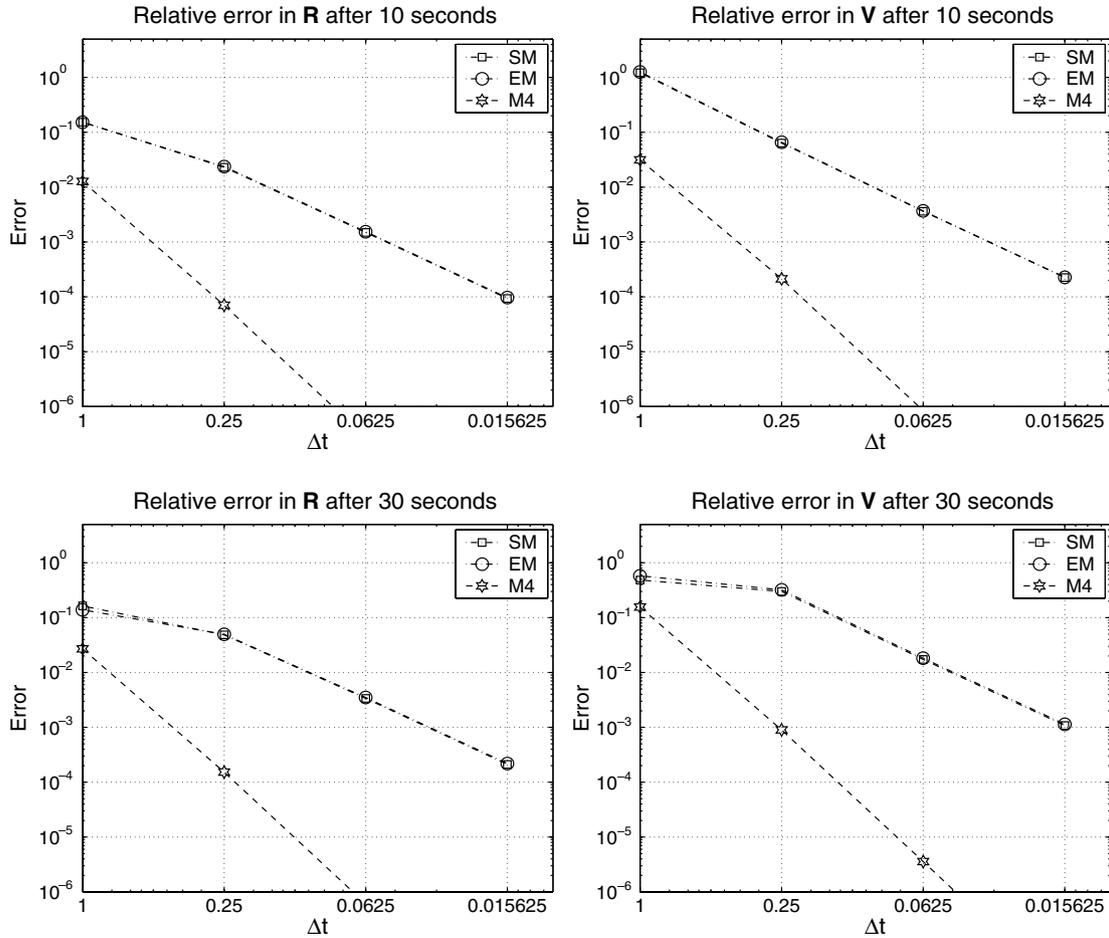


Fig. 4. Relative errors in  $R$  and  $V$  for the force-free spring-mass system.

From Fig. 4; we see that the respective second- and fourth-order accuracy characteristics of these three algorithms are borne out by the error graphs. In particular, the scheme M4 has now demonstrated fourth-order accuracy as well as conservation of momenta. The actual errors in the positions at  $\Delta t = 1$  are acceptably small, being roughly 12–15% for the second-order schemes and 2% for M4. The velocities are not accurately represented by the second-order schemes at this time-step size, however, with errors of over 100% incurred.

7.1.2. Decaying force

For this version of the problem, we add an arbitrary constant force  $F^i$  (acting on mass  $i$ ) for each of the four masses, multiplied by a damping factor in each coordinate direction. The constant components are

$$\begin{aligned} \hat{F}^1 &= \langle 0.0025, 0.01, 0.005 \rangle, & \hat{F}^2 &= \langle 0.005, 0.0075, 0.005 \rangle, \\ \hat{F}^3 &= \langle 0.0075, 0.005, 0.005 \rangle & \text{and } \hat{F}^4 &= \langle 0.01, 0.0025, 0.005 \rangle, \end{aligned}$$

and the actual forces  $F^i \equiv \langle F_x^i, F_y^i, F_z^i \rangle$  are

$$F^i := \langle e^{-t/5} \hat{F}_x^i, e^{-2t/5} \hat{F}_y^i, e^{-3t/5} \hat{F}_z^i \rangle$$

for each mass  $i$ , producing a different effect in each of the three directions. All the other data remain the same as for the force-free case. Therefore the system becomes conservative in the limit  $t \rightarrow \infty$ , with energy and momenta tending to constant values as shown in Fig. 5, where P.A.L. stands for the potential of the applied loads.

Fig. 6 shows the corresponding energy and momenta given by M4. Indeed, the values of the momentum components calculated at the end of the analysis, using the coarse time-step size  $\Delta t = 1$ , are remarkably close to the reference values. The differences in linear momentum are roughly 3% for each of the second-order algorithms; this similarity can be anticipated given that the velocity update for each algorithm is the same. The fourth-order scheme M4 gives the same linear

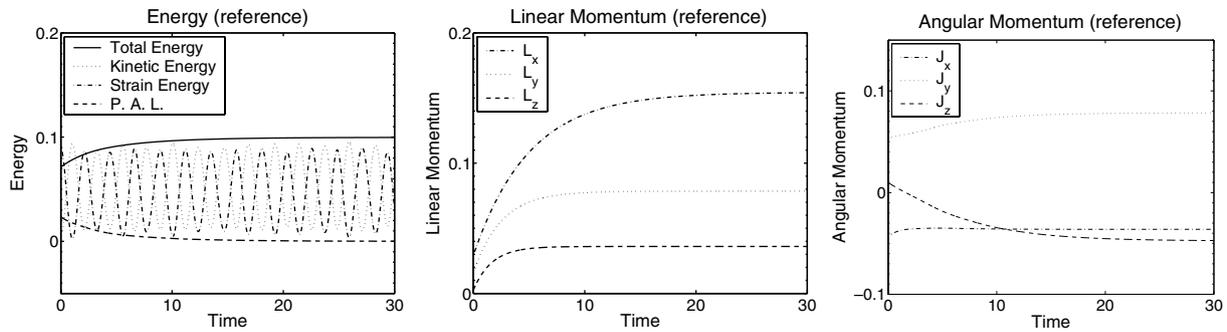


Fig. 5. Reference data for the spring-mass system with a decaying force.

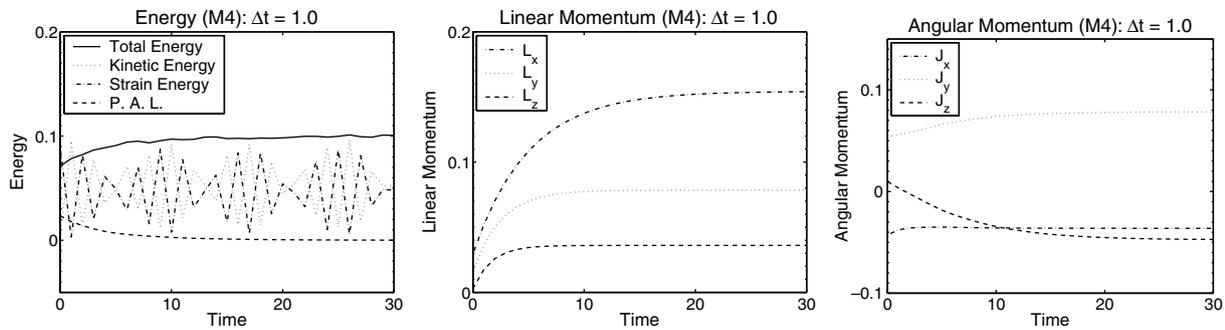


Fig. 6. Energy and momenta as given by M4.

momentum data as the reference solution to nine significant digits. The differences in angular momentum vary, at approximately 1% for EM and 0.06% for M4.

Fig. 7 shows that the relative errors in positions and velocities for the problem with a time-dependent decaying force are actually very similar to those in the force-free case shown in Fig. 4, and thus all observations made previously also apply here. Thus the time-dependence of the forcing term had no effect on accuracy here, and *fourth-order accuracy can be observed for M4 in the case of a non-constant force*. The other algorithms exhibit second-order accuracy, using the mid-point definition  $F_a := F_{1/2}$  for the algebraic force.

### 7.2. Truss structure in relative equilibrium

This example concerns the motion of a truss structure under relative equilibrium conditions, as shown in Fig. 8. It consists of eight bar elements that connect each of the five nodes at  $(0, 1, 0)$ ,  $(1, 0, 0)$ ,  $(0, -1, 0)$ ,  $(-1, 0, 0)$  and  $(0, 0, 0)$  to one another. The structure is therefore two-dimensional, existing entirely within the plane  $z = 0$ . The natural length of the elements is consequently  $\bar{l}_{ij} = 1$  for those aligned with the coordinate axes, and  $\bar{l}_{ij} = \sqrt{2}$  for the rest. The mass per unit length of each bar is equal to unity, hence the element masses are equal in magnitude to their natural lengths. Each element has the strain energy function

$$\phi_{ij}(l_{ij}) := \frac{1}{2} k_{ij} \left( \frac{l_{ij}^2 - \bar{l}_{ij}^2}{2\bar{l}_{ij}} \right)^2, \tag{7.2}$$

which gives rise to *Green's strain* (e.g. [44, Section 3.1.2]). Young's modulus of the material is such that  $E_{ij}A_{ij} = 10^2$  throughout the structure, and so the stiffness of the elements is  $k_{ij} \equiv \frac{E_{ij}A_{ij}}{\bar{l}_{ij}} = 10^2$  for those along an axis and  $k_{ij} = \frac{10^2}{\sqrt{2}}$  for the others.

The initial position and momentum vectors for relative equilibria are obtained from (2.37). In this example, the centre of mass has initial position  $r_0^c = \mathbf{0}$ , and we set its constant velocity to be  $v_0^c = \langle 0, 0, 0.75 \rangle$ , which is orthogonal to the plane of the structure, and thus aligned with an axis of inertia. The constant angular velocity is  $\omega_0 = \langle 0, 0, -1 \rangle$ , thereby giving  $\omega \times v^c = \mathbf{0}$  as specified in Proposition 2.1. Therefore the initial vectors  $r_0^i$  and  $v_0^i$  are calculated as

$$\begin{aligned} r_0^1 &= \langle 0, 0, 0 \rangle, & v_0^1 &= \langle 0, 0, 0.75 \rangle, \\ r_0^2 &= \langle 0, 1.0052720575, 0 \rangle, & v_0^2 &= \langle 1.0052720575, 0, 0.75 \rangle, \end{aligned}$$

$$\begin{aligned}
 \mathbf{r}_0^3 &= \langle -1.0052720575, 0, 0 \rangle, & \mathbf{v}_0^3 &= \langle 0, 1.0052720575, 0.75 \rangle, \\
 \mathbf{r}_0^4 &= \langle 0, -1.0052720575, 0 \rangle, & \mathbf{v}_0^4 &= \langle -1.0052720575, 0, 0.75 \rangle, \\
 \mathbf{r}_0^5 &= \langle 1.0052720575, 0, 0 \rangle & \text{and} & \mathbf{v}_0^5 &= \langle 0, -1.0052720575, 0.75 \rangle
 \end{aligned}$$

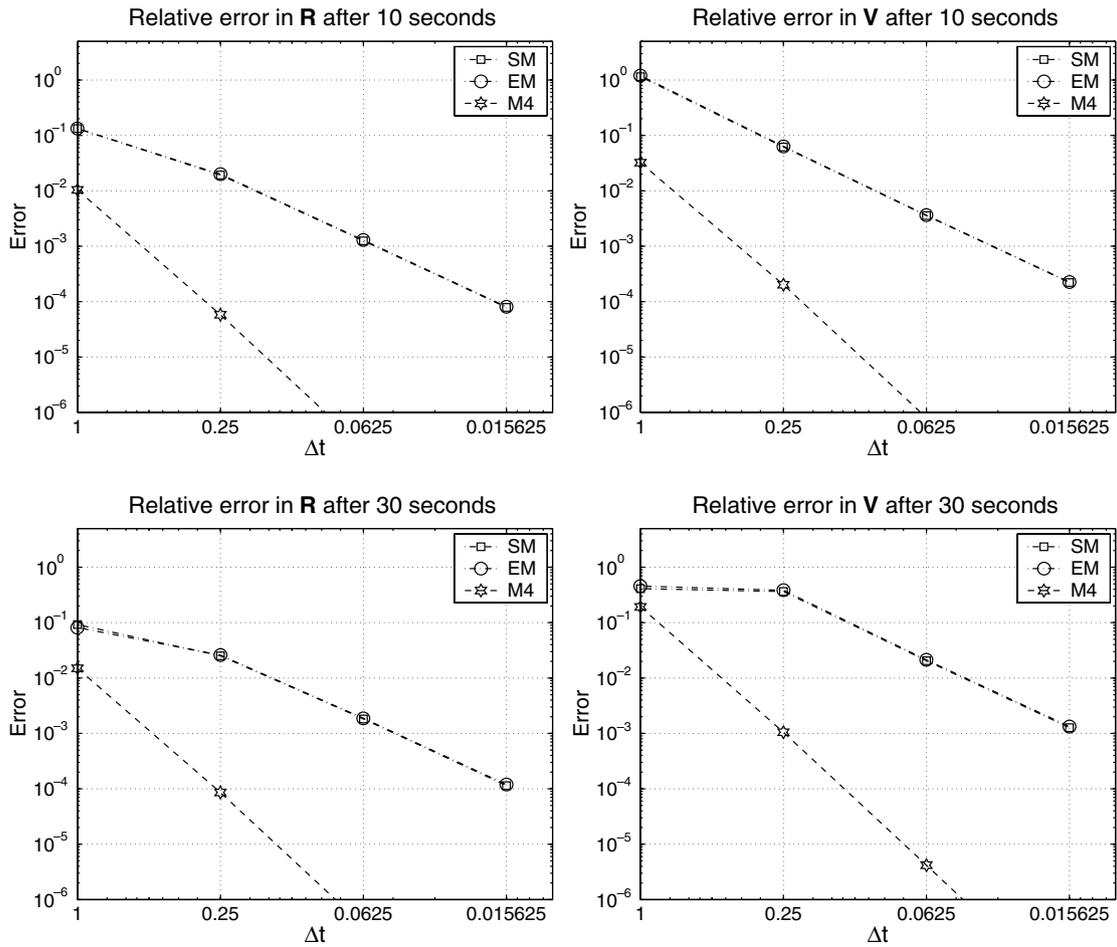


Fig. 7. Relative errors in  $\mathbf{R}$  and  $\mathbf{V}$  for the spring-mass system with a decaying force.

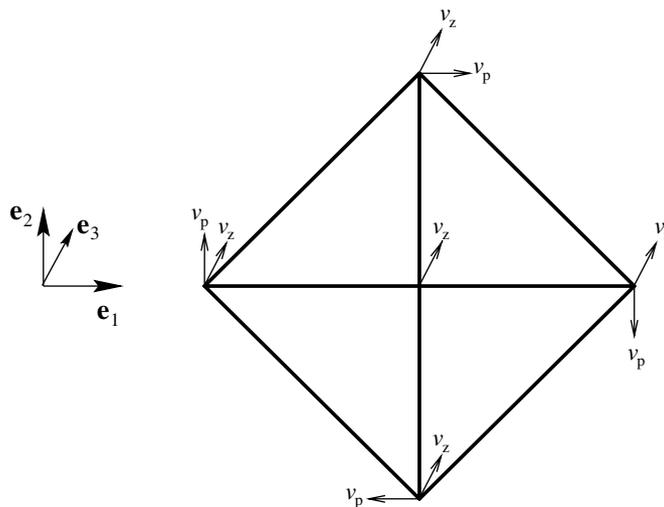


Fig. 8. Truss structure in its initial configuration under relative equilibrium conditions.

to 11 significant digits, with  $P_0 = MV_0$ . This amounts to the structure being stretched by a distance 0.0052720575 along both of its coordinate axes (with unit vectors  $e_1$  and  $e_2$  in the diagram). In the absence of an external force, these starting conditions imply that the structure will undergo uniform in-plane rotation, with simultaneous uniform translation in the direction of  $e_3$ . The image in Fig. 9 depicts this motion during the first few seconds, as given by the reference solution, with the darkest colour denoting the initial configuration.

We will analyse this problem for a total response time of 9 s, and measure the relative errors in the solutions obtained at  $t = 3$  and  $t = 9$  s. The reference solution was run using EM with a time-step size of  $\Delta t = 10^{-6}$  and Newton–Raphson convergence tolerance of  $10^{-15}$ . We note that this system is fully conservative. In addition, the relative equilibrium starting conditions mean that the *individual components of energy remain constant* for this problem, as does the length of each element. This fact is illustrated in Fig. 10, which shows the total energy and the distance  $\|r^2(t) - r^1(t)\|$  as given by the reference solution. Note that the kinetic energy comprises almost all of the total energy, and is thus hidden in the graph.

For the strain energy function given in (7.2), the energy–momentum and assumed distance method of Section 6 are distinct in general. For steady-state problems, however, they coincide [42]. The purpose of this example is to show the accuracy characteristics of the angle-preserving algorithms  $EM\theta$  and  $A\theta$ , so we will run tests with SM and EM for comparison. The time-step sizes remain as given in Section 7.1, and the experiments are conducted using double precision arithmetic, with a convergence tolerance of  $10^{-10}$ .

Fig. 11 shows the relative errors in positions and velocities for each of the algorithms tested. The errors are significantly larger for SM than the other second-order schemes, reflecting the fact that this algorithm does not preserve relative equilibria. A surprising aspect of these results is the similarity between the position errors from EM and  $EM\theta$ . This is to a large

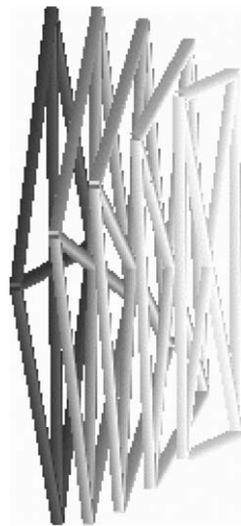


Fig. 9. Truss structure in motion under relative equilibrium conditions.

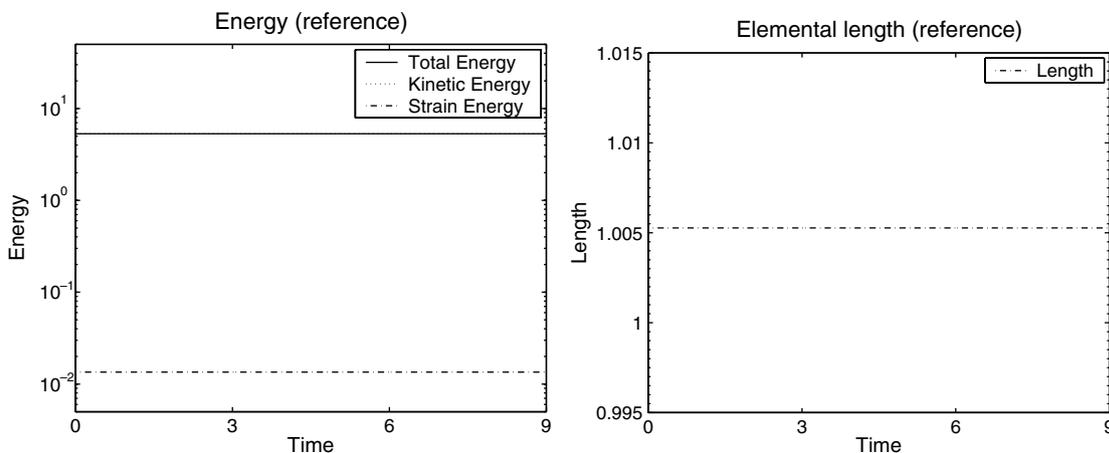


Fig. 10. Reference data for the truss structure in relative equilibrium.

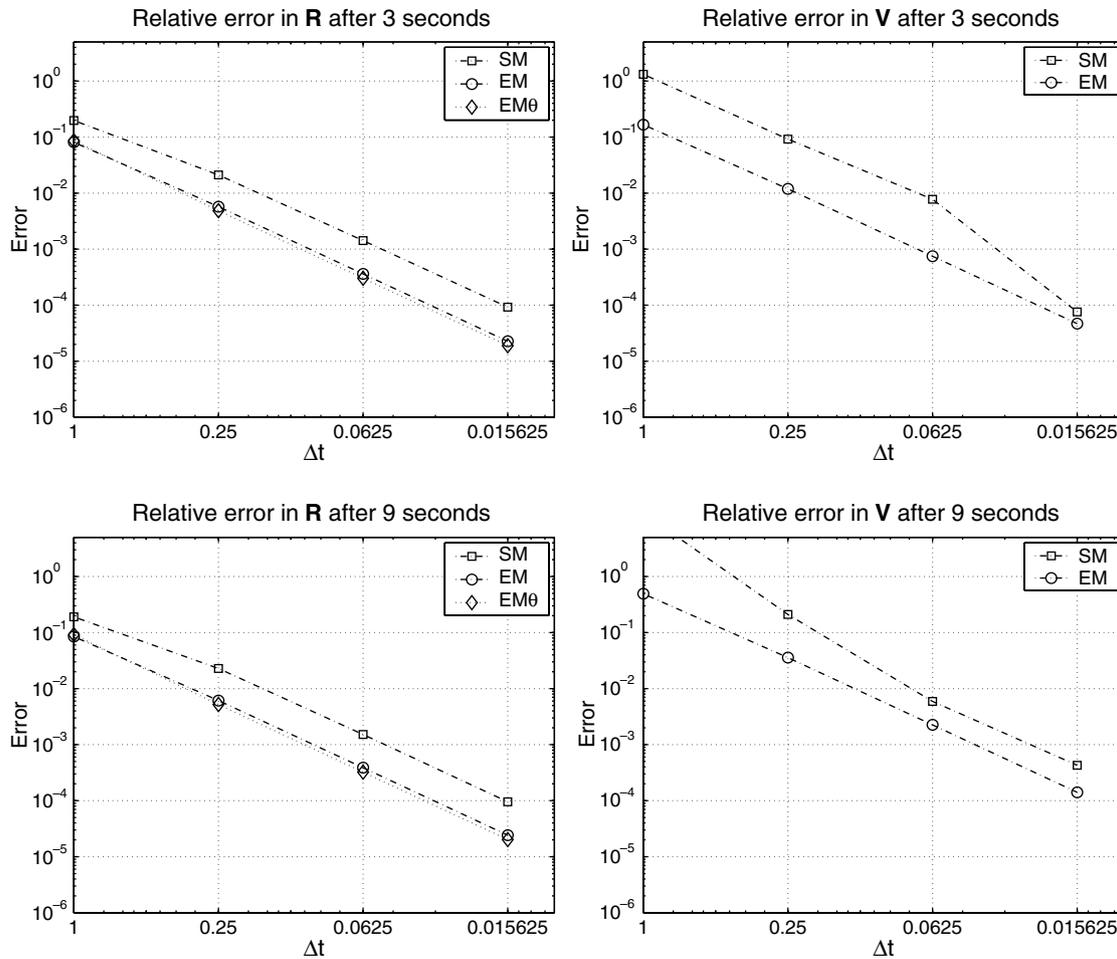


Fig. 11. Relative errors in  $R$  and  $V$  for the truss structure in relative equilibrium.

extent coincidental, given that the errors are from different sources, although representative of the fact that both modes of motion will be present in a general problem.

Table 2 shows that the overall position error for EM $\theta$  consists entirely of translational error: the angle of rotation is reproduced exactly, as expected. Conversely, SM and EM incur some error in the rotational mode, but maintain the exact translation of the structure. We see that both sets of errors (where non-zero) increase linearly with the sampling time, thus confirming that the translational and angular velocities returned by each algorithm are constant. Note that the non-zero period errors indicate that the algorithmic solutions from SM and EM lag *behind* the true solution, whereas the non-zero translational errors show that the solution obtained from EM $\theta$  is *ahead* of it.

We also note the absence of any results for EM $\theta$  in the velocity graphs, as this quantity is calculated exactly (to within machine precision) by this algorithm, as explained in Section 6.3; hence the errors are zero. On the other hand, the error in positions with EM $\theta$  depends on both  $v_0^c$  and  $\omega_0$ , and is unbounded, whereas for EM it depends on  $\omega_0$  only, from (6.5)<sub>1</sub> and (6.6), and stays bounded in time as a consequence of the fact that it is entirely caused by the (unbounded) rotational error. Only algorithm A $\theta$  reproduces both aspects of the motion exactly.

Table 2  
Period and translation errors in positions at  $\Delta t = 0.25$  for the truss structure in relative equilibrium

Algorithm	Period error Sampling time (s)			Translation error Sampling time (s)		
	3	6	9	3	6	9
SM	0.0567	0.114	0.172	0	0	0
EM	0.0155	0.031	0.0464	0	0	0
EM $\theta$	0	0	0	0.0118	0.0236	0.0354
A $\theta$	0	0	0	0	0	0

## 8. Conclusions

In this work, the design of conservative algorithms with higher-order accuracy was investigated for multi-element truss structures (or, equivalently, multi-particle dynamics), and conditions under which they can be developed were given that follow on from our earlier work on the central-force problem [32]. A general framework for algorithms that conserve linear and angular momenta was presented, and conditions for energy conservation were also given. Criteria for the preservation of relative equilibrium states and also for higher-order accuracy were elaborated in detail, and the exact solution to a general non-linear problem was derived (in power-series form).

Time-reversible algorithms that conserve linear and angular momenta can be designed to have arbitrarily high orders of accuracy; these algorithms retain their accuracy properties in the presence of non-conservative external forces. This can be achieved without recourse to extra stages or calculation or additional degrees of freedom in a manner proposed by Argyris et al. [1,29] and LaBudde and Greenspan [30,31], although the resulting systems of equations will no longer be as sparse for higher-order schemes.

Momentum-conserving algorithms in this framework can be designed to conserve energy or preserve relative equilibrium states, although they are unlikely to do both unless they meet the far more restrictive criteria for *elemental* energy conservation, satisfied only by a small group of algorithms based on the second-order energy–momentum algorithm of Simo et al. [8]. However, it is possible to design algorithms that capture relative equilibrium solutions exactly and, in future, we intend to explore the numerical properties of algorithms that can reproduce exact relative equilibrium solutions, in regard to achieving numerical stability and non-linear iterative convergence when solving stiff problems at large time-steps.

## Acknowledgement

This work has been financially supported by the Engineering and Physical Sciences Research Council of Great Britain under a Case Studentship and the Advanced Research Fellowship AF/100089.

## Appendix A. Proofs of algorithmic properties

### A.1. Proposition 1 (linear momentum conservation)

From Algorithm MC, we have

$$\mathcal{L}_\Delta = \mathcal{I}_N \mathbf{P}_\Delta = \mathcal{I}_N (\Delta t \mathbf{F}_a - \Delta t \mathcal{X} \mathbf{R}_{1/2} - \mathcal{G} \mathbf{P}_{1/2}).$$

Since  $\mathbf{R}_{1/2}$  and  $\mathbf{P}_{1/2}$  are arbitrary, we have  $\frac{1}{\Delta t} \mathcal{L}_\Delta = \mathcal{I}_N \mathbf{F}_a$  in general if and only if

$$\mathcal{I}_N \mathcal{X} = \mathcal{I}_N \mathcal{G} = \mathbf{0}_N.$$

Therefore we must have  $\sum_{i=1}^N \mathcal{G}^{ij} = \sum_{i=1}^N \mathcal{X}^{ij} = \mathbf{0}_3 \forall 1 \leq j \leq N$ , as given in (4.2).  $\square$

### A.2. Proposition 2 (angular momentum conservation)

We first note that, for any vectors  $\mathbf{r}, \mathbf{p} \in \mathbb{R}^3$ , we have the identity

$$\mathbf{r}_{n+1} \times \mathbf{p}_{n+1} - \mathbf{r}_n \times \mathbf{p}_n = \mathbf{r}_\Delta \times \mathbf{p}_{1/2} + \mathbf{r}_{1/2} \times \mathbf{p}_\Delta,$$

which can be verified directly. We also note that since  $\mathcal{M}$  comprises unit submatrices, the same is true for  $\mathcal{M}^{-1}$ , i.e.

$$(\mathcal{M}^{-1})^{ij} = \tilde{\mu}^{ij} \mathbf{I}_3.$$

From Algorithm MC and (3.9), we have

$$\begin{aligned} \mathcal{J}_\Delta &= \sum_{i=1}^N \left( \mathbf{r}_\Delta^i \times \mathbf{p}_{1/2}^i + \mathbf{r}_{1/2}^i \times \mathbf{p}_\Delta^i \right) \\ &= \sum_{i=1}^N \left( \left[ \sum_{j=1}^N \left( \Delta t \tilde{\mu}^{ij} \mathbf{p}_{1/2}^j + \mathbf{g}^{ji} \mathbf{r}_{1/2}^j \right) - \Delta t \mathbf{v}_a^i \right] \times \mathbf{p}_{1/2}^i - \mathbf{r}_{1/2}^i \times \left[ \sum_{j=1}^N \left( \Delta t x^{ij} \mathbf{r}_{1/2}^j + \mathbf{g}^{ij} \mathbf{p}_{1/2}^j \right) - \Delta t \mathbf{F}_a^i \right] \right) \\ &= \sum_{i=1}^N \left[ \sum_{j=1}^N \left( \Delta t \tilde{\mu}^{ij} \mathbf{p}_{1/2}^j \times \mathbf{p}_{1/2}^i + \mathbf{g}^{ji} \mathbf{r}_{1/2}^j \times \mathbf{p}_{1/2}^i - \mathbf{g}^{ij} \mathbf{r}_{1/2}^i \times \mathbf{p}_{1/2}^j - \Delta t x^{ij} \mathbf{r}_{1/2}^i \times \mathbf{r}_{1/2}^j \right) + \Delta t \left( \mathbf{r}_{1/2}^i \times \mathbf{F}_a^i + \mathbf{p}_{1/2}^i \times \mathbf{v}_a^i \right) \right] \\ &= \sum_{i=1}^N \left[ \sum_{j=1}^N \left( \Delta t \tilde{\mu}^{ij} \mathbf{p}_{1/2}^j \times \mathbf{p}_{1/2}^i - \Delta t x^{ij} \mathbf{r}_{1/2}^i \times \mathbf{r}_{1/2}^j \right) + \Delta t \left( \mathbf{r}_{1/2}^i \times \mathbf{F}_a^i + \mathbf{p}_{1/2}^i \times \mathbf{v}_a^i \right) \right], \end{aligned}$$

where the vectors  $\mathbf{p}_{1/2}^i \times \mathbf{p}_{1/2}^i$  and  $\mathbf{r}_{1/2}^i \times \mathbf{r}_{1/2}^i$  are unrelated for all  $i$  and  $j$ . Thus we have

$$\frac{1}{\Delta t} \mathcal{J}_\Delta = \sum_{i=1}^N \left( \mathbf{r}_{1/2}^i \times \mathbf{F}_a^i + \mathbf{p}_{1/2}^i \times \mathbf{v}_a^i \right)$$

in general if and only if  $\tilde{\mu}^{ij} = \tilde{\mu}^{ji}$ ,  $x^{ij} = x^{ji}$ . Thus  $\mathcal{M}^{-1}$  and  $\mathcal{X}$  are symmetric, and since

$$\mathcal{M}^{-1} \mathcal{M} = \mathbf{I}_{3N} = (\mathcal{M}^{-1})^\top \mathcal{M}^\top,$$

we see that  $\mathcal{M} = \mathcal{M}^\top \iff \mathcal{M}^{-1} = (\mathcal{M}^{-1})^\top$ ; hence  $\mathcal{M}$  itself must be symmetric.  $\square$

### A.3. Proposition 3 (energy conservation)

We first note the identity

$$\frac{1}{2} \mathbf{P}_{n+1} \cdot \mathbf{M}^{-1} \mathbf{P}_{n+1} - \frac{1}{2} \mathbf{P}_n \cdot \mathbf{M}^{-1} \mathbf{P}_n = \mathbf{P}_\Delta \cdot \mathbf{M}^{-1} \mathbf{P}_{1/2},$$

which can also be easily verified. Algorithm MC then gives

$$H_\Delta = \mathbf{P}_\Delta \cdot \mathbf{M}^{-1} \mathbf{P}_{1/2} + \phi_\Delta - U_\Delta = (\Delta t \mathbf{F}_a - \Delta t \mathcal{X} \mathbf{R}_{1/2} - \mathcal{G} \mathbf{P}_{1/2}) \cdot \mathbf{M}^{-1} \mathbf{P}_{1/2} + \phi_\Delta - U_\Delta, \tag{A.1}$$

from which it can be seen that

$$H_\Delta = -U_\Delta^{\text{NC}} \iff (\Delta t \mathbf{F}_a - \mathcal{G} \mathbf{P}_{1/2}) \cdot \mathbf{M}^{-1} \mathbf{P}_{1/2} + \phi_\Delta - U_\Delta^{\text{C}} = \Delta t \mathcal{X} \mathbf{R}_{1/2} \cdot \mathbf{M}^{-1} \mathbf{P}_{1/2}.$$

Writing  $\mathcal{X} = \kappa \overline{\mathcal{X}}$  now furnishes the result given in (4.8).  $\square$

### A.4. Proposition 4 (preservation of relative equilibria)

We first require a preliminary lemma:

**Lemma A.1.** Let  $\mathbf{A} \in \mathbb{R}^{m \times m}$  be an arbitrary real square matrix, and  $t, t_1, t_2$  be arbitrary real scalars. Define  $\mathbf{B} := \text{diag}(\mathbf{A}) \in \mathbb{R}^{mn \times mn}$  and let  $\mathbf{C} \in \mathbb{R}^{mn \times mn}$  be a real matrix with arbitrary unit submatrices  $\mathbf{C}^{ij} := c_{ij} \mathbf{I}_m \in \mathbb{R}^{m \times m}$ . Then the following statements are true:

- (i)  $\mathbf{A}$  commutes with  $\exp(t\mathbf{A})$ ,
- (ii)  $\exp[(t_1 + t_2)\mathbf{A}] = \exp(t_1\mathbf{A}) \exp(t_2\mathbf{A})$ ,
- (iii)  $\exp(t\mathbf{A})$  is non-singular, with  $\exp(t\mathbf{A})^{-1} = \exp(-t\mathbf{A})$ , and
- (iv)  $\mathbf{B}$  commutes with  $\mathbf{C}$ .

**Proof.** From (2.38), we have

$$\mathbf{A} \exp(t\mathbf{A}) = \mathbf{A} \sum_{s=0}^{\infty} \frac{t^s}{s!} \mathbf{A}^s = \sum_{s=0}^{\infty} \frac{t^s}{s!} \mathbf{A}^{s+1} = \left( \sum_{s=0}^{\infty} \frac{t^s}{s!} \mathbf{A}^s \right) \mathbf{A} = \exp(t\mathbf{A}) \mathbf{A}$$

which establishes (i). Also from (2.38), we have

$$\begin{aligned} \exp(t_1\mathbf{A}) \exp(t_2\mathbf{A}) &= \left( \sum_{q=0}^{\infty} \frac{(t_1)^q}{q!} \mathbf{A}^q \right) \left( \sum_{r=0}^{\infty} \frac{(t_2)^r}{r!} \mathbf{A}^r \right) = \sum_{s=0}^{\infty} \left( \sum_{q+r=s} \frac{(t_1)^q}{q!} \frac{(t_2)^r}{r!} \right) \mathbf{A}^s = \sum_{s=0}^{\infty} \frac{1}{s!} \left( \sum_{r=0}^s \frac{s!}{(s-r)!r!} (t_1)^{s-r} (t_2)^r \right) \mathbf{A}^s \\ &= \sum_{s=0}^{\infty} \frac{(t_1 + t_2)^s}{s!} \mathbf{A}^s = \exp[(t_1 + t_2)\mathbf{A}] \end{aligned}$$

which establishes (ii), and (iii) follows by putting  $t_1 = t$  and  $t_2 = -t$ , since  $\exp(\mathbf{0}_m) := \mathbf{I}_m$ . Finally, for  $\mathbf{B}$  and  $\mathbf{C}$  as given, we have

$$(\mathbf{BC})^{ij} = \sum_{k=1}^n \mathbf{B}^{ik} \mathbf{C}^{kj} = c_{ij} \mathbf{A} \quad \text{since } \mathbf{B}^{ij} = \mathbf{0}_m \text{ for } i \neq j,$$

and also

$$(\mathbf{CB})^{ij} = \sum_{k=1}^n \mathbf{C}^{ik} \mathbf{B}^{kj} = c_{ij} \mathbf{A};$$

hence  $\mathbf{BC} = \mathbf{CB}$  as given in (iv).  $\square$

We will also make use of the identity

$$\exp(\hat{\theta}) - \mathbf{I}_3 = \frac{\tan(\frac{1}{2}\theta)}{\theta} [\exp(\hat{\theta}) + \mathbf{I}_3] \hat{\theta} \quad \text{for } \theta := \|\theta\|, \tag{A.2}$$

which can be established using the formula

$$\exp(\hat{\theta}) = \mathbf{I}_3 + \frac{\sin(\theta)}{\theta} \hat{\theta} + 2 \frac{\sin^2(\frac{1}{2}\theta)}{\theta^2} \hat{\theta}^2 \tag{A.3}$$

given in [36]. From (2.32) we note that  $\hat{\theta}^3 = -\theta^2 \hat{\theta}$ , and so using standard trigonometric identities we get

$$\begin{aligned} [\exp(\hat{\theta}) + \mathbf{I}_3] \hat{\theta} &= 2 \left[ 1 - \sin^2\left(\frac{1}{2}\theta\right) \right] \hat{\theta} + \frac{\sin(\theta)}{\theta} \hat{\theta}^2 = \frac{\cos(\frac{1}{2}\theta)}{\sin(\frac{1}{2}\theta)} \left[ 2 \sin\left(\frac{1}{2}\theta\right) \cos\left(\frac{1}{2}\theta\right) \hat{\theta} + 2 \frac{\sin^2(\frac{1}{2}\theta)}{\theta} \hat{\theta}^2 \right] \\ &= \frac{\theta}{\tan(\frac{1}{2}\theta)} \left[ \frac{\sin(\theta)}{\theta} \hat{\theta} + 2 \frac{\sin^2(\frac{1}{2}\theta)}{\theta^2} \hat{\theta}^2 \right] = \frac{\theta}{\tan(\frac{1}{2}\theta)} [\exp(\hat{\theta}) - \mathbf{I}_3], \end{aligned}$$

which leads to (A.2).

We can now begin the proof of Proposition 4. First, we define the abbreviation

$$\mathbf{E}(k\Delta t) := \exp(\lambda_r k \Delta t \hat{\Omega}_0) \equiv \text{diag}[\exp(\lambda_r k \Delta t \hat{\omega}_0)]. \tag{A.4}$$

From (4.17), we then have

$$\mathbf{R}_k = \mathbf{R}_0^c + \lambda_t k \Delta t \mathbf{V}_0^c + \mathbf{E}(k\Delta t) \bar{\mathbf{R}}_0 \quad \text{and} \quad \mathbf{P}_k = \mathbf{M}[\mathbf{V}_0^c + \mathbf{E}(k\Delta t) \bar{\mathbf{V}}_0] \quad \forall k,$$

which leads to

$$\begin{aligned} \mathbf{R}_\Delta &= \lambda_t \Delta t \mathbf{V}_0^c + (\mathbf{E}[(n+1)\Delta t] - \mathbf{E}(n\Delta t)) \bar{\mathbf{R}}_0, \\ \mathbf{R}_{1/2} &= \mathbf{R}_0^c + \lambda_t \left( n + \frac{1}{2} \right) \Delta t \mathbf{V}_0^c + \frac{1}{2} (\mathbf{E}[(n+1)\Delta t] + \mathbf{E}(n\Delta t)) \bar{\mathbf{R}}_0, \\ \mathbf{P}_\Delta &= \mathbf{M}(\mathbf{E}[(n+1)\Delta t] - \mathbf{E}(n\Delta t)) \bar{\mathbf{V}}_0 \quad \text{and} \\ \mathbf{P}_{1/2} &= \mathbf{M} \mathbf{V}_0^c + \frac{1}{2} \mathbf{M} (\mathbf{E}[(n+1)\Delta t] + \mathbf{E}(n\Delta t)) \bar{\mathbf{V}}_0. \end{aligned} \tag{A.5}$$

From Lemma A.1(ii), we can write  $\mathbf{E}[(n+1)\Delta t] = \mathbf{E}(n\Delta t) \mathbf{E}(\Delta t)$ , and from Lemma A.1(iv) we see that  $\mathbf{E}(k\Delta t)$  and  $\hat{\Omega}_0$  commute with any matrix composed of scaled unit submatrices, which obviously includes  $\mathcal{G}$ ,  $\mathcal{X}$ ,  $\mathcal{M}$  and  $\mathbf{M}$  and their (relevant) inverses.

Thus by inserting (A.5) into (3.8), we obtain

$$\begin{aligned} &\frac{1}{\Delta t} \mathbf{E}(n\Delta t) \left( [\mathbf{E}(\Delta t) - \mathbf{I}_{3N}] + \frac{1}{2} [\mathbf{E}(\Delta t) + \mathbf{I}_{3N}] \mathcal{G}^{\text{RE}} \right) \mathbf{M} \bar{\mathbf{V}}_0 + \frac{1}{\Delta t} \mathcal{G}^{\text{RE}} \mathbf{M} \mathbf{V}_0^c \\ &= -\mathcal{X}^{\text{RE}} \left[ \mathbf{R}_0^c + \lambda_t \left( n + \frac{1}{2} \right) \Delta t \mathbf{V}_0^c \right] - \frac{1}{2} \mathbf{E}(n\Delta t) [\mathbf{E}(\Delta t) + \mathbf{I}_{3N}] \mathcal{X}^{\text{RE}} \bar{\mathbf{R}}_0, \\ &\lambda_t \mathbf{V}_0^c + \frac{1}{\Delta t} \mathbf{E}(n\Delta t) \left( [\mathbf{E}(\Delta t) - \mathbf{I}_{3N}] - \frac{1}{2} [\mathbf{E}(\Delta t) + \mathbf{I}_{3N}] \mathcal{G}^{\text{RET}} \right) \bar{\mathbf{R}}_0 - \frac{1}{\Delta t} \mathcal{G}^{\text{RET}} \left[ \mathbf{R}_0^c + \lambda_t \left( n + \frac{1}{2} \right) \Delta t \mathbf{V}_0^c \right] \\ &= (\mathcal{M}^{\text{RE}})^{-1} \mathbf{M} \mathbf{V}_0^c + \frac{1}{2} \mathbf{E}(n\Delta t) [\mathbf{E}(\Delta t) + \mathbf{I}_{3N}] (\mathcal{M}^{\text{RE}})^{-1} \mathbf{M} \bar{\mathbf{V}}_0 \end{aligned} \tag{A.6}$$

if  $\mathbf{F}_a^{\text{RE}}$  and  $\mathbf{V}_a^{\text{RE}}$  are zero. Now from Propositions 1 and 2 (linear and angular momentum conservation) we know that  $\sum_{j=1}^N \mathcal{X}^{ij} = \mathbf{0}_3 \quad \forall 1 \leq i \leq N$ ; hence

$$\mathcal{X}^{\text{RE}} \mathbf{R}_0^c = \left\langle \sum_{j=1}^N \mathcal{X}^{1j} \mathbf{r}_0^c \quad \dots \quad \sum_{j=1}^N \mathcal{X}^{Nj} \mathbf{r}_0^c \right\rangle = \mathbf{0} \quad \text{and similarly } \mathcal{X}^{\text{RE}} \mathbf{V}_0^c = \mathbf{0}.$$

Also, Proposition 1 shows that  $\mathcal{J}_N \mathcal{G} = \mathcal{C}_N \iff \mathcal{G}^T \mathcal{J}_N^T = \mathcal{C}_N^T$ , and hence

$$\mathcal{G}^{\text{RET}} \mathbf{R}_0^c = \mathcal{G}^{\text{RET}} \mathcal{J}_N^T \mathbf{r}_0^c = \mathbf{0} \quad \text{and similarly } \mathcal{G}^{\text{RET}} \mathbf{V}_0^c = \mathbf{0}. \tag{A.7}$$

Now, Lemma A.1(iii) tells us that  $\mathbf{E}(k\Delta t)^{-1} = \mathbf{E}(-k\Delta t)$ , thus we can multiply each equation in (A.6) by  $\mathbf{E}(-n\Delta t)$  to get, after rearranging,

$$\begin{aligned} \frac{1}{\Delta t} \left( [E(\Delta t) - I_{3N}] + \frac{1}{2} [E(\Delta t) + I_{3N}] \mathcal{G}^{\text{RE}} \right) M \bar{V}_0 + \frac{1}{2} [E(\Delta t) + I_{3N}] \mathcal{X}^{\text{RE}} \bar{R}_0 &= -\frac{1}{\Delta t} E(-n\Delta t) \mathcal{G}^{\text{RE}} M V_0^c, \\ \frac{1}{\Delta t} \left( [E(\Delta t) - I_{3N}] - \frac{1}{2} [E(\Delta t) + I_{3N}] \mathcal{G}^{\text{RE T}} \right) \bar{R}_0 - \frac{1}{2} [E(\Delta t) + I_{3N}] (\mathcal{M}^{\text{RE}})^{-1} M \bar{V}_0 &= E(-n\Delta t) [(\mathcal{M}^{\text{RE}})^{-1} M - \lambda_t I_{3N}] V_0^c. \end{aligned} \tag{A.8}$$

We now invoke the relationship

$$E(\Delta t) - I_{3N} = \frac{\tan(\frac{1}{2}\theta^{\text{RE}})}{\omega_0} [E(\Delta t) + I_{3N}] \hat{\Omega}_0 \quad \text{for } \omega_0 := \|\omega_0\| \text{ and } \theta^{\text{RE}} := \lambda_t \omega_0 \Delta t, \tag{A.9}$$

which is easily deduced from (A.2). Substituting (A.9) into (A.8) brings

$$\begin{aligned} [E(\Delta t) + I_{3N}] \left( \frac{\tan(\frac{1}{2}\theta^{\text{RE}})}{\omega_0 \Delta t} \hat{\Omega}_0 M \bar{V}_0 + \frac{1}{2\Delta t} \mathcal{G}^{\text{RE}} M \bar{V}_0 + \frac{1}{2} \mathcal{X}^{\text{RE}} \bar{R}_0 \right) &= -\frac{1}{\Delta t} E(-n\Delta t) \mathcal{G}^{\text{RE}} M V_0^c, \\ [E(\Delta t) + I_{3N}] \left( \frac{\tan(\frac{1}{2}\theta^{\text{RE}})}{\omega_0 \Delta t} \hat{\Omega}_0 \bar{R}_0 - \frac{1}{2\Delta t} \mathcal{G}^{\text{RE T}} \bar{R}_0 - \frac{1}{2} (\mathcal{M}^{\text{RE}})^{-1} M \bar{V}_0 \right) &= E(-n\Delta t) [(\mathcal{M}^{\text{RE}})^{-1} M - \lambda_t I_{3N}] V_0^c. \end{aligned} \tag{A.10}$$

Now it can be shown by direct calculation that  $[\exp(\lambda_t \Delta t \hat{\omega}_0) + I_3]^{-1}$  exists provided that  $\theta^{\text{RE}} \neq (2m + 1)\pi; m \in \mathbb{Z}$ . Hence for  $-\pi < \theta^{\text{RE}} < \pi$  we can define

$$\tilde{E} = [E(\Delta t) + I_{3N}]^{-1} E(-n\Delta t)$$

and multiply (A.10) by  $[E(\Delta t) + I_{3N}]^{-1}$  to get

$$\begin{aligned} \frac{\tan(\frac{1}{2}\theta^{\text{RE}})}{\omega_0 \Delta t} \hat{\Omega}_0 M \bar{V}_0 + \frac{1}{2\Delta t} \mathcal{G}^{\text{RE}} M \bar{V}_0 + \frac{1}{2} \mathcal{X}^{\text{RE}} \bar{R}_0 &= -\frac{1}{\Delta t} \tilde{E} \mathcal{G}^{\text{RE}} M V_0^c, \\ \frac{\tan(\frac{1}{2}\theta^{\text{RE}})}{\omega_0 \Delta t} \hat{\Omega}_0 \bar{R}_0 - \frac{1}{2\Delta t} \mathcal{G}^{\text{RE T}} \bar{R}_0 - \frac{1}{2} (\mathcal{M}^{\text{RE}})^{-1} M \bar{V}_0 &= \tilde{E} [(\mathcal{M}^{\text{RE}})^{-1} M - \lambda_t I_{3N}] V_0^c. \end{aligned} \tag{A.11}$$

Pre-multiplying either equation in (A.11) by  $\hat{\Omega}_0$  makes the right-hand side disappear, and dot-multiplying either by  $\langle \omega_0 \cdots \omega_0 \rangle$  makes the  $\bar{V}_0$  (or  $\hat{\Omega}_0 \bar{R}_0$ ) terms disappear. Hence each of the  $\bar{R}_0$ ,  $\bar{V}_0$  and  $V_0^c$  terms in (A.11) must be individually zero, i.e.

$$\begin{aligned} \frac{\tan(\frac{1}{2}\theta^{\text{RE}})}{\frac{1}{2}\omega_0 \Delta t} \hat{\Omega}_0 M \bar{V}_0 + \mathcal{X}^{\text{RE}} \bar{R}_0 = 0, \quad \frac{\tan(\frac{1}{2}\theta^{\text{RE}})}{\frac{1}{2}\omega_0 \Delta t} \hat{\Omega}_0 \bar{R}_0 - (\mathcal{M}^{\text{RE}})^{-1} M \bar{V}_0 = 0, \\ \mathcal{G}^{\text{RE}} = \mathbf{0}_{3N} \quad \text{and} \quad [(\mathcal{M}^{\text{RE}})^{-1} M - \lambda_t I_{3N}] V_0^c = 0, \end{aligned} \tag{A.12}$$

since  $\tilde{E}^{-1}$  has been shown to exist. Introducing the initial conditions (2.37) into (A.12)<sub>1</sub> and multiplying (A.12)<sub>2</sub> by  $\mathcal{M}^{\text{RE}}$  results in (4.18). Finally, we note as before that (4.18) assures only that (4.17) is a possible solution returned by Algorithm MC; to guarantee that it will be returned, the algorithm must also give a unique solution for  $R_{n+1}$  and  $P_{n+1}$  given  $R_n$ ,  $P_n$  and  $\Delta t$ , for all  $n$ .  $\square$

#### A.5. Lemma 1 (preservation of relative equilibria and conservation of energy)

Along a relative equilibrium path, we have the difference in energy as

$$H_\Delta = \frac{1}{2} P_{n+1} \cdot M^{-1} P_{n+1} - \frac{1}{2} P_n \cdot M^{-1} P_n + \phi_\Delta$$

from (4.5), since there is no external force. We also have, from (4.17) and (A.4),

$$P_k \cdot M^{-1} P_k = M[V_0^c + E(k\Delta t)\bar{V}_0] \cdot [V_0^c + E(k\Delta t)\bar{V}_0] = M V_0^c \cdot V_0^c + 2ME(k\Delta t)\bar{V}_0 \cdot V_0^c + ME(k\Delta t)\bar{V}_0 \cdot E(k\Delta t)\bar{V}_0$$

by symmetry of  $M$ . From (A.3) it can be shown, using standard trigonometric identities, that for any  $\hat{\theta} \in \mathbb{R}^{3 \times 3}$  we have

$$\exp(\hat{\theta})^{-1} = \exp(\hat{\theta})^T,$$

as is widely known for rotation matrices. Thus  $E(k\Delta t)^{-1} = E(k\Delta t)^T$ , and from (Lemma A.1)(iv) we have

$$ME(k\Delta t)\bar{V}_0 \cdot E(k\Delta t)\bar{V}_0 = E(k\Delta t)^T E(k\Delta t) M \bar{V}_0 \cdot \bar{V}_0 = M \bar{V}_0 \cdot \bar{V}_0;$$

thus we can now write

$$\mathbf{P}_k \cdot \mathbf{M}^{-1} \mathbf{P}_k = \mathbf{M} \mathbf{V}_0^c \cdot \mathbf{V}_0^c + 2\mathbf{M} \mathbf{E}(k\Delta t) \bar{\mathbf{V}}_0 \cdot \mathbf{V}_0^c + \mathbf{M} \bar{\mathbf{V}}_0 \cdot \bar{\mathbf{V}}_0. \quad (\text{A.13})$$

Also, from (2.32), (2.35)<sub>1</sub> and (A.3) we have

$$\exp(\alpha \hat{\omega}_0) \mathbf{v}_0^c = \left[ \mathbf{I}_3 + \frac{\sin(\alpha \omega_0)}{\omega_0} \hat{\omega}_0 + 2 \frac{\sin^2(\frac{1}{2} \alpha \omega_0)}{\omega_0^2} \hat{\omega}_0^2 \right] \mathbf{v}_0^c = \mathbf{v}_0^c \quad \forall \alpha,$$

and since  $\mathbf{E}(k\Delta t)^{-1} = \mathbf{E}(-k\Delta t)$ , we can write

$$\mathbf{M} \mathbf{E}(k\Delta t) \bar{\mathbf{V}}_0 \cdot \mathbf{V}_0^c = \mathbf{M} \bar{\mathbf{V}}_0 \cdot \mathbf{E}(-k\Delta t) \mathbf{V}_0^c = \mathbf{M} \bar{\mathbf{V}}_0 \cdot \mathbf{V}_0^c.$$

Therefore (A.13) becomes

$$\mathbf{P}_k \cdot \mathbf{M}^{-1} \mathbf{P}_k = \mathbf{M} \mathbf{V}_0^c \cdot \mathbf{V}_0^c + 2\mathbf{M} \bar{\mathbf{V}}_0 \cdot \mathbf{V}_0^c + \mathbf{M} \bar{\mathbf{V}}_0 \cdot \bar{\mathbf{V}}_0 = \mathbf{P}_0 \cdot \mathbf{M}^{-1} \mathbf{P}_0 \quad \forall k,$$

and thus  $\frac{1}{2} \mathbf{P}_{n+1} \cdot \mathbf{M}^{-1} \mathbf{P}_{n+1} - \frac{1}{2} \mathbf{P}_n \cdot \mathbf{M}^{-1} \mathbf{P}_n = 0$ .

Finally, since  $\phi = \sum_{i,j=1}^N \phi_{ij}$  is a function of the element lengths  $l_{ij} = \|\mathbf{r}^j - \mathbf{r}^i\|$  only, the difference in strain energy is solely dependent on the differences  $\|\mathbf{r}^j - \mathbf{r}^i\|_\Delta$  in the lengths of the bars. From (4.17) we have

$$\mathbf{r}_k^j = \mathbf{r}_0^j + \lambda_t k \Delta t \mathbf{v}_0^j + \exp(\lambda_r k \Delta t \hat{\omega}_0) (\mathbf{r}_0^j - \mathbf{r}_0^i) \Rightarrow \mathbf{r}_k^j - \mathbf{r}_k^i = \exp(\lambda_r k \Delta t \hat{\omega}_0) (\mathbf{r}_0^j - \mathbf{r}_0^i),$$

and so

$$\|\mathbf{r}_k^j - \mathbf{r}_k^i\| = \sqrt{\exp(\lambda_r k \Delta t \hat{\omega}_0) (\mathbf{r}_0^j - \mathbf{r}_0^i) \cdot \exp(\lambda_r k \Delta t \hat{\omega}_0) (\mathbf{r}_0^j - \mathbf{r}_0^i)} = \|\mathbf{r}_0^j - \mathbf{r}_0^i\| \quad \forall k.$$

Hence  $\phi_\Delta = 0$ , and thus  $H_\Delta = 0$ .  $\square$

#### A.6. Proposition 5 (time reversibility)

From (4.20), we see that Algorithm MC is time-reversible if and only if

$$\mathcal{B}_{n+1}^{\text{TR}} = \mathcal{B}_{n+1}^{-1} \quad \text{and} \quad \mathbf{Z}_F^{\text{TR}} = -\mathcal{B}_{n+1}^{-1} \mathbf{Z}_F. \quad (\text{A.14})$$

Using (3.12), we can write (A.14) as

$$\begin{aligned} \left( \mathbf{I}_{6N} - \frac{1}{2} \mathcal{C}_{n+1}^{\text{TR}} \right)^{-1} \left( \mathbf{I}_{6N} + \frac{1}{2} \mathcal{C}_{n+1}^{\text{TR}} \right) &= \left( \mathbf{I}_{6N} + \frac{1}{2} \mathcal{C}_{n+1} \right)^{-1} \left( \mathbf{I}_{6N} - \frac{1}{2} \mathcal{C}_{n+1} \right) \quad \text{and} \\ \left( \mathbf{I}_{6N} - \frac{1}{2} \mathcal{C}_{n+1}^{\text{TR}} \right)^{-1} \mathbf{Z}_a^{\text{TR}} &= - \left( \mathbf{I}_{6N} + \frac{1}{2} \mathcal{C}_{n+1} \right)^{-1} \mathbf{Z}_a, \end{aligned}$$

which leads to

$$\mathcal{C}_{n+1}^{\text{TR}} = -\mathcal{C}_{n+1} \quad \text{and} \quad \mathbf{Z}_a^{\text{TR}} = -\mathbf{Z}_a.$$

Eq. (4.21) then follows from (3.11).  $\square$

#### A.7. Proposition 6

If an algorithm is  $p$ th-order accurate, then (5.1) shows that

$$\boldsymbol{\epsilon} = \mathbf{Z}_{n+1} - \mathbf{Z}(t_{n+1}) \in \mathcal{O}(\Delta t^{p+1}) \quad (\text{A.15})$$

when  $\mathbf{Z}_n = \mathbf{Z}(t_n)$ . Therefore we have

$$\mathbf{f}(\mathbf{Z}_{n+1}) = \mathbf{f}[\mathbf{Z}(t_{n+1}) + \boldsymbol{\epsilon}] = \mathbf{f}[\mathbf{Z}(t_{n+1})] + [\mathbf{f} \otimes \nabla](t_{n+1}) \cdot \boldsymbol{\epsilon} + \mathcal{O}(\|\boldsymbol{\epsilon}\|^2)$$

from Taylor's theorem, and since  $\mathbf{f}[\mathbf{Z}(t)]$  is a constant of motion, we have  $\mathbf{f}[\mathbf{Z}(t_{n+1})] = \mathbf{f}[\mathbf{Z}(t_n)] = \mathbf{f}(\mathbf{Z}_n)$ . Thus we can write

$$\mathbf{f}(\mathbf{Z}_{n+1}) - \mathbf{f}(\mathbf{Z}_n) = [\mathbf{f} \otimes \nabla](t_{n+1}) \cdot \boldsymbol{\epsilon} + \mathcal{O}(\|\boldsymbol{\epsilon}\|^2),$$

where the matrix  $[\mathbf{f} \otimes \nabla](t_{n+1})$  is independent of  $\Delta t$  (since it is constant). Hence  $\mathbf{f}(\mathbf{Z}_{n+1}) - \mathbf{f}(\mathbf{Z}_n)$  can be at most  $\mathcal{O}(\Delta t^{p+1})$  due to (A.15).  $\square$

## References

- [1] J.H. Argyris, P.C. Dunne, T. Angelopoulos, Dynamic response by large step integration, *Earthquake Engrg. Struct. Dynam.* 2 (1973) 185–203.
- [2] C. Lanczos, *The Variational Principles of Mechanics*, Dover, New York, 1970.
- [3] G. Dahlquist, A special stability problem for linear multistep methods, *BIT* 3 (1963) 27–43.
- [4] T.J.R. Hughes, *The Finite Element method: Linear Static and Dynamic Finite Element Analysis*, Prentice-Hall, Englewood Cliffs, NJ, 1987.
- [5] T. Belytschko, D.F. Schoeberle, On the unconditional stability of an implicit algorithm for nonlinear structural dynamics, *J. Appl. Mech.* 97 (1975) 865–869.
- [6] T.J.R. Hughes, A note on the stability of Newmark’s algorithm in nonlinear structural dynamics, *Int. J. Numer. Methods Engrg.* 6 (1976) 383–386.
- [7] R.D. Krieg, Unconditional stability in numerical time integration methods, *J. Appl. Mech.* (1973) 417–421.
- [8] J.C. Simo, N. Tarnow, The discrete energy–momentum method. Conserving algorithms for non-linear elastodynamics, *Journal of Applied Mathematics and Physics (ZAMP)*. 43 (1992) 757–792.
- [9] J.E. Marsden, M. West, Discrete mechanics and variational integrators, *Acta Numer.* (2001) 357–514.
- [10] G. Zhong, J.E. Marsden, Lie–Poisson Hamilton–Jacobi theory and Lie–Poisson integrators, *Phys. Lett. A* 133 (3) (1988) 134–139.
- [11] J.C. Simo, O. Gonzalez, Assessment of energy–momentum and symplectic schemes for stiff dynamical systems. *Papers—American Society of Mechanical Engineers—All Series*, 93(4), 1993. Presented at the ASME Winter Annual Meeting, New Orleans, Louisiana, November 28–December 3, 1993.
- [12] O. Gonzalez, J.C. Simo, On the stability of symplectic and energy–momentum algorithms for non-linear Hamiltonian systems with symmetry, *Comput. Methods Appl. Mech. Engrg.* 134 (1996) 197–222.
- [13] M. Ortiz, A note on energy conservation and stability of nonlinear time-stepping algorithms, *Comput. Struct.* 24 (1) (1986) 167–168.
- [14] P. Betsch, P. Steinmann, Conservation properties of a time FE method. Part II: Time-stepping schemes for non-linear elastodynamics, *Int. J. Numer. Methods Engrg.* 50 (2001) 1931–1955.
- [15] J.C. Simo, N. Tarnow, K.K. Wong, Exact energy–momentum conserving algorithms and symplectic schemes for non-linear dynamics, *Comput. Methods Appl. Mech. Engrg.* 100 (1992) 63–116.
- [16] J.C. Simo, D. Lewis, J.E. Marsden, Stability of relative equilibria. Part I: The reduced energy–momentum method, *Arch. Rat. Mech. Anal.* 115 (1991) 15–59.
- [17] J.C. Simo, T.A. Posbergh, J.E. Marsden, Stability of relative equilibria. Part II: Application to nonlinear elasticity, *Arch. Rat. Mech. Anal.* 115 (1991) 61–100.
- [18] F. Armero, I. Romero, On the formulation of high-frequency dissipative time-stepping algorithms for nonlinear dynamics. Part I: Low order methods for two model problems and nonlinear elastodynamics, *Comput. Methods Appl. Mech. Engrg.* 190 (2001) 2603–2649.
- [19] F. Armero, E. Petőcz, Formulation and analysis of conserving algorithms for frictionless dynamic contact/impact problems, *Comput. Methods Appl. Mech. Engrg.* 158 (3–4) (1998) 269–300.
- [20] F. Armero, I. Romero, On the formulation of high-frequency dissipative time-stepping algorithms for nonlinear dynamics. Part II: Second order methods, *Comput. Methods Appl. Mech. Engrg.* 190 (2001) 6783–6824.
- [21] H. Yoshida, Construction of higher order symplectic integrators, *Phys. Lett. A* 150 (5) (1990) 262–268.
- [22] S. Reich, Enhancing energy conserving methods, *BIT* 36 (1996) 122–134.
- [23] J.M. Sanz-Serna, M.P. Calvo, *Numerical Hamiltonian Problems*, Chapman and Hall, London, 1994.
- [24] E. Hairer, G. Wanner, *Solving Ordinary Differential Equations II*, second ed., Springer-Verlag, Berlin, Germany, 1996.
- [25] P. Betsch, P. Steinmann, Conservation properties of a time FE method. Part I: Time-stepping schemes for  $N$ -body problems, *Int. J. Numer. Methods Engrg.* 49 (2000) 599–638.
- [26] E. Forest, Sixth-order Lie group integrators, *Journal of Computational Physics* 99 (1992) 209–213.
- [27] N. Tarnow, J.C. Simo, How to render second order accurate time-stepping algorithms fourth order accurate while retaining the stability and conservation properties, *Comput. Methods Appl. Mech. Engrg.* 115 (1994) 233–252.
- [28] P. Betsch, P. Steinmann, Inherently energy conserving time finite elements for classical mechanics, *J. Comput. Phys.* 160 (2000) 88–116.
- [29] J.H. Argyris, P.C. Dunne, T. Angelopoulos, Non-linear oscillations using the finite element technique, *Comput. Methods Appl. Mech. Engrg.* 2 (1973) 203–250.
- [30] R.A. LaBudde, D. Greenspan, Energy and momentum conserving methods of arbitrary order for the numerical integration of equations of motion I. Motion of a single particle, *Numer. Math.* 25 (1976) 323–346.
- [31] R.A. LaBudde, D. Greenspan, Energy and momentum conserving methods of arbitrary order for the numerical integration of equations of motion II. Motion of a system of particles, *Numer. Math.* 26 (1976) 1–16.
- [32] E. Graham, G. Jelenić, A general framework for conservative single-step time-integration schemes with higher-order accuracy for a central-force system, *Comput. Methods Appl. Mech. Engrg.* 192 (2003) 3585–3618.
- [33] J.E. Marsden, T.J.R. Hughes, *Mathematical Foundations of Elasticity*, Dover, Mineola, New York, 1994.
- [34] M.L. Boas, *Mathematical Methods in the Physical Sciences*, John Wiley & Sons, Chichester, England, 1983.
- [35] R.A. Horn, C.R. Johnson, *Matrix Analysis*, first ed., Cambridge University Press, Cambridge, England, 1996.
- [36] J. Argyris, An excursion into large rotations, *Comput. Methods Appl. Mech. Engrg.* 32 (1982) 85–155.
- [37] U.M. Ascher, S. Reich, On some difficulties in integrating highly oscillatory Hamiltonian systems, *Lecture Notes Comput. Sci. Engrg.* 4 (1998) 281–296.
- [38] F. Kang, Difference schemes for Hamiltonian formalism and symplectic geometry, *J. Comput. Math.* 4 (1986) 279–289.
- [39] D. Kuhl, E. Ramm, Constraint energy momentum algorithm and its application to non-linear dynamics of shells, *Comput. Methods Appl. Mech. Engrg.* 136 (1996) 293–315.

- [40] E. Graham, Higher-order accuracy in implicit, conservative, single-step time-integration schemes for non-linear structural dynamics, Ph.D. thesis, Imperial College of Science, Technology and Medicine, London, England, October 2003.
- [41] D. Greenspan, Completely conservative covariant numerical methodology, *Comput. Math. Appl.* 29 (4) (1995) 37–43.
- [42] E. Graham, G. Jelenić, M.A. Crisfield, A note on the equivalence of two recent time-integration schemes for  $N$ -body problems, *Commun. Numer. Methods Engrg.* 18 (2002) 615–620.
- [43] J.C. Simo, N. Tarnow, M. Doblare, Non-linear dynamics of three-dimensional rods: exact energy and momentum conserving algorithms, *Int. J. Numer. Methods Engrg.* 38 (1995) 1431–1473.
- [44] M.A. Crisfield, *Non-linear Finite Element Analysis of Solids and Structures*, vol. 1, John Wiley & Sons, Chichester, UK, 1991.