

# NUMERICALLY STABLE FORMULAS FOR A PARTICLE-BASED EXPLICIT EXPONENTIAL INTEGRATOR\*

PRASHANTH NADUKANDI<sup>†</sup>

**Abstract.** Numerically stable formulas are presented for the closed-form analytical solution of the X-IVAS scheme in 3D. This scheme is a state-of-the-art particle-based explicit exponential integrator developed for the Particle Finite Element Method. Algebraically, this scheme involves two steps: 1) the solution of tangent curves for piecewise linear vector fields defined on simplicial meshes and 2) the solution of line integrals of piecewise linear vector-valued functions along these tangent curves. Hence, the stable formulas presented here have general applicability, e.g. exact integration of trajectories in particle-based (Lagrangian-type) methods, flow visualization and computer graphics. The Newton form of the polynomial interpolation definition is used to express exponential functions of matrices which appear in the analytical solution of the X-IVAS scheme. The divided difference coefficients in these expressions are defined in a piecewise manner, i.e. in a prescribed neighbourhood of removable singularities their series approximations are computed. An optimal series approximation of divided differences is presented which plays a critical role in this methodology. At least ten significant decimal digits in the formula computations are guaranteed to be exact using double-precision floating-point arithmetic. The worst case scenarios occur in the neighbourhood of removable singularities found in fourth-order divided differences of the exponential function.

**Key words.** X-IVAS scheme, particle finite element method, explicit exponential integrators, tangent curves, closed-form analytical solutions, finite arithmetic, loss of significance, numerically stable formulas

**AMS subject classifications.** Primary: 65-04; Secondary: 34A05, 39-04, 65G30, 70B05, 76M10, 76M28

**1. Introduction.** The particle finite element method [11, 16] (PFEM) is a versatile numerical method in which each fluid particle is followed in a Lagrangian manner. It is shown to successfully simulate a wide variety of complex engineering problems, e.g. free-surface/multi-fluid flows with violent interface motions, polymer melting and burning simulations, multi-fluid mixing and buoyancy driven segregation problems, etc. A recent development within the framework of the PFEM is the X-IVAS (eXplicit Integration along the Velocity and Acceleration Streamlines) scheme [10]. Its development was motivated in the need for a faster and more accurate time integrator for incompressible flows.

The X-IVAS scheme targets the explicit time integration of the kinematics of the fluid particles using large time steps. The equations of motion are obtained from the incompressible Navier–Stokes equations subjected to a second-order pressure segregation method. In the Lagrangian formulation the segregated momentum balance equations define the acceleration of the fluid particles. The adopted hypothesis is that a streamline in the configuration at the start of the time step (reference configuration) is a good approximation to the pathline of the fluid particle for a relatively large time step. For a chosen time step each fluid particle is advected along the streamline passing through its position in the reference configuration. The particle velocity is updated for this time step by doing a line integral of the acceleration existing in the reference configuration along the streamline. It is expected that the explicit time integration of the particle position and velocity along the streamline yields a better and more stable approximation than doing so via standard finite difference time integrators.

The X-IVAS hypothesis has been tested [10, 9] by successfully simulating some benchmark CFD and FSI examples using very large time steps, e.g. 10–15 times the standard Courant–Friedrichs–Lewy stability limit. Further, the simulation results of multi-fluid flows

---

\* This study was partially supported by the SAFECON project of the European Research Council (European Commission) and the WAM-V project funded under the Navy Grant N62909-12-1-7101 issued by Office of Naval Research Global. The United States Government has a royalty-free license throughout the world in all copyrightable material contained herein.

<sup>†</sup> Centre Internacional de Mètodes Numèrics en Enginyeria (CIMNE), Edifici C1, Gran Capitan s/n, 08034 Barcelona, Spain (npras@cimne.upc.edu).

were compared [9] with those obtained using OpenFOAM [5] and for similar accuracy the PFEM+X-IVAS method took nearly half the simulation time taken by OpenFOAM.

In the PFEM+X-IVAS method the data is stored in a discrete manner, i.e. all the essential variables are stored as intrinsic properties with the particles. To perform the explicit time integration we need a spatially continuous description of the velocity and acceleration. For this purpose a fixed auxiliary simplicial mesh is used to interpolate the variables of interest in a piecewise linear manner. Thus, the integration of the position and velocity of the fluid particles is actually driven by *approximate piecewise linear vector fields*. Nevertheless, the higher-order fine-scale details of the initial solution and of the integration results are retained with the particles. This is essential for the accuracy of the scheme and avoids the so-called erosion artefact associated to the Eulerian formulations.

Despite the possibility to compute analytically the positions and velocities of the particles, numerical sub-stepping methods based on simple finite-difference schemes (e.g. forward Euler) were used for this purpose in the original proposal and in the subsequent developments until now. On the one hand, this introduces an additional source of numerical error and the repercussions of the same are not well understood. One has to carefully choose an appropriate sub-time-step for stability. On the other hand, numerical sub-stepping might weaken the X-IVAS hypothesis. For instance, a situation originally used to motivate the X-IVAS hypothesis is that the streamlines of a flow never cross fixed impermeable domain boundaries (e.g. vortices near corners of wall bounded flows). The X-IVAS hypothesis guarantees that the material points never leave impermeable domain boundaries. Numerical sub-stepping procedures clearly compromise this guarantee.

Diachin and Herzog reported [3] that analytical solvers provide faster more accurate results for streamline calculations on a linear tetrahedra than the forward Euler and the fourth-order Runge–Kutta methods. The matrix functions were computed therein using a procedure based on matrix decomposition methods. The singular value decomposition was used to determine the matrix rank which in turn was used to classify the calculation procedure into four cases in 3D. Using a Schur decomposition the matrix functions were transformed to equivalent functions of upper triangular matrices. The calculation of the latter was done using a recursive relation proposed by Parlett [18]. Nielson and Jung presented [15] formulas in 2D and 3D to compute the closed-form analytical solution of tangent curves for linearly varying vector fields over tetrahedral domains. We infer from the algebraic structure of these formulas that the matrix functions were expressed therein using the Lagrange form of the polynomial interpolation definition. The analytical solution of the tangent curves were classified into five cases in 2D and nine cases in 3D which depend on the eigenvalues of the system matrix.

Idelsohn et al. presented [10] a procedure to compute the analytical solution of the X-IVAS scheme in 2D. The matrix functions are expressed therein using the Jordan canonical form definition. The solution in 3D was omitted pointing out that the extension to 3D is straightforward. Unlike in 2D where the procedure to express matrices in the Jordan canonical form is straightforward, in 3D (and for matrices of larger dimensions) this procedure is arduous as repeated eigenvalues with different Jordan blocks might exist. As the Jordan structure of a 3D matrix involves multiple cases, we beg to differ with Idelsohn et al. [10] that it is not trivial to derive and implement (code) their approach to compute the analytical solution of the X-IVAS scheme in 3D.

Moreover, using finite precision arithmetic the *as is* computation of the analytical solution procedure of Idelsohn et al. [10] and the formulas of Nielson and Jung [15] are conditionally stable. Errors creep into the computations in the neighbourhood of removable singularities where subtractive cancellations in finite precision arithmetic are brought to prominence. This leads to a *gradual* loss of significant digits (errors gradually build up) in the compu-

tations as we approach the points of removable singularities. In the absence of numerically stable formulas/computational techniques, the accuracy promise of analytical solvers is lost and what is worse, the loss of significance will go unnoticed or misdiagnosed. In the aforesaid papers [10, 15] we find expressions with removable singularities, e.g.  $[\exp(\lambda t) - 1]/\lambda$ , which are typical examples used to demonstrate [12, 6] loss of significance in finite precision arithmetic.

The conditional stability issue also extends to the analytical solution procedure presented by Diachin and Herzog [3]. Although matrix decompositions are robust/stable with respect to rounding errors, the recursive algorithm used to compute the exponential function of block upper triangular matrices breaks down in certain situations. Moreover Parlett [18, p. 199] warned that implementations in finite precision arithmetic could be expected to give inaccurate results in particular situations. The message is clear: irrespective of the choice of the solution procedure, issues related to numerical instability exist and they need to be addressed.

Further, once such instabilities are identified, it is often not trivial to localize the terms in these formulas that participate to obtain a finite limit at removable singularities. Identifying such terms is crucial to control numerical instabilities and bound the loss of significant digits. We discuss these issues here and present algebraically equivalent yet numerically stable formulas for the X-IVAS scheme in 2D and 3D.

**2. Preliminaries.** Here we describe briefly the convention used in the description of the flow. The independent variables in Lagrangian kinematics are  $(\chi, t)$ , where  $\chi$  represents a label to identify particles (material points) and  $t$  represents the time elapsed after labeling. The primary dependent variable is the fluid particle trajectory denoted as  $\mathbf{X}(\chi, t)$ . The initial particle positions denoted by  $\mathbf{X}^0 := \mathbf{X}(\chi, 0)$  are assumed to be given. A natural choice for the label  $\chi$  is the ordered triple  $\mathbf{X}^0$ . The Lagrangian velocity and acceleration, denoted as  $\dot{\mathbf{X}}(\chi, t)$  and  $\ddot{\mathbf{X}}(\chi, t)$ , respectively are defined as follows.

$$(2.1) \quad \dot{\mathbf{X}}(\chi, t) := \frac{d}{dt}\mathbf{X}(\chi, t), \quad \ddot{\mathbf{X}}(\chi, t) := \frac{d^2}{dt^2}\mathbf{X}(\chi, t)$$

On the other hand, the independent variables in Eulerian kinematics are  $(\mathbf{x}, t)$ . Here  $\mathbf{x}$  denotes the spatial coordinate. The primary dependent variable is the fluid velocity  $\mathbf{u}(\mathbf{x}, t)$ . The so-called fundamental principle of kinematics [19] states that the velocity  $\mathbf{u}(\mathbf{x}, t)$  and acceleration  $\mathbf{a}(\mathbf{x}, t)$  at a given time  $t$  and fixed position  $\mathbf{x}$  (Eulerian description) is equal to the velocity  $\dot{\mathbf{X}}(\chi, t)$  and acceleration  $\ddot{\mathbf{X}}(\chi, t)$  of a particle that is present at that position and at that instant (Lagrangian description). Thus,

$$(2.2) \quad \mathbf{u}(\mathbf{x}, t) = \left. \frac{d}{dt}\mathbf{X}(\chi, t) \right|_{\mathbf{X}(\chi, t)=\mathbf{x}}, \quad \mathbf{a}(\mathbf{x}, t) = \left. \frac{d^2}{dt^2}\mathbf{X}(\chi, t) \right|_{\mathbf{X}(\chi, t)=\mathbf{x}}$$

As a corollary we have the following exact but implicit equations of particle motion.

$$(2.3) \quad \frac{d}{dt}\mathbf{X}(\chi, t) = \mathbf{u}(\mathbf{X}(\chi, t), t), \quad \frac{d^2}{dt^2}\mathbf{X}(\chi, t) = \frac{d}{dt}\dot{\mathbf{X}}(\chi, t) = \mathbf{a}(\mathbf{X}(\chi, t), t)$$

### 3. The X-IVAS scheme.

**3.1. Introduction.** The explicit time integration of the particle position and velocity along the streamline results in the following equations of motion.

$$(3.1) \quad \frac{d}{dt}\mathbf{X}(\chi, t) = \mathbf{u}(\mathbf{X}(\chi, t), t^n), \quad \frac{d}{dt}\dot{\mathbf{X}}(\chi, t) = \mathbf{a}(\mathbf{X}(\chi, t), t^n)$$

Recall that the data corresponding to the dependent variables is stored with the particles which form a sufficiently large yet finite set. In other words, the data at any given time is available as discrete samples at the spatial locations occupied by the particles. Data interpolation is inevitable to have spatially continuous vector fields and to solve for the particle motion. Hence in the equations of motion  $\mathbf{u}(\mathbf{x}, t^n)$  and  $\mathbf{a}(\mathbf{x}, t^n)$ , which are unknown for an arbitrary  $\mathbf{x}$  are replaced by the interpolated counterparts  $\mathbf{u}^h(\mathbf{x}, t^n)$  and  $\mathbf{a}^h(\mathbf{x}, t^n)$ , respectively. The superscript  $h$  represents the discretization size associated to the interpolation. It follows that the trajectory obtained from these interpolated vector fields needs to be represented as  $\mathbf{X}^h(\chi, t)$ .

In the following we describe the X-IVAS scheme to integrate the equations of particle motion from time  $t^n$  to  $t^{n+1}$  as a four step process.

*Step 1: Projection.* This step involves the projection of vector fields stored with the particles onto a simplicial mesh. Consider a simplicial mesh over the problem domain and a set of characteristic domains corresponding to every node of the mesh. Let  $\mathcal{P}^i$  be an operator that projects data onto a mesh node with index  $i$  from a set of sample points in the corresponding characteristic domain. Using this projection operator we calculate the velocity  $\bar{\mathbf{u}}^i(t^n)$  and acceleration  $\bar{\mathbf{a}}^i(t^n)$  vector fields at the mesh nodes as follows.

$$(3.2) \quad \bar{\mathbf{u}}^i(t^n) := \mathcal{P}^i[\dot{\mathbf{X}}(\chi, t^n)], \quad \bar{\mathbf{a}}^i(t^n) := \mathcal{P}^i[\ddot{\mathbf{X}}(\chi, t^n)]$$

This projection step is unnecessary when  $t^n = 0$  where we can obtain  $\bar{\mathbf{u}}^i(0)$  and  $\bar{\mathbf{a}}^i(0)$  directly from the prescribed initial conditions.

*Step 2: Interpolation.* In this step we do a piecewise linear interpolation of vector fields projected onto the mesh nodes. Using the velocity  $\bar{\mathbf{u}}^i(t^n)$  and acceleration  $\bar{\mathbf{a}}^i(t^n)$  vector fields at the mesh nodes we construct a piecewise linear interpolation of these vector fields as follows.

$$(3.3) \quad \mathbf{u}^h(\mathbf{x}, t^n) := \mathbf{N}^i(\mathbf{x})\bar{\mathbf{u}}^i(t^n), \quad \mathbf{a}^h(\mathbf{x}, t^n) := \mathbf{N}^i(\mathbf{x})\bar{\mathbf{a}}^i(t^n)$$

In the above equation  $\mathbf{N}^i(\mathbf{x})$  represents the piecewise linear shape function corresponding to the node  $i$ . Let  $\mathbf{x}^j$  denote the spatial coordinate of node  $j$ ,  $\langle \cdot \rangle_j$  denote the average operator over the index  $j$  and  $\delta^{ij}$  denote the Kronecker delta. For a given simplex, we can express  $\mathbf{N}^i(\mathbf{x})$  in terms of its gradient  $\nabla \mathbf{N}^i$  (which is constant within the simplex) and the spatial coordinate  $\mathbf{x}$  as follows.

$$(3.4) \quad \mathbf{N}^i(\mathbf{x}) := \nabla \mathbf{N}^i \cdot (\mathbf{x} - \langle \mathbf{x}^j \rangle_j) + \frac{1}{\delta^{kk}}$$

Using the above equation  $\mathbf{u}^h(\mathbf{x}, t^n)$  and  $\mathbf{a}^h(\mathbf{x}, t^n)$  can be expressed within each simplex as follows.

$$(3.5) \quad \mathbf{u}^h(\mathbf{x}, t^n) = [\bar{\mathbf{u}}^i(t^n) \otimes \nabla \mathbf{N}^i] \cdot (\mathbf{x} - \langle \mathbf{x}^j \rangle_j) + \langle \bar{\mathbf{u}}^j(t^n) \rangle_j = \mathbf{A}^n \cdot \mathbf{x} + \mathbf{b}^n$$

$$(3.6) \quad \mathbf{a}^h(\mathbf{x}, t^n) = [\bar{\mathbf{a}}^i(t^n) \otimes \nabla \mathbf{N}^i] \cdot (\mathbf{x} - \langle \mathbf{x}^j \rangle_j) + \langle \bar{\mathbf{a}}^j(t^n) \rangle_j = \mathbf{C}^n \cdot \mathbf{x} + \mathbf{d}^n$$

Here  $\otimes$  denotes the tensor product. Further,  $\mathbf{A}^n, \mathbf{b}^n, \mathbf{C}^n$  and  $\mathbf{d}^n$  are constant tensors evaluated for each simplex at time  $t^n$  and are defined as follows.

$$(3.7) \quad \mathbf{A}^n := [\bar{\mathbf{u}}^i(t^n) \otimes \nabla \mathbf{N}^i], \quad \mathbf{b}^n := \langle \bar{\mathbf{u}}^i(t^n) \rangle_i - \mathbf{A}^n \cdot \langle \mathbf{x}^i \rangle_i$$

$$(3.8) \quad \mathbf{C}^n := [\bar{\mathbf{a}}^i(t^n) \otimes \nabla \mathbf{N}^i], \quad \mathbf{d}^n := \langle \bar{\mathbf{a}}^i(t^n) \rangle_i - \mathbf{C}^n \cdot \langle \mathbf{x}^i \rangle_i$$

*Step 3: Integration.* Here we describe the time integration of the approximate equations of particle motion. The approximate equations of motion for the particles in the X-IVAS scheme can be written as follows.

$$(3.9) \quad \frac{d}{dt} \mathbf{X}^h(\chi, t) = \mathbf{u}^h(\mathbf{X}^h(\chi, t), t^n), \quad \frac{d}{dt} \dot{\mathbf{X}}^h(\chi, t) = \mathbf{a}^h(\mathbf{X}^h(\chi, t), t^n)$$

Recall that the above equations are expressed in a piecewise manner as both  $\mathbf{u}^h$  and  $\mathbf{a}^h$  are defined in this manner. To be precise, within each simplex the particle motion is driven by the following equations which vary from one simplex to another.

$$(3.10) \quad \frac{d}{dt} \mathbf{X}^h(\chi, t) = \mathbf{A}^n \cdot \mathbf{X}^h(\chi, t) + \mathbf{b}^n, \quad \frac{d}{dt} \dot{\mathbf{X}}^h(\chi, t) = \mathbf{C}^n \cdot \mathbf{X}^h(\chi, t) + \mathbf{d}^n$$

Likewise, the time integration of these equations should also be done in a piecewise manner. In other words, if a particle tends to exit the current simplex prior to the end of the time step, its subsequent motion is driven by the equations written for the simplex in which it tends to enter and so forth until the end of the time step.

Further, certain relationships that existed between the dependent variables in the exact equations of motions no longer hold for the corresponding variables in the approximate equations of motion. That is,

$$(3.11) \quad \dot{\mathbf{X}}^h(\chi, t) \neq \frac{d}{dt} \mathbf{X}^h(\chi, t), \quad \ddot{\mathbf{X}}^h(\chi, t) := \frac{d}{dt} \dot{\mathbf{X}}^h(\chi, t) \neq \frac{d^2}{dt^2} \mathbf{X}^h(\chi, t),$$

Nevertheless, the X-IVAS scheme is consistent in the sense that these relations are recovered as the discretization size  $h \rightarrow 0$ . The analytical solution to the pair of equations given in Eq. (3.10) can be written within each simplex as follows.

$$(3.12) \quad \mathbf{X}^h(\chi, t) = e^{(t-t^n)\mathbf{A}^n} \cdot \mathbf{X}^h(\chi, t^n) + \left[ \int_{t^n}^t e^{(t-\tau)\mathbf{A}^n} d\tau \right] \cdot \mathbf{b}^n$$

$$(3.13) \quad \dot{\mathbf{X}}^h(\chi, t) = \dot{\mathbf{X}}^h(\chi, t^n) + \mathbf{C}^n \cdot \left[ \int_{t^n}^t \mathbf{X}^h(\chi, \tau) d\tau \right] + (t - t^n) \mathbf{d}^n$$

Note that the particle motion is restricted to the tangent curve of  $\mathbf{u}^h(\mathbf{x}, t^n)$  (i.e. the streamline) on which it was located at time  $t^n$  and is accelerated along this curve up to time  $t^{n+1}$ . Note that the solution for the particle velocity  $\dot{\mathbf{X}}^h(\chi, t)$  is given as a line integral along the curve  $\mathbf{X}^h(\chi, t)$ . This integral is left here *as is* for compactness and its evaluated form will be given in the following section.

Equations 3.12 and 3.13 justify the classification of the X-IVAS scheme as a *particle-based explicit exponential integrator*. Unlike classical exponential integrators [8, 17, 2] which integrate a global system of equations, the X-IVAS scheme integrates analytically a small and fixed-size local system of equations for the particles. This is an innovative approach that combines concepts of exponential integrators with the particle percept. On the one hand, it inherits the stability properties of exponential integrators which allow us to choose larger time steps. On the other hand, the associated algebra is computationally intensive, i.e. it is not memory bound. This conceptual setting is ideal for parallel computation which is an important strategy for faster simulations.

*Step 4: Update.* In this step we update the dependent variables at time  $t^{n+1}$  and repeat the process. At the end of the time step we obtain  $\mathbf{X}^h(\chi, t^{n+1})$  and  $\dot{\mathbf{X}}^h(\chi, t^{n+1})$  which are governed by the kinematics of the flow. The state of  $\ddot{\mathbf{X}}^h(\chi, t^{n+1})$  is governed by the dynamics of the internal and the external force terms that appear in the momentum balance equation of the flow.

**3.2. Remarks on the analytical solution.** In this section we simplify the analytical solution given in Eq. (3.12) and Eq. (3.13) and identify relationships among the terms that appear therein, if any. Consider three matrices  $\mathbf{P}$ ,  $\mathbf{Q}$  and  $\mathbf{R}$  which in turn are defined as functions of a given matrix  $\mathbf{A}$  and a scalar  $\tau$  as follows.

$$(3.14) \quad \mathbf{P}(\tau, \mathbf{A}) := e^{\tau \mathbf{A}}, \quad \mathbf{Q}(\tau, \mathbf{A}) := \int_0^\tau e^{\xi \mathbf{A}} d\xi, \quad \mathbf{R}(\tau, \mathbf{A}) := \int_0^\tau \int_0^\eta e^{\xi \mathbf{A}} d\xi d\eta$$

As it can be seen from the above equation, the considered matrices are exponential functions of the given matrix  $\mathbf{A}$ . The matrix  $\mathbf{P}$  is usually called the propagator [4]. The following relationships can be identified between the matrices  $\mathbf{P}$  and  $\mathbf{Q}$ .

$$(3.15) \quad \mathbf{Q}(\tau, \mathbf{A}) = \int_0^\tau \mathbf{P}(\xi, \mathbf{A}) d\xi = [e^{\tau \mathbf{A}} - \mathbf{I}] \cdot \text{inv}(\mathbf{A}) = [\mathbf{P}(\tau, \mathbf{A}) - \mathbf{I}] \cdot \text{inv}(\mathbf{A})$$

$$(3.16) \quad \Rightarrow \quad \mathbf{P}(\tau, \mathbf{A}) = \mathbf{Q}(\tau, \mathbf{A}) \cdot \mathbf{A} + \mathbf{I}$$

Likewise, the matrices  $\mathbf{Q}$  and  $\mathbf{R}$  satisfy the following relationships.

$$(3.17) \quad \mathbf{R}(\tau, \mathbf{A}) = \int_0^\tau \mathbf{Q}(\xi, \mathbf{A}) d\xi = \left[ (e^{\tau \mathbf{A}} - \mathbf{I}) \cdot \text{inv}(\mathbf{A}) - \tau \mathbf{I} \right] \cdot \text{inv}(\mathbf{A})$$

$$(3.18) \quad \mathbf{R}(\tau, \mathbf{A}) = [\mathbf{Q}(\tau, \mathbf{A}) - \tau \mathbf{I}] \cdot \text{inv}(\mathbf{A}), \quad \Rightarrow \quad \mathbf{Q}(\tau, \mathbf{A}) = \mathbf{R}(\tau, \mathbf{A}) \cdot \mathbf{A} + \tau \mathbf{I}$$

In the above equations  $\text{inv}(\mathbf{A})$  denotes the matrix inverse of  $\mathbf{A}$ . Further, the products involving  $\text{inv}(\mathbf{A})$  and  $\mathbf{A}$  in these equations are commutative, i.e. the order in which they appear are irrelevant. Using these definitions, we can express the analytical solution of the equations of motion in the X-IVAS scheme as follows.

$$(3.19) \quad \mathbf{X}^h(\chi, t) = \mathbf{P}(t - t^n, \mathbf{A}^n) \cdot \mathbf{X}^h(\chi, t^n) + \mathbf{Q}(t - t^n, \mathbf{A}^n) \cdot \mathbf{b}^n$$

$$(3.20) \quad \dot{\mathbf{X}}^h(\chi, t) = \dot{\mathbf{X}}^h(\chi, t^n) + (t - t^n) \mathbf{d}^n \\ + \mathbf{C}^n \cdot [\mathbf{Q}(t - t^n, \mathbf{A}^n) \cdot \mathbf{X}^h(\chi, t^n) + \mathbf{R}(t - t^n, \mathbf{A}^n) \cdot \mathbf{b}^n]$$

Recall that a nodal projection of the data carried by the particles onto the background mesh is done after every time step and the tensors  $\mathbf{A}^n$ ,  $\mathbf{b}^n$ ,  $\mathbf{C}^n$  and  $\mathbf{d}^n$  have to be recalculated for each element using the projected data. It follows that the matrices  $\mathbf{P}$ ,  $\mathbf{Q}$  and  $\mathbf{R}$  also need to be recalculated for each element after every time step.

Note that we need the exit points of the particles on the simplex boundary to perform the piecewise integration of particle motions described earlier in the paragraph following Eq. (3.10). To find the exit points we need to solve the intersection of their trajectories with the simplex boundary. A procedure to solve for the exit points using Newton linearisation was presented by Kipfer et al. [13]. In this procedure the matrices  $\mathbf{P}$  and  $\mathbf{Q}$  have to be evaluated at every iteration as the time increments need not be uniform. Should one decide to use *analytical* sub-stepping procedures to arrive at the exit point and if a constant *sub* time step is used for all the particles throughout the sub-stepping procedure, then we need to compute these matrices for each element just once. The incremental method [15] for computing tangent curves and the analytical time stepping algorithm called ANTS [3] are based on this idea.

## 4. Functions of matrices.

**4.1. Introduction.** Functions of matrices can be defined in various yet equivalent ways, a comprehensive presentation of which is made by Higham [7]. Let the size of  $\mathbf{A}$  be  $n \times n$

and assume that a given scalar function  $f(\lambda)$  takes well-defined values (including values associated with derivatives where appropriate) at the eigenvalues of  $\mathbf{A}$  denoted by the sequence  $Z := \{\lambda_1, \lambda_2, \dots, \lambda_n\}$ . We will use the Newton form of the polynomial interpolation definition for  $f(\mathbf{A})$ . The coefficients in the Newton form of the interpolating polynomial have the algebraic structure of divided differences.

We denote the  $k$ th divided difference of  $f(\lambda)$  on the subsequence  $Z_i^k := \{\lambda_i, \lambda_{i+1}, \dots, \lambda_{i+k}\}$  as  $f[\lambda_i; \lambda_{i+1}; \dots; \lambda_{i+k}]$  and define it using the following recurrence equations.

$$(4.1) \quad f[\lambda_i] := f(\lambda_i)$$

$$(4.2) \quad f[\lambda_i; \lambda_{i+1}; \dots; \lambda_{i+k}] := \frac{f[\lambda_{i+1}; \lambda_{i+2}; \dots; \lambda_{i+k}] - f[\lambda_i; \lambda_{i+1}; \dots; \lambda_{i+k-1}]}{\lambda_{i+k} - \lambda_i}$$

$$(4.3) \quad \lambda_i = \lambda_{i+1} = \dots = \lambda_{i+k} \Rightarrow f[\lambda_i; \lambda_{i+1}; \dots; \lambda_{i+k}] := \frac{1}{k!} \frac{\partial^k}{\partial \lambda^k} f(\lambda) \Big|_{\lambda=\lambda_i}$$

It is a well-known result that the value of  $f[\lambda_i; \lambda_{i+1}; \dots; \lambda_{i+k}]$  does not depend on the order of  $\lambda_i, \lambda_{i+1}, \dots, \lambda_{i+k}$  in  $Z_i^k$ . The Newton form of the polynomial interpolation definition for  $f(\mathbf{A})$  is

$$(4.4) \quad f(\mathbf{A}) = f[\lambda_1] \mathbf{I} + \sum_{k=1}^{n-1} f[\lambda_1; \lambda_2; \dots; \lambda_{1+k}] (\mathbf{A} - \lambda_1 \mathbf{I}) (\mathbf{A} - \lambda_2 \mathbf{I}) \cdots (\mathbf{A} - \lambda_k \mathbf{I})$$

The computation of *nearly confluent* divided differences are known to suffer from subtractive cancellations in floating point arithmetic. Thus, this definition gives *a priori* warning about the *gradual* loss of significance in the computation of  $f(\mathbf{A})$  in the neighbourhood of removable singularities. By grouping terms prone to loss of significance (as divided differences) it also paves way to systematically design procedures for the stable computation of  $f(\mathbf{A})$ . Moreover this definition is independent of the Jordan structure of  $\mathbf{A}$  [7, p. 6] which makes it convenient to implement in a computer program.

It is possible to write the  $k$ th divided difference  $f[\lambda_1; \lambda_2; \dots; \lambda_{1+k}]$  as follows.

$$(4.5) \quad f[\lambda_1; \lambda_2; \dots; \lambda_{1+k}] = \sum_{i=1}^{1+k} \frac{f(\lambda_i)}{\prod_{j \neq i} (\lambda_i - \lambda_j)}, \quad j \in \{1, 2, \dots, 1+k\}$$

Using the above identity we can transform the formulas given in this article to compute tangent curves to the ones presented by Nielson and Jung [15]. Algebraic rearrangements using this identity do not avoid existing issues related to loss of significance and makes matters worse by obscuring them. The example in §5.1 drives this point home.

To express the closed-form solution of the X-IVAS scheme we need to consider the cases where  $n = 2$  (2D) and  $n = 3$  (3D). For  $n = 3$ , we can express  $f(\mathbf{A})$  as follows.

$$(4.6) \quad f(\mathbf{A}) = f(\lambda_1) \mathbf{I} + f[\lambda_1; \lambda_2] (\mathbf{A} - \lambda_1 \mathbf{I}) + f[\lambda_1; \lambda_2; \lambda_3] (\mathbf{A} - \lambda_1 \mathbf{I}) (\mathbf{A} - \lambda_2 \mathbf{I})$$

Without loss of generality we assume that the eigenvalue  $\lambda_3$  is a real number and the eigenvalues  $\lambda_1$  and  $\lambda_2$  might be complex numbers. Complex eigenvalues will always occur in conjugate pairs, i.e.  $\{\lambda_1, \lambda_2\} = \{\lambda_c, \lambda_c^*\}$ . The subscript  $c$  indicates that it is a complex number and the superscript  $*$  indicates that it is a complex conjugate.

Although Eq. (4.6) holds for all eigenvalues, this form is convenient to implement in a computer program when the eigenvalues are real numbers. In the case of complex eigenvalues

Eq. (4.6) can be simplified to evaluate  $f(\mathbf{A})$  as follows.

$$(4.7) \quad f(\mathbf{A}) = \frac{\operatorname{Im}[f^*(\lambda_c)\lambda_c]}{\operatorname{Im}[\lambda_c]}\mathbf{I} + \frac{\operatorname{Im}[f(\lambda_c)]}{\operatorname{Im}[\lambda_c]}\mathbf{A} \\ + \left[ \frac{f(\lambda_3)\operatorname{Im}(\lambda_c) - \lambda_3\operatorname{Im}[f(\lambda_c)] - \operatorname{Im}[f^*(\lambda_c)\lambda_c]}{\operatorname{Im}[\lambda_c]} \right] \left[ \frac{\mathbf{A}^2 - 2\operatorname{Re}(\lambda_c)\mathbf{A} + |\lambda_c|^2\mathbf{I}}{\lambda_3^2 - 2\operatorname{Re}(\lambda_c)\lambda_3 + |\lambda_c|^2} \right]$$

Here, the functions  $\operatorname{Re}(\lambda_c)$  and  $\operatorname{Im}(\lambda_c)$  return the real and imaginary parts of a complex argument  $\lambda_c$ , respectively.

**4.2. Formulas for exponential functions of  $3 \times 3$  matrices.** In this section we consider the case when  $f(\lambda) := \exp(\tau\lambda)$  and write the expressions for the matrices  $\mathbf{P}(\tau, \mathbf{A})$ ,  $\mathbf{Q}(\tau, \mathbf{A})$  and  $\mathbf{R}(\tau, \mathbf{A})$  which were defined earlier in Eq. (3.14). It is straightforward to verify the following results.

$$(4.8) \quad \int_0^\tau e^{\xi\lambda} d\xi = \frac{e^{\tau\lambda} - 1}{\lambda} = \tau \exp[0; \tau\lambda]$$

$$(4.9) \quad \int_0^\tau \int_0^\eta e^{\xi\lambda} d\xi d\eta = \frac{e^{\tau\lambda} - 1 - \tau\lambda}{\lambda^2} = \tau^2 \exp[0; 0; \tau\lambda]$$

Following this line we define two auxiliary functions  $q(x)$  and  $r(x)$  which are divided differences of the exponential function.

$$(4.10) \quad q(x) := \exp[0; x] = \begin{cases} \frac{e^x - 1}{x} & \text{if } x \neq 0, \\ 1 & \text{if } x = 0. \end{cases}$$

$$(4.11) \quad r(x) := \exp[0; 0; x] = q[0; x] = \begin{cases} \frac{e^x - 1 - x}{x^2} & \text{if } x \neq 0, \\ \frac{1}{2} & \text{if } x = 0. \end{cases}$$

Using these auxiliary functions we can express  $\mathbf{P}(\tau, \mathbf{A})$ ,  $\mathbf{Q}(\tau, \mathbf{A})$  and  $\mathbf{R}(\tau, \mathbf{A})$  as follows.

(4.12)

$$\mathbf{P}(\tau, \mathbf{A}) = e^{\tau\lambda_1}\mathbf{I} + \tau \exp[\tau\lambda_1; \tau\lambda_2](\mathbf{A} - \lambda_1\mathbf{I}) + \tau^2 \exp[\tau\lambda_1; \tau\lambda_2; \tau\lambda_3](\mathbf{A} - \lambda_1\mathbf{I})(\mathbf{A} - \lambda_2\mathbf{I})$$

(4.13)

$$\mathbf{Q}(\tau, \mathbf{A}) = \tau q(\tau\lambda_1)\mathbf{I} + \tau^2 q[\tau\lambda_1; \tau\lambda_2](\mathbf{A} - \lambda_1\mathbf{I}) + \tau^3 q[\tau\lambda_1; \tau\lambda_2; \tau\lambda_3](\mathbf{A} - \lambda_1\mathbf{I})(\mathbf{A} - \lambda_2\mathbf{I})$$

(4.14)

$$\mathbf{R}(\tau, \mathbf{A}) = \tau^2 r(\tau\lambda_1)\mathbf{I} + \tau^3 r[\tau\lambda_1; \tau\lambda_2](\mathbf{A} - \lambda_1\mathbf{I}) + \tau^4 r[\tau\lambda_1; \tau\lambda_2; \tau\lambda_3](\mathbf{A} - \lambda_1\mathbf{I})(\mathbf{A} - \lambda_2\mathbf{I})$$

The expression for  $\mathbf{P}(\tau, \mathbf{A})$  is a trivial specialization of Eq. (4.6) for  $f(\lambda) := \exp(\tau\lambda)$ . We obtain  $\mathbf{Q}(\tau, \mathbf{A})$  by integrating the terms in  $\mathbf{P}(\xi, \mathbf{A})$  with respect to  $\xi$ ; cf. Eq. (3.15). Likewise,  $\mathbf{R}(\tau, \mathbf{A})$  is obtained by integrating the terms in  $\mathbf{Q}(\xi, \mathbf{A})$  with respect to  $\xi$ ; cf. Eq. (3.17). We



have used the following results to arrive at these equations.

$$(4.15) \quad \int_0^\tau \xi \exp[\xi \lambda_1; \xi \lambda_2] d\xi = \tau \frac{\exp[0; \tau \lambda_2] - \exp[0; \tau \lambda_1]}{\lambda_2 - \lambda_1} = \tau^2 q[\tau \lambda_1; \tau \lambda_2]$$

$$(4.16) \quad \int_0^\tau \xi^2 \exp[\xi \lambda_1; \xi \lambda_2; \xi \lambda_3] d\xi = \int_0^\tau \xi \frac{\exp[\xi \lambda_2; \xi \lambda_3] - \exp[\xi \lambda_1; \xi \lambda_2]}{\lambda_3 - \lambda_1} d\xi \\ = \tau^3 q[\tau \lambda_1; \tau \lambda_2; \tau \lambda_3]$$

$$(4.17) \quad \int_0^\tau \int_0^\eta \xi \exp[\xi \lambda_1; \xi \lambda_2] d\xi d\eta = \tau^2 \frac{\exp[0; 0; \tau \lambda_2] - \exp[0; 0; \tau \lambda_1]}{\lambda_2 - \lambda_1} = \tau^3 r[\tau \lambda_1; \tau \lambda_2]$$

$$(4.18) \quad \int_0^\tau \int_0^\eta \xi^2 \exp[\xi \lambda_1; \xi \lambda_2; \xi \lambda_3] d\xi d\eta = \int_0^\tau \int_0^\eta \xi \frac{\exp[\xi \lambda_2; \xi \lambda_3] - \exp[\xi \lambda_1; \xi \lambda_2]}{\lambda_3 - \lambda_1} d\xi d\eta \\ = \tau^4 r[\tau \lambda_1; \tau \lambda_2; \tau \lambda_3]$$

Let  $\alpha, \beta$  be real numbers and consider a complex number  $\lambda_c$  as defined below.

$$(4.19) \quad i := \sqrt{-1}, \quad \lambda_c := \alpha + i\beta, \quad \Rightarrow \lambda_c^* = \alpha - i\beta$$

The cardinal sine function  $\text{sinc}(x)$  is defined as follows.

$$(4.20) \quad \text{sinc}(x) := \begin{cases} \frac{\sin(x)}{x} & \text{if } x \neq 0, \\ 1 & \text{if } x = 0. \end{cases}$$

Further, define two auxiliary functions  $\Psi(x, y)$  and  $\Phi(x, y)$  as follows.

$$(4.21) \quad \Psi(x, y) := \cos(y) - x \text{sinc}(y)$$

$$(4.22) \quad \Phi(x, y) := \exp[-iy; iy; x] = \begin{cases} \frac{e^x - \Psi(-x, y)}{x^2 + y^2} & \text{if } (x, y) \neq (0, 0), \\ \frac{1}{2} & \text{if } (x, y) = (0, 0). \end{cases}$$

In the case of complex eigenvalues, i.e.  $\{\lambda_1, \lambda_2\} = \{\lambda_c, \lambda_c^*\}$  we can evaluate  $\mathbf{P}(\tau, \mathbf{A})$ ,  $\mathbf{Q}(\tau, \mathbf{A})$  and  $\mathbf{R}(\tau, \mathbf{A})$  as follows.

$$(4.23) \quad \mathbf{P}(\tau, \mathbf{A}) = e^{\tau\alpha} [\Psi(\tau\alpha, \tau\beta)\mathbf{I} + \tau \text{sinc}(\tau\beta)\mathbf{A} + \tau^2 \Phi(\tau\lambda_3 - \tau\alpha, \tau\beta)[(\mathbf{A} - \alpha\mathbf{I})^2 + \beta^2\mathbf{I}]$$

$$(4.24) \quad \mathbf{Q}(\tau, \mathbf{A}) = \tau e^{\tau\alpha} [\text{sinc}(\tau\beta)\mathbf{I} + \tau \Phi(-\tau\alpha, \tau\beta)(\mathbf{A} - 2\alpha\mathbf{I}) \\ + \tau^2 \Phi(\star, \tau\beta)[- \tau\alpha; \tau\lambda_3 - \tau\alpha][(\mathbf{A} - \alpha\mathbf{I})^2 + \beta^2\mathbf{I}]$$

$$(4.25) \quad \mathbf{R}(\tau, \mathbf{A}) = \tau^2 e^{\tau\alpha} [\Phi(-\tau\alpha, \tau\beta)\mathbf{I} + \tau \Phi(\star, \tau\beta)[- \tau\alpha; -\tau\alpha](\mathbf{A} - 2\alpha\mathbf{I}) \\ + \tau^2 \Phi(\star, \tau\beta)[- \tau\alpha; -\tau\alpha; \tau\lambda_3 - \tau\alpha][(\mathbf{A} - \alpha\mathbf{I})^2 + \beta^2\mathbf{I}]$$

The expression for  $\mathbf{P}(\tau, \mathbf{A})$  is a trivial specialization of Eq. (4.7) for  $f(\lambda) := \exp(\tau\lambda)$  and the choice of the auxiliary functions  $\Psi(x, y)$  and  $\Phi(x, y)$  is motivated by the structure of the same. The notation  $\Phi(\star, y)_{[x_1; x_2]}$  means that the divided differences are to be taken with respect to the variable in whose place the symbol  $\star$  appears. The rest of this section describes some results which were used to arrive at the expressions for  $\mathbf{Q}(\tau, \mathbf{A})$  and  $\mathbf{R}(\tau, \mathbf{A})$  from the expression for  $\mathbf{P}(\tau, \mathbf{A})$ .

The following integrals are straightforward.

$$(4.26) \quad \int_0^\tau e^{\xi\alpha} \cos(\xi\beta) d\xi = \frac{e^{\tau\alpha}[\alpha \cos(\tau\beta) + \beta \sin(\tau\beta)] - \alpha}{\alpha^2 + \beta^2} = \tau e^{\tau\alpha} [\text{sinc}(\tau\beta) - \tau\alpha \Phi(-\tau\alpha, \tau\beta)]$$

$$(4.27) \quad \int_0^\tau e^{\xi\alpha} \sin(\xi\beta) d\xi = \frac{e^{\tau\alpha}[\alpha \sin(\tau\beta) - \beta \cos(\tau\beta)] + \beta}{\alpha^2 + \beta^2} = \tau^2 \beta e^{\tau\beta} \Phi(-\tau\alpha, \tau\beta)$$

$$(4.28) \quad \Rightarrow \int_0^\tau e^{\xi\alpha} \Psi(\xi\alpha, \xi\beta) d\xi = \tau e^{\tau\alpha} [\text{sinc}(\tau\beta) - 2\tau\alpha \Phi(-\tau\alpha, \tau\beta)]$$

Using Eqs. (4.27) and (4.28) we obtain the following result.

$$(4.29) \quad \int_0^\tau e^{\xi\alpha} \xi^2 \Phi(\xi\lambda_3 - \xi\alpha, \xi\beta) d\xi = \int_0^\tau \frac{e^{\xi\lambda_3} - e^{\xi\alpha} [\Psi(\xi\alpha, \xi\beta) + \xi\lambda_3 \text{sinc}(\xi\beta)]}{(\lambda_3 - \alpha)^2 + \beta^2} d\xi$$

$$(4.30) \quad = \tau \left[ \frac{q(\tau\lambda_3) - e^{\tau\alpha} [\text{sinc}(\tau\beta) + \tau(\lambda_3 - 2\alpha) \Phi(-\tau\alpha, \tau\beta)]}{(\lambda_3 - \alpha)^2 + \beta^2} \right]$$

$$(4.31) \quad = \tau e^{\tau\alpha} \left[ \frac{[e^{\tau\lambda_3} - 1]e^{-\tau\alpha} - \tau\lambda_3 [\text{sinc}(\tau\beta) + \tau(\lambda_3 - 2\alpha) \Phi(-\tau\alpha, \tau\beta)]}{\tau\lambda_3 [(\lambda_3 - \alpha)^2 + \beta^2]} \right]$$

$$(4.32) \quad = \tau e^{\tau\alpha} \left[ \frac{e^{\tau\lambda_3 - \tau\alpha} - \Psi(\tau\alpha - \tau\lambda_3, \tau\beta) - \tau^2 [(\lambda_3 - \alpha)^2 + \beta^2] \Phi(-\tau\alpha, \tau\beta)}{\tau\lambda_3 [(\lambda_3 - \alpha)^2 + \beta^2]} \right]$$

$$(4.33) \quad = \tau^3 e^{\tau\alpha} \left[ \frac{\Phi(\tau\lambda_3 - \tau\alpha, \tau\beta) - \Phi(-\tau\alpha, \tau\beta)}{(\tau\lambda_3 - \tau\alpha) - (-\tau\alpha)} \right] = \tau^3 e^{\tau\alpha} \Phi(\star, \tau\beta)[- \tau\alpha, \tau\lambda_3 - \tau\alpha]$$

The results given in Eqs. (4.27), (4.28) and (4.33) are used to obtain  $\mathbf{Q}(\tau, \mathbf{A})$  from  $\mathbf{P}(\tau, \mathbf{A})$ . Substituting  $\lambda_3 = 0$  in Eq. (4.33) we get the following result.

$$(4.34) \quad \int_0^\tau e^{\xi\alpha} \xi^2 \Phi(-\xi\alpha, \xi\beta) d\xi = \tau^3 e^{\tau\alpha} \Phi(\star, \tau\beta)[- \tau\alpha; -\tau\alpha]$$

Using Eqs. (4.33) and (4.34) we arrive at the following result.

$$(4.35)$$

$$\int_0^\tau e^{\xi\alpha} \xi^3 \Phi(\star, \xi\beta)[- \xi\alpha; \xi\lambda_3 - \xi\alpha] d\xi = \int_0^\tau e^{\xi\alpha} \xi^2 \frac{\Phi(\xi\lambda_3 - \xi\alpha, \xi\beta) - \Phi(-\xi\alpha, \xi\beta)}{\lambda_3} d\xi$$

$$(4.36) \quad = \tau^3 e^{\tau\alpha} \left[ \frac{\Phi(\star, \tau\beta)[- \tau\alpha; \tau\lambda_3 - \tau\alpha] - \Phi(\star, \tau\beta)[- \tau\alpha; -\tau\alpha]}{\lambda_3} \right]$$

$$(4.37) \quad = \tau^4 e^{\tau\alpha} \Phi(\star, \tau\beta)[- \tau\alpha; -\tau\alpha; \tau\lambda_3 - \tau\alpha]$$

The results given in Eqs. (4.27), (4.34) and (4.37) are used to obtain  $\mathbf{R}(\tau, \mathbf{A})$  from  $\mathbf{Q}(\tau, \mathbf{A})$ .

**REMARK:** Note that the equations for  $\mathbf{P}(\tau, \mathbf{A})$ ,  $\mathbf{Q}(\tau, \mathbf{A})$  and  $\mathbf{R}(\tau, \mathbf{A})$  are expressed (3D problems where  $n = 3$ ) as the sum of three terms. Due to the properties of the polynomial in the Newton's form, the corresponding equations for  $n = 2$  (2D problems) can be obtained from the equations for  $n = 3$  by dropping out the third term.

**4.3. Formulas for the eigenvalues of  $3 \times 3$  matrices.** Let  $\det(\mathbf{A})$  and  $\text{tr}(\mathbf{A})$  denote the determinant and trace of the matrix  $\mathbf{A}$ , respectively. When  $n = 3$ , the characteristic equation of the matrix  $\mathbf{A}$  is given by the following.

$$(4.38) \quad \det(\mathbf{A} - \lambda \mathbf{I}) = 0$$

$$(4.39) \quad \Rightarrow \lambda^3 - \text{tr}(\mathbf{A})\lambda^2 + \frac{\text{tr}(\mathbf{A})^2 - \text{tr}(\mathbf{A}^2)}{2}\lambda - \det(\mathbf{A}) = 0$$

The solution of the above cubic equation can be found by Cardano's method (see [20]). The calculation steps of the same are summarized below.

$$(4.40) \quad \mathbf{B} := \mathbf{A} - \frac{\text{tr}(\mathbf{A})}{3} \mathbf{I}, \quad Q := \frac{\text{tr}(\mathbf{B}^2)}{6}, \quad R := \frac{\det(\mathbf{B})}{2}$$

$$(4.41) \quad \lambda_1 = \frac{\text{tr}(\mathbf{A})}{3} + \sqrt[3]{R - \sqrt{R^2 - Q^3}} e^{-i(2\pi/3)} + \sqrt[3]{R + \sqrt{R^2 - Q^3}} e^{i(2\pi/3)}$$

$$(4.42) \quad \lambda_2 = \frac{\text{tr}(\mathbf{A})}{3} + \sqrt[3]{R - \sqrt{R^2 - Q^3}} e^{-i(4\pi/3)} + \sqrt[3]{R + \sqrt{R^2 - Q^3}} e^{i(4\pi/3)}$$

$$(4.43) \quad \lambda_3 = \frac{\text{tr}(\mathbf{A})}{3} + \sqrt[3]{R - \sqrt{R^2 - Q^3}} + \sqrt[3]{R + \sqrt{R^2 - Q^3}}$$

We follow the convention that the cube roots that appear in the above expressions are real and single valued. The three admissible solutions to the cube root function are already taken into consideration in the above formula.

Note that when the discriminant  $(R^2 - Q^3) > 0$ , we obtain complex eigenvalues. In this case, the formulas are already in a suitable format for implementation. When  $(R^2 - Q^3) \leq 0$  we obtain real eigenvalues and the formulas for the same can be written in a form better suited for implementation as follows.

$$(4.44) \quad \theta := \arccos\left(\frac{R}{\sqrt{Q^3}}\right), \quad \lambda_n = \frac{\text{tr}(\mathbf{A})}{3} + 2\sqrt{Q} \cos\left(\frac{2\pi n + \theta}{3}\right)$$

where  $\arccos()$  denotes the inverse cosine function whose range is defined to be the closed interval  $[0, \pi]$ . The formula for the real eigenvalues given in Eq. (4.44) guarantees  $\lambda_1 \leq \lambda_2 \leq \lambda_3$ . This can be verified using the following results.

$$(4.45) \quad \begin{aligned} & -1 \leq \cos\left(\frac{2\pi + \theta}{3}\right) \leq \frac{-1}{2} \quad \frac{\text{tr}(\mathbf{A})}{3} - 2\sqrt{Q} \leq \lambda_1 \leq \frac{\text{tr}(\mathbf{A})}{3} - \sqrt{Q} \\ 0 \leq \theta \leq \pi \Rightarrow & \frac{-1}{2} \leq \cos\left(\frac{4\pi + \theta}{3}\right) \leq \frac{1}{2} \quad \Rightarrow \quad \frac{\text{tr}(\mathbf{A})}{3} - \sqrt{Q} \leq \lambda_2 \leq \frac{\text{tr}(\mathbf{A})}{3} + \sqrt{Q} \\ & \frac{1}{2} \leq \cos\left(\frac{6\pi + \theta}{3}\right) \leq 1 \quad \frac{\text{tr}(\mathbf{A})}{3} + \sqrt{Q} \leq \lambda_3 \leq \frac{\text{tr}(\mathbf{A})}{3} + 2\sqrt{Q} \end{aligned}$$

Note that in the case of two equal eigenvalues, it will be either  $\lambda_1 = \lambda_2$  or  $\lambda_2 = \lambda_3$ . In all the situations the eigenvalue  $\lambda_3$  is always a real number.

## 5. Stable computation of formulas using finite precision.

**5.1. Introduction.** The issue with stable computation of formulas is best explained by an example. The example consists in the naïve computation of the second-order divided difference  $\exp[1; 1 + \varepsilon; 1 + 2\varepsilon]$ . We denote by Formula1 the *as is* expression of the second-order divided difference.

$$(5.1) \quad x_1 = 1, \quad x_2 = 1 + \varepsilon, \quad x_3 = 1 + 2\varepsilon$$

$$(5.2) \quad \exp[x_1; x_2; x_3] = \frac{1}{x_3 - x_1} \left[ \frac{e^{x_3} - e^{x_2}}{x_3 - x_2} - \frac{e^{x_2} - e^{x_1}}{x_2 - x_1} \right]$$

Using Eq. (4.5) the above equation can be rearranged in an algebraically equivalent form which we denote as Formula2.

$$(5.3) \quad \exp[x_1; x_2; x_3] = \frac{e^{x_1}}{(x_1 - x_2)(x_1 - x_3)} + \frac{e^{x_2}}{(x_2 - x_1)(x_2 - x_3)} + \frac{e^{x_3}}{(x_3 - x_1)(x_3 - x_2)}$$

TABLE 1  
Loss of significant digits in the naïve computations of  $\exp[1; 1 + \varepsilon; 1 + 2\varepsilon]$ .

$\varepsilon$	Formula1 computation	Exact 16 digits	Formula2 computation
$10^{-01}$	1.503 335 165 136 320	1.503 335 165 136 325	1.503 335 165 136 292
$10^{-02}$	1.372 811 947 550 877	1.372 811 947 550 820	1.372 811 947 550 871
$10^{-03}$	1.360 500 848 424 467	1.360 500 848 315 854	1.360 500 848 386 436
$10^{-04}$	1.359 276 824 430 971	1.359 276 836 249 607	1.359 276 831 150 054
$10^{-05}$	1.359 152 790 283 402	1.359 154 505 717 948	1.359 151 840 209 960
$10^{-06}$	1.359 135 026 857 928	1.359 142 273 371 229	1.359 375
$10^{-07}$	1.359 132 267 628 772 320	1.359 141 050 143 621	1.359 375
$10^{-08}$	2.220 446 084 949 470	1.359 140 927 820 931	4
$10^{-09}$	0	1.359 140 915 588 663	256
$10^{-10}$	0	1.359 140 914 365 436	0
$10^{-11}$	-2 220 445.681 810 107	1.359 140 914 243 114	-4 194 304
$10^{-12}$	-222 005 130.399 6447	1.359 140 914 230 881	-268 435 456
$10^{-13}$	0	1.359 140 914 229 658	17 179 869 184
$10^{-14}$	2 223 999 815 985.422	1.359 140 914 229 536	4 398 046 511 104
$10^{-15}$	0	1.359 140 914 229 523	0

All the expressions that appear in the formulas given by Nielson and Jung [15] are expressed in the above simplified form.

Table 1 illustrates the results of naïve computations of both formulas using double precision floating point arithmetic as  $\varepsilon \rightarrow 0$ . The exact values up to 16 digits of precision are given in the third column. The significant digits in both formula computations that coincide with the exact values are highlighted in green colour. We observe a gradual loss of significant digits in both formula computations which deteriorates as  $\varepsilon \rightarrow 0$ . For  $\varepsilon \leq 10^{-8}$  we lose all the significant digits in both formula computations.

This example demonstrates two features: 1) algebraic rearrangements using Eq. (4.5) does not avoid loss of significance in Formula2 and 2) the bad fame of nearly confluent divided differences is a blessing in disguise as it gives a priori warning about loss of significance in Formula1. In other words, due to the algebraic structure of Formula2 the loss of significance in the computations might go unnoticed or misdiagnosed.

In the analytical solution of the X-IVAS scheme, the following expressions might suffer from cancellation errors in a straight-forward (naïve) computation of the same using finite precision arithmetic.

$$(5.4) \quad \exp[\tau\lambda_1; \tau\lambda_2], \quad \exp[\tau\lambda_1; \tau\lambda_2; \tau\lambda_3], \quad q(\tau\lambda_1), \quad q[\tau\lambda_1; \tau\lambda_2], \quad q[\tau\lambda_1; \tau\lambda_2; \tau\lambda_3], \\ r(\tau\lambda_1), \quad r[\tau\lambda_1; \tau\lambda_2], \quad r[\tau\lambda_1; \tau\lambda_2; \tau\lambda_3], \quad \Phi(-\tau\alpha, \tau\beta), \quad \Phi(\tau\lambda_3 - \tau\alpha, \tau\beta), \\ \Phi(\star, \tau\beta)[- \tau\alpha; - \tau\alpha], \quad \Phi(\star, \tau\beta)[- \tau\alpha; \tau\lambda_3 - \tau\alpha], \quad \Phi(\star, \tau\beta)[- \tau\alpha; - \tau\alpha; \tau\lambda_3 - \tau\alpha]$$

The above expressions can be identified as the elements of the following nested set of divided differences.

$$(5.5) \quad \left\{ \left\{ \exp[x_1; x_2], q(x) \right\}, \left\{ \exp[x_1; x_2; x_3], q[x_1; x_2], r(x), \Phi(x, y) \right\}, \right. \\ \left. \left\{ q[x_1; x_2; x_3], r[x_1; x_2], \Phi(\star, y)[x_1; x_2] \right\}, \left\{ r[x_1; x_2; x_3], \Phi(\star, y)[x_1; x_1; x_2] \right\} \right\}$$

The order of the divided differences gradually increase from first-order in the first subset to fourth-order in the last subset. All elements of a subset are particular cases of the first element of that subset. For instance,

$$(5.6) \quad q[x_1;x_2] = \exp[0;x_1;x_2], \quad r(x) = \exp[0;0;x], \quad \Phi(x,y) = \exp[-iy;iy;x]$$

Following this line, it is possible to express all the divided differences in Eq. (5.5) as the divided differences of the exponential function; The details of the same are given in §5.5. As  $\varepsilon \rightarrow 0$ , the rate of loss of significant digits in a naïve computation of divided differences is generally equal to the order of the same. In the considered example, i.e.  $\exp[x_1;x_2;x_3]$  we loose significant digits at a second order rate. Following this line, naïve computation of the third and the fourth subsets in Eq. (5.5) are meaningless for  $\varepsilon \leq 10^{-5}$  and  $\varepsilon \leq 10^{-4}$ , respectively.

An algorithm for the accurate computation of divided differences of the exponential function was presented by McCurdy et.al [14]. Following this line, a similar algorithm for the accurate computation of divided differences of the auxiliary functions  $q()$  and  $r()$ , cf. Eq. (4.10) and Eq. (4.11), was presented by Caliri [1]. These algorithms have a wider scope, i.e. they were designed to evaluate functions of  $n \times n$  matrices appearing in exponential integrators for large systems of equations (ordinary or differential). A user who already has these algorithms implemented, might just invoke them to evaluate the divided differences listed in Eq. (5.4) and use them in the formulas for  $\mathbf{P}(\tau, \mathbf{A})$ ,  $\mathbf{Q}(\tau, \mathbf{A})$  and  $\mathbf{R}(\tau, \mathbf{A})$  given in §4.2. This would address the numerical stability issues in the formula computations.

In what follows, we present a simple yet stable piecewise definitions for divided differences. The methodology used to arrive at these piecewise definitions is of limited scope, i.e. this approach is not suitable for arbitrary  $k$ th order divided differences. Nevertheless, it is well suited for the at most fourth order divided differences found in the formulas for  $\mathbf{P}(\tau, \mathbf{A})$ ,  $\mathbf{Q}(\tau, \mathbf{A})$  and  $\mathbf{R}(\tau, \mathbf{A})$ .

**5.2. Optimal series approximation of divided differences.** In this section we establish optimal series approximation of divided differences of a given function  $f(x)$ . Consider the sequence  $\{x_1, x_2, \dots, x_n\}$  and some definitions related to this sequence.

$$(5.7) \quad x_a := \frac{1}{n} \sum_{i=1}^n x_i, \quad \tilde{x}_i := x_i - x_a, \quad \mathcal{X} := \{\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n\}$$

$$(5.8) \quad \mathcal{X}_p := \text{choose}(\mathcal{X}, 2), \quad \binom{n}{k} := \frac{n!}{k!(n-k)!}, \quad x_p^2 := \sum_{i=1}^{\binom{n}{2}} \prod_{j=1}^2 \mathcal{X}_p(i, j)$$

where  $x_a$  is the mean value of the sequence and  $\tilde{x}_i$  is the fluctuation of  $x_i$  about the mean. The function  $\text{choose}(\mathcal{X}, 2)$  returns a sequence  $\mathcal{X}_p$  consisting of pair-combinations (2-combinations) of elements from  $\mathcal{X}$ . In  $\mathcal{X}_p(i, j)$  the index  $i$  points to a combination and index  $j$  points to an element within this combination. The sum of the product of the pairs in  $\mathcal{X}_p$  is stored as the square of the auxiliary variable  $x_p$ . The result is stored as  $x_p^2$  to highlight the fact that it is a second order term. Likewise, the triple, quadruple and quintuple combinations of  $\mathcal{X}$  are denoted as  $\mathcal{X}_t$ ,  $\mathcal{X}_q$  and  $\mathcal{X}_v$ , respectively. Further, the sum of the product of the triples, quadruples and quintuples are stored in  $x_t^3$ ,  $x_q^4$  and  $x_v^5$ , respectively. Following  $x_p^2$ , the superscripts (which are ordinary powers) in  $x_t^3$ ,  $x_q^4$  and  $x_v^5$  highlight the fact that they are third,

fourth and fifth order terms, respectively. Thus,

$$(5.9) \quad \mathcal{X}_t := \text{choose}(\mathcal{X}, 3), \quad \mathcal{X}_q := \text{choose}(\mathcal{X}, 4), \quad \mathcal{X}_v := \text{choose}(\mathcal{X}, 5)$$

$$(5.10) \quad x_t^3 := \sum_{i=1}^{\binom{n}{3}} \prod_{j=1}^3 \mathcal{X}_t(i, j), \quad x_q^4 := \sum_{i=1}^{\binom{n}{4}} \prod_{j=1}^4 \mathcal{X}_q(i, j), \quad x_v^5 := \sum_{i=1}^{\binom{n}{5}} \prod_{j=1}^5 \mathcal{X}_v(i, j)$$

Using the above definitions, we can derive<sup>1</sup> the following identity for the divided differences of  $f(x)$ . The mean value theorem guarantees the existence of a  $\xi$  in the smallest interval containing  $\{x_1, x_2, \dots, x_n\}$  such that,

$$(5.11) \quad f^{(n)}(\xi) := \left. \frac{\partial^n}{\partial \lambda^n} f(\lambda) \right|_{\lambda=\xi}, \quad f(x_n) = f(\xi) + \sum_{n=1}^{\infty} (x_n - \xi)^n \frac{f^{(n)}(\xi)}{n!}$$

$$(5.12) \quad f[x_1; x_2; \dots; x_n] = \frac{f^{(n-1)}(x_a)}{(n-1)!} - x_p^2 \frac{f^{(n+1)}(x_a)}{(n+1)!} + x_t^3 \frac{f^{(n+2)}(x_a)}{(n+2)!} + (x_p^4 - x_q^4) \frac{f^{(n+3)}(x_a)}{(n+3)!} \\ + (x_v^5 - 2x_p^2 x_t^3) \frac{f^{(n+4)}(\xi)}{(n+4)!}$$

Note that the first term in the above series expansion provides a second-order approximation to  $f[x_1; x_2; \dots; x_n]$ . If the series is expanded with respect to any point other than  $x_a$ , the first-order terms are resurrected. Thus, the approximation is optimal for the choice  $x_a$ . For the first-order divided difference  $f[x_1; x_2]$ , the above equation can be simplified and easily extended to any number of terms as shown below.

$$(5.13) \quad h := \frac{x_2 - x_1}{2}, \quad x_p^2 = -h^2, \quad x_t^2 = 0, \quad x_q^2 = 0, \quad x_v^2 = 0$$

$$(5.14) \quad f[x_1; x_2] = f^{(1)}(x_a) + h^2 \frac{f^{(3)}(x_a)}{3!} + \dots + h^{2n-2} \frac{f^{(2n-1)}(x_a)}{(2n-1)!} + h^{2n} \frac{f^{(2n+1)}(\xi)}{(2n+1)!}$$

Likewise, for the second-order divided difference  $f[x_1; x_2; x_3]$ , Eq. (5.12) can be simplified to the following.

$$(5.15) \quad x_p^2 = -\frac{3}{2}x_\sigma^2, \quad x_\sigma^2 := \frac{\tilde{x}_1^2 + \tilde{x}_2^2 + \tilde{x}_3^2}{3}, \quad x_t^3 = \tilde{x}_1 \tilde{x}_2 \tilde{x}_3, \quad x_q^2 = 0, \quad x_v^2 = 0$$

$$(5.16) \quad f[x_1; x_2; x_3] = \frac{1}{2} f^{(2)}(x_a) + \frac{x_\sigma^2}{16} f^{(4)}(x_a) + \frac{x_t^3}{120} f^{(5)}(x_a) + \frac{x_\sigma^4}{320} f^{(6)}(x_a) + 3x_\sigma^2 x_t^3 \frac{f^{(7)}(\xi)}{7!}$$

where  $x_\sigma$  is the standard deviation of the considered sequence. It is possible to relate  $x_p$  and  $x_\sigma$  for all  $n$  and in this work we exploit this relationship as it reduces the number of arithmetic operations.

$$(5.17) \quad x_p^2 = -\frac{n}{2}x_\sigma^2$$

**5.3. Double precision floating point numbers.** We briefly describe how double precision floating point numbers are stored in a computer as per the IEEE 754 standard. Any

<sup>1</sup>As the algebra involved is overwhelming and error-prone, we have used the computer algebra system `Maple` to perform the simplifications and verifications. Thus, human intervention is dedicated to identify patterns and to discover abstract expressions such as  $x_p$ ,  $x_t$ ,  $x_q$ , etc.

decimal floating point number within the range of the double can be written in the normalized form as follows.

$$(5.18) \quad \text{Decimal form} \rightarrow (-1)^s 2^e 1.f \approx \underbrace{\boxed{s}}_1 \underbrace{\boxed{(e+1023)_b}}_{11} \underbrace{\boxed{0.f_b}}_{52} \leftarrow \text{Binary form}$$

$\leftarrow \text{No. of bits stored}$

In the above equation the boolean  $s \in \{0, 1\}$  is called the *sign bit*, the integer  $e$  is called the *exponent*;  $-1022 \leq e \leq 1023$  and the fraction  $f$  is called the *significand*. The numbers with a subscript b are expressed in the binary format. When the binary expression  $0.f_b$  is not exactly representable using 52 bits, it is rounded to the nearest representable number.

**5.4. Piecewise definition of divided differences.** To control (bound) the loss of significant digits in the computations of the divided differences in Eq. (5.5), we present piecewise definitions for the same. In this approach, we switch the computations to the corresponding series expansions of the same should the difference of the independent variables be less than some threshold. These threshold values are chosen such that we retain as many significant digits as possible. This methodology (technique) is simple and systematic. It is explained in full detail using an example by means of which we answer the three questions raised in this context by Kahan and Darcy [12]: 1) What value should be assigned to the threshold in this technique? 2) How many terms in the series approximation should this technique retain? and 3) How accurate is this technique?

The elements of the first subset of Eq. (5.5) can be computed [14] to machine precision by rearranging them to the following functional form<sup>2</sup>.

$$(5.19) \quad \text{sinhc}(x) := \begin{cases} \frac{\sinh(x)}{x} & \text{if } x \neq 0, \\ 1 & \text{if } x = 0. \end{cases}$$

$$(5.20) \quad \exp[x_1; x_2] = e^{(x_1+x_2)/2} \text{sinhc}\left(\frac{x_2-x_1}{2}\right), \quad q(x) = e^{x/2} \text{sinhc}\left(\frac{x}{2}\right)$$

We now explain the methodology used to arrive at piecewise definitions of divided differences using  $\exp[x_1; x_2; x_3]$  as an example. This term can be written as

$$(5.21) \quad \exp[x_1; x_2; x_3] = e^{x_2} \exp[x_1 - x_2; 0; x_3 - x_2] = e^{x_2} \frac{q(x_3 - x_2) - q(x_1 - x_2)}{x_3 - x_1}$$

where the function  $q(x)$  is evaluated as shown in Eq. (5.20). Without loss of generality, we assume  $x_1 \leq x_2 \leq x_3$ . Consequently we have,

$$(5.22) \quad \forall \xi \in [x_1, x_3], \quad |\xi - x_2| \leq (x_3 - x_1)$$

$$(5.23) \quad x_\sigma^2 \leq (x_3 - x_1)^2, \quad |x_\tau^3| \leq (x_3 - x_1)^3$$

In the computations of divided differences, the loss of significance is due to the subtractive cancellations that occur in the dependent variables which is brought to prominence after a division by the difference of the independent variables. Particularly, in Eq. (5.21) the loss of significant digits is due to the cancellations that occur in the term  $q(x_3 - x_2) - q(x_1 - x_2)$ . This term admits the following series expansion.

$$(5.24) \quad q(x_3 - x_2) - q(x_1 - x_2) = \frac{x_3 - x_1}{2} [1 + (x_a - x_2) + \dots]$$

<sup>2</sup>In this form the difference of the independent variables appear symbolically as input to a function that could be evaluated to machine precision

Let  $x_3 - x_1 = 2^{-m}$  where  $m \geq 1$  is an integer. Then, Eq. (5.22) implies that the higher order terms in Eq. (5.24) tend to zero as  $m \rightarrow \infty$ . Thus,

$$(5.25) \quad \left. \begin{array}{l} x_3 - x_1 = 2^{-m} \\ x_1 \leq x_2 \leq x_3 \end{array} \right\} \Rightarrow \tilde{q}(x_3 - x_2) - q(x_1 - x_2) = O(2^{-(m+1)})$$

When written in the normalized decimal form (cf. Eq. (5.18)), the exponent of  $q(x_3 - x_2)$  and  $q(x_1 - x_2)$  will be 0 and  $-1$ , respectively. This can be inferred using Eq. (5.22) as shown below.

$$(5.26) \quad 0 \leq x \leq 2^{-m} \Rightarrow 2^0 \leq q(x) < 2^1, \quad -2^m \leq x \leq 0 \Rightarrow 2^{-1} \leq q(x) < 2^0$$

$$(5.27) \quad 0 \leq (x_3 - x_2) \leq 2^{-m} \Rightarrow q(x_3 - x_2) = (-1)^0 2^0 1.\hat{f} \approx \boxed{0 \mid 1023_b \mid 0.\hat{f}_b}$$

$$(5.28) \quad -2^{-m} \leq (x_1 - x_2) \leq 0 \Rightarrow q(x_1 - x_2) = (-1)^0 2^{-1} 1.\tilde{f} \approx \boxed{0 \mid 1022_b \mid 0.\tilde{f}_b}$$

where  $\hat{f}$  and  $\tilde{f}$  denote the significands of  $q(x_3 - x_2)$  and  $q(x_1 - x_2)$ , respectively. The subtraction  $q(x_3 - x_2) - q(x_1 - x_2)$  can be described schematically as follows.

$$(5.29) \quad \begin{aligned} q(x_3 - x_2) - q(x_1 - x_2) &= (-1)^0 2^0 1.\hat{f} - (-1)^0 2^{-1} 1.\tilde{f} \\ &\approx \boxed{0 \mid 1023_b \mid 0.\hat{f}_b} - \boxed{0 \mid 1022_b \mid 0.\tilde{f}_b} && \text{normalized form} \\ &= \boxed{0 \mid 1023_b \mid 0.\hat{f}_b} - \boxed{0 \mid 1023_b \mid 0.1\tilde{f}_b} && \text{align radix points} \\ &= \boxed{0 \mid 1023_b \mid 0.\{0\}_m 1\hat{f}_b} && O(2^{-(m+1)}) \\ &= \boxed{0 \mid (1022 - m)_b \mid 0.\hat{f}_b} && \text{normalized form} \end{aligned}$$

The notation  $\{0\}_m$  means that the bits within the braces are repeated  $m$  times. We see that among the stored 52 bits of  $\hat{f}_b$  and  $\tilde{f}_b$ , the first  $m$  bits are lost due to cancellation. After subtraction, the unit bit at the  $m + 1^{\text{th}}$  place will become the implicit bit of the result which is not stored, cf. Eq. (5.18). The exponent of the result will become  $-(m + 1)$ . The significand of the result will become the remaining bits denoted in Eq. (5.29) as  $\hat{f}_b$  of which only  $51 - m$  bits are significant.

We see that if  $\exp[x_1; x_2; x_3]$  is evaluated using Eq. (5.21) we loose significant bits at a first order rate. We will call this form of computation as the *direct computation*. If  $x_1 \leq x_2 \leq x_3$  and  $x_3 - x_1 = 2^{-m}$  then in the direct computation of  $\exp[x_1; x_2; x_3]$  we are left with  $51 - m$  significant bits in the significand.

Using Eq. (5.16) we can write the series expansion for  $\exp[x_1; x_2; x_3]$  as follows.

$$(5.30) \quad \exp[x_1; x_2; x_3] = \frac{e^{x_a}}{2} \mathcal{S}, \quad \mathcal{S} := \left[ 1 + \frac{x_\sigma^2}{8} + \frac{x_t^3}{60} + \frac{x_\sigma^4}{160} + \frac{x_\sigma^2 x_t^3}{840} + \dots \right]$$

The above form to evaluate  $\exp[x_1; x_2; x_3]$  will be called as the *series computation*. Clearly, in the series computation we do not find removable singularities which imply that subtractive cancellations (if any) are not brought to prominence. However, the truncation of the series will introduce an error which will limit the number of significant digits in the series computation that match those in an exact computation. When the series  $\mathcal{S}$  is truncated after the first  $n$  terms it will be denoted as  $\mathcal{S}_n$ . As  $\exp(x_a)/2$  can be evaluated to machine precision, the number of significant digits in the series computation is essentially limited by the term



$\mathcal{S}_n$ . The series  $\mathcal{S}$  when written in the normalized decimal form has a zero exponent when  $x_1 \leq x_2 \leq x_3$  and  $x_3 - x_1 = 2^{-m}$ . This can be inferred using Eqs. (5.16) and (5.23) as follows.

$$(5.31) \quad \exists \xi \in [x_1, x_3] \text{ such that, } \frac{e^{x_a}}{2} \left[ 1 + \frac{x_\sigma^2}{8} + \frac{x_t^3}{60} + \dots \right] = \frac{e^{x_a}}{2} + \frac{x_\sigma^2}{16} e^\xi$$

$$(5.32) \quad \Rightarrow 1 + \frac{x_\sigma^2}{8} + \frac{x_t^3}{60} + \dots = 1 + \frac{x_\sigma^2}{8} e^{\xi - x_a}$$

$$(5.33) \quad \Rightarrow 2^0 \leq 1 + \frac{x_\sigma^2}{8} e^{\xi - x_a} \leq 1 + 2^{-(2m+3)} e^{2^{-m}} < 2^1$$

It follows that all terms except the first one contribute to the significand of  $\mathcal{S}$ . Thus,

$$(5.34) \quad \mathcal{S} = (-1)^0 2^0 1.f \approx \boxed{0 \mid 1023_b \mid 0.f_b}$$

Hence, when  $\mathcal{S}$  is replaced by  $\mathcal{S}_n$ , the associated truncation error can be understood as to limit the number of significant digits in the series computation. The truncation error associated to  $\mathcal{S}_n$  is denoted as  $\mathcal{E}_n$ . From Eqs. (5.23) and (5.30) we infer,

$$(5.35) \quad \mathcal{E}_1 = O\left(\frac{x_\sigma^2}{8}\right) \leq O(2^{-(2m+3)}), \quad \mathcal{E}_2 = O\left(\frac{x_t^3}{60}\right) \leq O(2^{-(3m+6)})$$

$$(5.36) \quad \mathcal{E}_3 = O\left(\frac{x_\sigma^4}{160}\right) \leq O(2^{-(4m+8)}), \quad \mathcal{E}_4 = O\left(\frac{x_\sigma^2 x_t^3}{840}\right) \leq O(2^{-(5m+10)})$$

Expressing  $\mathcal{S}_n = \mathcal{S} - \mathcal{E}_n$  in the double storage format we get,

$$(5.37) \quad \mathcal{S}_1 \approx \boxed{0 \mid 1023_b \mid 0.f_b} - \boxed{0 \mid 1023_b \mid 0.\{0\}_{2m+2}1\dots}$$

$$(5.38) \quad \mathcal{S}_2 \approx \boxed{0 \mid 1023_b \mid 0.f_b} - \boxed{0 \mid 1023_b \mid 0.\{0\}_{3m+5}1\dots}$$

$$(5.39) \quad \mathcal{S}_3 \approx \boxed{0 \mid 1023_b \mid 0.f_b} - \boxed{0 \mid 1023_b \mid 0.\{0\}_{4m+7}1\dots}$$

$$(5.40) \quad \mathcal{S}_4 \approx \boxed{0 \mid 1023_b \mid 0.f_b} - \boxed{0 \mid 1023_b \mid 0.\{0\}_{5m+9}1\dots}$$

where  $\mathcal{E}_n$  is written after the alignment of radix points and the remaining digits in the significands are denoted by ellipsis. This implies that we have  $(2m+2)$ ,  $(3m+5)$ ,  $(4m+7)$  and  $(5m+9)$  significant digits in  $\mathcal{S}_1$ ,  $\mathcal{S}_2$ ,  $\mathcal{S}_3$  and  $\mathcal{S}_4$ , respectively.

For each  $\mathcal{S}_n$  we solve for  $m$  by matching the accuracy of the series computation with the one obtained in the direct computation. In this way, we obtain the threshold value of  $(x_3 - x_1) = 2^{-m}$  and the lower bound for the number of significant digits nsd in a piecewise computation of  $\exp[x_1; x_2; x_3]$ . Thus,

$$(5.41) \quad \mathcal{S}_1: \quad 51 - m = 2m + 2 \quad \Rightarrow \quad m = 16, \quad \text{nsd} = 34 \text{ bits} \approx 11 \text{ decimal digits}$$

$$(5.42) \quad \mathcal{S}_2: \quad 51 - m = 3m + 5 \quad \Rightarrow \quad m = 12, \quad \text{nsd} = 39 \text{ bits} \approx 12 \text{ decimal digits}$$

$$(5.43) \quad \mathcal{S}_3: \quad 51 - m = 4m + 7 \quad \Rightarrow \quad m = 9, \quad \text{nsd} = 42 \text{ bits} \approx 13 \text{ decimal digits}$$

$$(5.44) \quad \mathcal{S}_4: \quad 51 - m = 5m + 9 \quad \Rightarrow \quad m = 7, \quad \text{nsd} = 44 \text{ bits} \approx 14 \text{ decimal digits}$$

In the above equations, the solution for  $m$  is rounded to the nearest integer. Using this rounded  $m$  we estimate nsd as the minimum of the number of significant digits found in the direct and the series computations. As the loss of significant digits is bounded, the piecewise computation of  $\exp[x_1; x_2; x_3]$  is stable.

TABLE 2  
Loss of significance controlled in the piecewise computations of  $\exp[1; 1 + \varepsilon; 1 + 2\varepsilon]$ .

$\varepsilon$	Formula3 computation	Exact 16 digits	Formula4 computation
$10^{-01}$	1.503 335 165 136 323	1.503 335 165 136 325	1.503 335 165 136 323
$10^{-02}$	1.372 811 947 550 791	1.372 811 947 550 820	1.372 811 947 550 791
$10^{-03}$	1.360 500 848 316 010	1.360 500 848 315 854	1.360 500 848 315 854
$10^{-04}$	1.359 276 836 253 229	1.359 276 836 249 607	1.359 276 836 249 607
$10^{-05}$	1.359 154 505 691 532	1.359 154 505 717 948	1.359 154 505 717 948
$10^{-06}$	1.359 142 273 371 116	1.359 142 273 371 229	1.359 142 273 371 229
$10^{-07}$	1.359 141 050 143 620	1.359 141 050 143 621	1.359 141 050 143 622
$10^{-08}$	1.359 140 927 820 931	1.359 140 927 820 931	1.359 140 927 820 931
$10^{-09}$	1.359 140 915 588 663	1.359 140 915 588 663	1.359 140 915 588 663
$10^{-10}$	1.359 140 914 365 436	1.359 140 914 365 436	1.359 140 914 365 436
$10^{-11}$	1.359 140 914 243 114	1.359 140 914 243 114	1.359 140 914 243 114
$10^{-12}$	1.359 140 914 230 881	1.359 140 914 230 881	1.359 140 914 230 881
$10^{-13}$	1.359 140 914 229 658	1.359 140 914 229 658	1.359 140 914 229 658
$10^{-14}$	1.359 140 914 229 536	1.359 140 914 229 536	1.359 140 914 229 536
$10^{-15}$	1.359 140 914 229 524	1.359 140 914 229 523	1.359 140 914 229 524

The numerical test presented in §5.1 is repeated here and the computations of the piecewise definitions are shown in Table 2. The piecewise definitions considering  $\mathcal{S}_1$  and  $\mathcal{S}_4$  for the series computations are called Formula3 and Formula4, respectively. The significant digits in both formula computations that differ from the exact values are highlighted in green color. The lower bounds for the number of significant digits given in Eqs. (5.41) and (5.44) are reproduced in this test for Formula3 and Formula4, respectively.

**5.5. Stable formulas for exponential divided differences.** Following the methodology described in the previous section we present stable piecewise definitions of all the expressions that belong to the subsets in Eq. (5.5). In the series computation of each piecewise definition, we consider the first four terms in the corresponding series expansion. Recall that the exponential function is its own derivative. This result along with the abstraction (e.g.  $x_p, x_t$  etc.) in the optimal series expansion permits us to use multiple terms in the series expansion without incurring substantial computational cost.

The elements of the first subset in Eq. (5.5) can be evaluated to machine precision without resorting to a series computation, cf. Eq. (5.20). The first element of the second subset, i.e.  $\exp[x_1; x_2; x_3]$  was used as an example to describe the details of the piecewise computation technique in the previous section. The stable piecewise definition of the same when  $x_1 \leq x_2 \leq x_3$  can be summarized as follows.

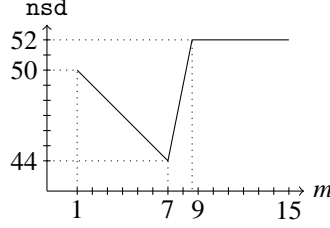
$$(5.45) \quad \exp[x_1; x_2; x_3] = \begin{cases} e^{x_2} \frac{q(x_3 - x_2) - q(x_1 - x_2)}{x_3 - x_1} & \text{if } (x_3 - x_1) > 2^{-7} \\ \frac{e^{x_a}}{2} \left[ 1 + \frac{x_\sigma^2}{8} + \frac{x_t^3}{60} + \frac{x_\sigma^4}{160} \right] & \text{else} \end{cases}$$

It is essential to sort the arguments lest the series computation should incur a significant truncation error. The variation in the number of significant digits  $\text{nsd}$  in  $\exp[x_1; x_2; x_3]$  with respect to  $m$ , where  $(x_3 - x_1) = 2^{-m}$ , is denoted as  $\text{nsd}(\exp[x_1; x_2; x_3], m)$ . Using the above

stable formula for  $\exp[x_1; x_2; x_3]$  we obtain,

(5.46)

$$\text{nsd}(\exp[x_1; x_2; x_3], m) = \begin{cases} 51 - m & \text{if } m < 7 \\ 5m + 9 & \text{if } 7 \leq m < 8.6, \\ 52 & \text{if } m \geq 8.6 \end{cases}$$



Recall that  $q[x_1; x_2] = \exp[x_1; 0; x_2]$ ,  $r(x) = \exp[0; 0; x]$  and  $\Phi(x, y) = \exp[-iy; iy; x]$ . Let  $\text{sort}$  be a sorting function and  $\hat{x}_1 \leq \hat{x}_2 \leq \hat{x}_3$ . Then, using Eq. (5.45) a stable formula for  $q[x_1; x_2]$  is

$$(5.47) \quad \{\hat{x}_1, \hat{x}_2, \hat{x}_3\} = \text{sort}(\{x_1, 0, x_2\}), \quad q[x_1; x_2] = \exp[\hat{x}_1; \hat{x}_2; \hat{x}_3];$$

Likewise a stable formula for  $r(x)$  is

$$(5.48) \quad \{\hat{x}_1, 0, \hat{x}_3\} = \text{sort}(\{0, 0, x\}), \quad r(x) = \exp[\hat{x}_1; 0; \hat{x}_3];$$

As  $\exp[-iy; iy; x]$  involves complex numbers we give it special attention. Recall that the exponential function is holomorphic, i.e. it is complex differentiable in a neighbourhood of every point in its domain. This implies that it is infinitely differentiable and is equal to its own Taylor series. Thus, the optimal series approximation of divided differences presented in §5.2 naturally extends to  $\exp[-iy; iy; x]$ . Following this line, a stable formula for  $\Phi(x, y)$  can be obtained as shown below.

$$(5.49) \quad z := x + iy, \quad z_a := \frac{x}{3}, \quad z_\sigma^2 := 2z_a^2 - \frac{2}{3}y^2, \quad z_t^3 := 2z_a(z_a^2 + y^2)$$

$$(5.50) \quad \Phi(x, y) = \exp[-iy; iy; x] = \begin{cases} \frac{e^{z/2} \text{sinhc}(z^*/2) - \text{sinc}(y)}{z} & \text{if } |z| > 2^{-7} \\ \frac{e^{z_a}}{2} \left[ 1 + \frac{z_\sigma^2}{8} + \frac{z_t^3}{60} + \frac{z_\sigma^4}{160} \right] & \text{else} \end{cases}$$

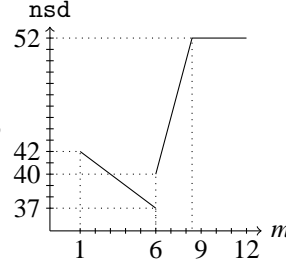
The above definition assumes the availability of a complex math library which provides an interface for stable computation of common arithmetic operations, elementary and transcendental functions. This assumption holds for the C++ programming language which is equipped with the standard math library `<complex>`.

The divided difference  $\exp[x_1; x_2; x_3; x_4]$  is a template for the elements of the third subset in Eq. (5.5). The proposed piecewise definition of  $\exp[x_1; x_2; x_3; x_4]$  when  $x_1 \leq x_2 \leq x_3 \leq x_4$  is

$$(5.51) \quad \exp[x_1; x_2; x_3; x_4] = \begin{cases} \frac{\exp[x_2; x_3; x_4] - \exp[x_1; x_2; x_3]}{x_4 - x_1} & \text{if } (x_4 - x_1) > 2^{-6} \\ \frac{e^{x_a}}{3!} \left[ 1 - \frac{x_p^2}{20} + \frac{x_t^3}{120} + \frac{x_p^4 - x_q^4}{840} \right] & \text{else} \end{cases}$$

For the above piecewise definition of  $\exp[x_1; x_2; x_3; x_4]$  we obtain  
(5.52)

$$\text{nsd}(\exp[x_1; x_2; x_3; x_4], m) = \begin{cases} 43 - m & \text{if } m < 6 \\ 5m + 10 & \text{if } 6 \leq m < 8.4, \\ 52 & \text{if } m \geq 8.4 \end{cases}$$



It is straightforward to verify that  $q[x_1; x_2; x_3] = \exp[0; x_1; x_2; x_3]$ ,  $r[x_1; x_2] = \exp[0; 0; x_1; x_2]$  and  $\Phi(\star, y)[x_1; x_2] = \exp[-iy; iy; x_1; x_2]$ . Using Eq. (5.51) a stable formula for  $q[x_1; x_2; x_3]$  is

$$(5.53) \quad \{\hat{x}_1, \hat{x}_2, \hat{x}_3, \hat{x}_4\} = \text{sort}(\{0, x_1, x_2, x_3\}), \quad q[x_1; x_2; x_3] = \exp[\hat{x}_1; \hat{x}_2; \hat{x}_3; \hat{x}_4];$$

Likewise a stable formula for  $r[x_1; x_2]$  is

$$(5.54) \quad \{\hat{x}_1, \hat{x}_2, \hat{x}_3, \hat{x}_4\} = \text{sort}(\{0, 0, x_1, x_2\}), \quad r[x_1; x_2] = \exp[\hat{x}_1; \hat{x}_2; \hat{x}_3; \hat{x}_4];$$

The proposed piecewise definition of  $\Phi(\star, y)[x_1; x_2]$  is

$$(5.55) \quad \{\hat{x}_1, \hat{x}_2\} = \text{sortabs}(\{x_1, x_2\}), \quad z_1 := \hat{x}_1 + iy, \quad z_2 := \hat{x}_2 + iy, \quad z_a := \frac{\hat{x}_1 + \hat{x}_2}{4}$$

$$(5.56) \quad z_p^2 := y^2 + \hat{x}_1 \hat{x}_2 - 6z_a^2, \quad z_t^3 := 2z_a(y^2 - \hat{x}_1 \hat{x}_2 + 4z_a^2), \quad z_q^4 := (y^2 + z_a^2)(\hat{x}_1 \hat{x}_2 - 3z_a^2)$$

$$(5.57) \quad \Phi(\star, y)[x_1; x_2] = \begin{cases} \left[ \frac{\exp[\hat{x}_1; \hat{x}_2] - e^{z_1/2} \text{sinhc}(z_1^*/2)}{z_2^*} - \Phi(\hat{x}_1, y) \right] \frac{1}{z_2} & \text{if } |z_2| > 2^{-6} \\ \frac{e^{z_a}}{3!} \left[ 1 - \frac{z_p^2}{20} + \frac{z_t^3}{120} + \frac{z_p^4 - z_q^4}{840} \right] & \text{else} \end{cases}$$

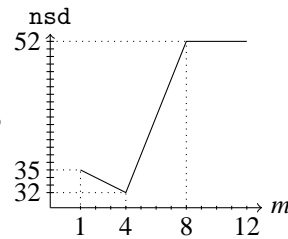
where  $\text{sortabs}$  is a function that sorts its arguments with respect to its absolute value, i.e.  $|\hat{x}_1| \leq |\hat{x}_2|$  in Eq. (5.55). The term  $\Phi(\hat{x}_1, y)$  is evaluated using the stable formula given in Eq. (5.50).

The divided difference  $\exp[x_1; x_2; x_3; x_4; x_5]$  is a template for the elements of the third subset in Eq. (5.5). The proposed piecewise definition of  $\exp[x_1; x_2; x_3; x_4; x_5]$  when  $x_1 \leq x_2 \leq x_3 \leq x_4 \leq x_5$  is

$$(5.58) \quad \exp[x_1; x_2; x_3; x_4; x_5] = \begin{cases} \frac{\exp[x_2; x_3; x_4; x_5] - \exp[x_1; x_2; x_3; x_4]}{x_5 - x_1} & \text{if } (x_5 - x_1) > 2^{-4} \\ \frac{e^{x_a}}{4!} \left[ 1 - \frac{x_p^2}{30} + \frac{x_t^3}{210} + \frac{x_p^4 - x_q^4}{1680} \right] & \text{else} \end{cases}$$

For the above piecewise definition of  $\exp[x_1; x_2; x_3; x_4; x_5]$  we obtain  
(5.59)

$$\text{nsd}(\exp[x_1; x_2; x_3; x_4; x_5], m) = \begin{cases} 36 - m & \text{if } m < 4 \\ 5m + 12 & \text{if } 4 \leq m < 8, \\ 52 & \text{if } m \geq 8 \end{cases}$$



It is straightforward to verify that  $r[x_1; x_2; x_3] = \exp[0; 0; x_1; x_2; x_3]$  and  $\Phi(\star, y)[x_1; x_2; x_3] = \exp[-iy; iy; x_1; x_2; x_3]$ . Using Eq. (5.58) a stable formula for  $r[x_1; x_2; x_3]$  can be written as,

$$(5.60) \quad \{\widehat{x}_1, \widehat{x}_2, \widehat{x}_3, \widehat{x}_4, \widehat{x}_5\} = \text{sort}(\{0, 0, x_1, x_2, x_3\})$$

$$(5.61) \quad r[x_1; x_2; x_3] = \exp[\widehat{x}_1; \widehat{x}_2; \widehat{x}_3; \widehat{x}_4; \widehat{x}_5];$$

The proposed piecewise definition of  $\Phi(\star, y)[x_1; x_2; x_3]$  is

$$(5.62) \quad \{\widehat{x}_1, \widehat{x}_2, \widehat{x}_3\} = \text{sortabs}(\{x_1, x_2, x_3\}), \quad z_3 := \widehat{x}_3 + iy$$

$$(5.63) \quad z_a := \frac{\widehat{x}_1 + \widehat{x}_2 + \widehat{x}_3}{5}, \quad z_q^4 := 6z_a^4 + \frac{1}{2}[11z_a^2 - (\widehat{x}_1^2 + \widehat{x}_2^2 + \widehat{x}_3^2)](y^2 + 3z_a^2) - 2z_a\widehat{x}_1\widehat{x}_2\widehat{x}_3$$

$$(5.64)$$

$$z_p^2 := \frac{1}{2}[2y^2 + 5z_a^2 - (\widehat{x}_1^2 + \widehat{x}_2^2 + \widehat{x}_3^2)], \quad z_t^3 := \frac{z_a}{2}[4y^2 - 35z_a^2 + 3(\widehat{x}_1^2 + \widehat{x}_2^2 + \widehat{x}_3^2)] + \widehat{x}_1\widehat{x}_2\widehat{x}_3$$

$$(5.65)$$

$$\Phi(\star, y)[x_1; x_2; x_3] = \begin{cases} \left[ \frac{\exp[\widehat{x}_1; \widehat{x}_2; \widehat{x}_3] - \exp[iy; \widehat{x}_1; \widehat{x}_2]}{z_3^*} - \Phi(\star, y)[\widehat{x}_1; \widehat{x}_2] \right] \frac{1}{z_3} & \text{if } |z_3| > 2^{-4} \\ \frac{e^{z_a}}{4!} \left[ 1 - \frac{z_p^2}{30} + \frac{z_t^3}{210} + \frac{z_q^4}{1680} \right] & \text{else} \end{cases}$$

where the term  $\Phi(\star, y)[\widehat{x}_1; \widehat{x}_2]$  is evaluated using the stable formula given in Eq. (5.57). Note that in the stable formula for  $\Phi(\star, y)[\widehat{x}_1; \widehat{x}_2]$ , just the direct computation of  $\exp[iy; \widehat{x}_1; \widehat{x}_2]$  is sufficient as the threshold value of  $|z_2|$  to switch to a series computation is larger for the former than the latter. This means that the series computation of  $\exp[iy; \widehat{x}_1; \widehat{x}_2]$  will never be used in the stable computation of  $\Phi(\star, y)[\widehat{x}_1; \widehat{x}_2]$ . On the contrary, in the stable computation of  $\Phi(\star, y)[x_1; x_2; x_3]$  the switch to the series computation is governed by some threshold value of  $|z_3|$  which includes the possibility  $|z_2| \rightarrow 0$ . Therefore, in Eq. (5.65) it is necessary to evaluate the term  $\exp[iy; \widehat{x}_1; \widehat{x}_2]$  in a piecewise manner.

Following Eqs. (5.45) and (5.50), a stable formula for  $\exp[iy; \widehat{x}_1; \widehat{x}_2]$  can be obtained as follows.

$$(5.66) \quad z_1 := \widehat{x}_1 + iy, \quad z_2 := \widehat{x}_2 + iy, \quad z_a := \frac{\widehat{x}_1 + \widehat{x}_2 + iy}{3}, \quad \widetilde{z}_1 := iy - z_a$$

$$(5.67) \quad \widetilde{z}_2 := \widehat{x}_1 - z_a, \quad \widetilde{z}_3 := \widehat{x}_2 - z_a, \quad z_\sigma^2 := \frac{\widetilde{z}_1^2 + \widetilde{z}_2^2 + \widetilde{z}_3^2}{3}, \quad z_t^3 := \widetilde{z}_1\widetilde{z}_2\widetilde{z}_3$$

$$(5.68) \quad \exp[iy; \widehat{x}_1; \widehat{x}_2] = \begin{cases} \frac{\exp[\widehat{x}_1; \widehat{x}_2] - e^{z_1/2} \text{sinhc}(z_1^*/2)}{z_2^*} & \text{if } |z_2| > 2^{-7} \\ \frac{e^{z_a}}{2} \left[ 1 + \frac{z_\sigma^2}{8} + \frac{z_t^3}{60} + \frac{z_\sigma^4}{160} \right] & \text{else} \end{cases}$$

**6. Examples.** We present two examples to validate the numerical stability in the computation of the proposed formulas for the X-IVAS scheme. In these examples the eigenvalues of the matrix  $\mathbf{A}$  gradually tends to zero. The symbolic computation of the formulas for the chosen eigenvalues are done using Maple and the first 16 significant decimal digits are stored as reference solutions. These reference solutions are used to measure the relative error in the formula computations using double precision floating point arithmetic.

**6.1. Example 1.** Consider the case when two of the eigenvalues of the matrix  $\mathbf{A}$  are complex numbers. Let  $\varepsilon := 10^{-n}$  and choose  $n \in \{1, 2, 3, \dots, 15\}$ . For each  $\varepsilon$  define the

TABLE 3  
Relative errors in the usual and stable computation when  $\lambda = \{\varepsilon \pm i2\varepsilon, 4\varepsilon\}$ .

$\varepsilon$	$\frac{\ \mathbf{R}^u - \mathbf{R}\ }{\ \mathbf{R}\ }$	$\frac{\ \mathbf{R}^s - \mathbf{R}\ }{\ \mathbf{R}\ }$	$\frac{\ \mathbf{Q}^u - \mathbf{Q}\ }{\ \mathbf{Q}\ }$	$\frac{\ \mathbf{Q}^s - \mathbf{Q}\ }{\ \mathbf{Q}\ }$	$\frac{\ \mathbf{P}^u - \mathbf{P}\ }{\ \mathbf{P}\ }$	$\frac{\ \mathbf{P}^s - \mathbf{P}\ }{\ \mathbf{P}\ }$
$10^{-01}$	$6.1 \cdot 10^{-15}$	$7.4 \cdot 10^{-15}$	$5.3 \cdot 10^{-16}$	$1.1 \cdot 10^{-15}$	$1.1 \cdot 10^{-16}$	$2.6 \cdot 10^{-16}$
$10^{-02}$	$6.7 \cdot 10^{-12}$	$4.8 \cdot 10^{-13}$	$6.9 \cdot 10^{-14}$	$7.5 \cdot 10^{-15}$	$1.1 \cdot 10^{-15}$	$9.7 \cdot 10^{-17}$
$10^{-03}$	$5.7 \cdot 10^{-09}$	$4.3 \cdot 10^{-17}$	$6.1 \cdot 10^{-12}$	$0.0 \cdot 10^{+00}$	$1.1 \cdot 10^{-14}$	$9.9 \cdot 10^{-17}$
$10^{-04}$	$4.9 \cdot 10^{-06}$	$6.9 \cdot 10^{-17}$	$5.2 \cdot 10^{-10}$	$0.0 \cdot 10^{+00}$	$9.8 \cdot 10^{-14}$	$0.0 \cdot 10^{+00}$
$10^{-05}$	$9.4 \cdot 10^{-04}$	$1.2 \cdot 10^{-16}$	$1.0 \cdot 10^{-08}$	$0.0 \cdot 10^{+00}$	$1.8 \cdot 10^{-13}$	$2.1 \cdot 10^{-21}$
$10^{-06}$	$2.4 \cdot 10^{+00}$	$2.1 \cdot 10^{-16}$	$2.6 \cdot 10^{-06}$	$5.9 \cdot 10^{-17}$	$4.9 \cdot 10^{-12}$	$9.9 \cdot 10^{-17}$
$10^{-07}$	$2.8 \cdot 10^{+03}$	$3.0 \cdot 10^{-17}$	$3.0 \cdot 10^{-04}$	$8.3 \cdot 10^{-17}$	$5.7 \cdot 10^{-11}$	$9.9 \cdot 10^{-17}$
$10^{-08}$	$2.8 \cdot 10^{+06}$	$0.0 \cdot 10^{+00}$	$3.0 \cdot 10^{-02}$	$8.3 \cdot 10^{-17}$	$5.7 \cdot 10^{-10}$	$9.9 \cdot 10^{-17}$
$10^{-09}$	$2.2 \cdot 10^{+09}$	$3.0 \cdot 10^{-17}$	$2.4 \cdot 10^{+00}$	$5.9 \cdot 10^{-17}$	$4.5 \cdot 10^{-09}$	$9.9 \cdot 10^{-17}$
$10^{-10}$	$7.2 \cdot 10^{+11}$	$0.0 \cdot 10^{+00}$	$7.7 \cdot 10^{+01}$	$0.0 \cdot 10^{+00}$	$1.4 \cdot 10^{-08}$	$0.0 \cdot 10^{+00}$
$10^{-11}$	$7.2 \cdot 10^{+13}$	$0.0 \cdot 10^{+00}$	$7.7 \cdot 10^{+02}$	$1.2 \cdot 10^{-27}$	$1.4 \cdot 10^{-08}$	$2.0 \cdot 10^{-27}$
$10^{-12}$	$1.9 \cdot 10^{+18}$	$2.1 \cdot 10^{-16}$	$2.0 \cdot 10^{+06}$	$5.9 \cdot 10^{-17}$	$3.8 \cdot 10^{-06}$	$9.9 \cdot 10^{-17}$
$10^{-13}$	$2.7 \cdot 10^{+21}$	$2.1 \cdot 10^{-16}$	$2.9 \cdot 10^{+08}$	$5.9 \cdot 10^{-17}$	$5.4 \cdot 10^{-05}$	$9.9 \cdot 10^{-17}$
$10^{-14}$	$6.9 \cdot 10^{+23}$	$6.2 \cdot 10^{-31}$	$7.4 \cdot 10^{+09}$	$0.0 \cdot 10^{+00}$	$1.4 \cdot 10^{-04}$	$1.9 \cdot 10^{-30}$
$10^{-15}$	$6.9 \cdot 10^{+25}$	$3.0 \cdot 10^{-17}$	$7.4 \cdot 10^{+10}$	$1.4 \cdot 10^{-31}$	$1.4 \cdot 10^{-04}$	$0.0 \cdot 10^{+00}$

matrix  $\mathbf{A}$  and a corresponding auxiliary matrix  $\mathbf{Z}$  as follows.

$$(6.1) \quad \mathbf{A} := \begin{bmatrix} a + \varepsilon & -2\varepsilon & b \\ 2\varepsilon & a + \varepsilon & c \\ 0 & 0 & a + d\varepsilon \end{bmatrix} \Rightarrow \text{eigs}(\mathbf{A}) := \begin{bmatrix} \alpha \pm i\beta \\ \lambda_3 \end{bmatrix} = \begin{bmatrix} a + \varepsilon \pm i2\varepsilon \\ a + d\varepsilon \end{bmatrix}$$

$$(6.2) \quad \mathbf{Z} := (\mathbf{A} - \alpha\mathbf{I})^2 + \beta^2\mathbf{I} = \varepsilon \begin{bmatrix} 0 & 0 & b(d-1) + c\varepsilon \\ 0 & 0 & c(d-1) + b\varepsilon \\ 0 & 0 & d(d-2)\varepsilon + 5\varepsilon \end{bmatrix} \Rightarrow \begin{aligned} \mathbf{Z}\mathbf{A} &= (a + d\varepsilon)\varepsilon\mathbf{Z} \\ \mathbf{Z}\mathbf{A}^2 &= (a + d\varepsilon)^2\varepsilon\mathbf{Z} \end{aligned}$$

where  $\text{eigs}(\mathbf{A})$  represents the eigenvalues of  $\mathbf{A}$ . We can drive all the eigenvalues and/or the gap between them to zero by appropriately choosing the parameters  $a$  and  $d$ . For each  $\mathbf{A}$  we compute  $\mathbf{R}(\tau, \mathbf{A})$  using the stable formulas summarized in the previous section. The matrices  $\mathbf{P}(\tau, \mathbf{A})$  and  $\mathbf{Q}(\tau, \mathbf{A})$  are computed from  $\mathbf{R}(\tau, \mathbf{A})$  using the relationships given in Eqs. (3.16) and (3.18).

Table 3 illustrates the relative errors in the computations of the matrices  $\mathbf{P}(\tau, \mathbf{A})$ ,  $\mathbf{Q}(\tau, \mathbf{A})$  and  $\mathbf{R}(\tau, \mathbf{A})$  choosing  $\tau = 1$ ,  $a = 0$ ,  $b = c = 1$  and  $d = 4$ . Therein  $\mathbf{R}^u$  and  $\mathbf{R}^s$  denotes the usual (naïve) and the stable (piecewise) computations of the matrix  $\mathbf{R}$ , respectively. The Frobenius norm is used in  $\|\mathbf{R}\|$ . The relative errors  $\|\mathbf{P}^s - \mathbf{P}\|/\|\mathbf{P}\|$ ,  $\|\mathbf{Q}^s - \mathbf{Q}\|/\|\mathbf{Q}\|$  and  $\|\mathbf{R}^s - \mathbf{R}\|/\|\mathbf{R}\|$  are found to be within the guaranteed computation accuracies established for the same and reflect the robustness of the stable formulas. The gradual loss of significance as  $\varepsilon \rightarrow 0$  is reflected as a gradual increase in the relative error (from machine epsilon to values intolerably high) in the usual computations of the considered matrices. The maximum relative error in the computations of  $\mathbf{P}^u$ ,  $\mathbf{Q}^u$  and  $\mathbf{R}^u$  are of the order of  $10^{-4}$ ,  $10^{10}$  and  $10^{25}$ , respectively. As  $\varepsilon \rightarrow 0$  we observe  $(\|\mathbf{P}^u - \mathbf{P}\|/\|\mathbf{P}\|)$  is  $O(\varepsilon)$  times smaller than  $(\|\mathbf{Q}^u - \mathbf{Q}\|/\|\mathbf{Q}\|)$  which in turn

TABLE 4  
Relative errors in the usual and stable computation when  $\lambda = \{\varepsilon, 2\varepsilon, 3\varepsilon\}$ .

$\varepsilon$	$\frac{\ \mathbf{R}^u - \mathbf{R}\ }{\ \mathbf{R}\ }$	$\frac{\ \mathbf{R}^s - \mathbf{R}\ }{\ \mathbf{R}\ }$	$\frac{\ \mathbf{Q}^u - \mathbf{Q}\ }{\ \mathbf{Q}\ }$	$\frac{\ \mathbf{Q}^s - \mathbf{Q}\ }{\ \mathbf{Q}\ }$	$\frac{\ \mathbf{P}^u - \mathbf{P}\ }{\ \mathbf{P}\ }$	$\frac{\ \mathbf{P}^s - \mathbf{P}\ }{\ \mathbf{P}\ }$
$10^{-01}$	$3.9 \cdot 10^{-14}$	$7.1 \cdot 10^{-14}$	$1.7 \cdot 10^{-15}$	$3.2 \cdot 10^{-15}$	$5.4 \cdot 10^{-16}$	$4.0 \cdot 10^{-16}$
$10^{-02}$	$1.0 \cdot 10^{-11}$	$7.1 \cdot 10^{-14}$	$1.2 \cdot 10^{-13}$	$2.1 \cdot 10^{-14}$	$1.0 \cdot 10^{-15}$	$2.1 \cdot 10^{-15}$
$10^{-03}$	$1.1 \cdot 10^{-07}$	$1.3 \cdot 10^{-16}$	$6.8 \cdot 10^{-11}$	$7.8 \cdot 10^{-17}$	$6.6 \cdot 10^{-14}$	$1.1 \cdot 10^{-16}$
$10^{-04}$	$3.8 \cdot 10^{-05}$	$1.3 \cdot 10^{-16}$	$1.4 \cdot 10^{-09}$	$7.9 \cdot 10^{-17}$	$9.3 \cdot 10^{-14}$	$8.2 \cdot 10^{-17}$
$10^{-05}$	$3.5 \cdot 10^{-02}$	$5.2 \cdot 10^{-17}$	$4.2 \cdot 10^{-07}$	$0.0 \cdot 10^{+00}$	$8.1 \cdot 10^{-12}$	$0.0 \cdot 10^{+00}$
$10^{-06}$	$8.0 \cdot 10^{+01}$	$0.0 \cdot 10^{+00}$	$4.4 \cdot 10^{-05}$	$5.5 \cdot 10^{-17}$	$5.8 \cdot 10^{-11}$	$0.0 \cdot 10^{+00}$
$10^{-07}$	$4.5 \cdot 10^{-02}$	$1.7 \cdot 10^{-16}$	$6.7 \cdot 10^{-05}$	$7.9 \cdot 10^{-17}$	$7.8 \cdot 10^{-11}$	$8.2 \cdot 10^{-17}$
$10^{-08}$	$4.0 \cdot 10^{+07}$	$5.2 \cdot 10^{-17}$	$4.7 \cdot 10^{-01}$	$0.0 \cdot 10^{+00}$	$7.5 \cdot 10^{-09}$	$8.2 \cdot 10^{-17}$
$10^{-09}$	$4.0 \cdot 10^{+10}$	$0.0 \cdot 10^{+00}$	$5.5 \cdot 10^{+01}$	$5.5 \cdot 10^{-17}$	$1.0 \cdot 10^{-07}$	$8.2 \cdot 10^{-17}$
$10^{-10}$	$3.4 \cdot 10^{-01}$	$3.0 \cdot 10^{-17}$	$8.3 \cdot 10^{-02}$	$0.0 \cdot 10^{+00}$	$1.5 \cdot 10^{-08}$	$0.0 \cdot 10^{+00}$
$10^{-11}$	$3.4 \cdot 10^{-01}$	$3.0 \cdot 10^{-17}$	$8.3 \cdot 10^{-02}$	$0.0 \cdot 10^{+00}$	$1.5 \cdot 10^{-08}$	$0.0 \cdot 10^{+00}$
$10^{-12}$	$4.0 \cdot 10^{+19}$	$3.0 \cdot 10^{-17}$	$5.5 \cdot 10^{+07}$	$7.9 \cdot 10^{-17}$	$1.0 \cdot 10^{-04}$	$8.2 \cdot 10^{-17}$
$10^{-13}$	$4.0 \cdot 10^{+22}$	$0.0 \cdot 10^{+00}$	$5.5 \cdot 10^{+09}$	$5.5 \cdot 10^{-17}$	$1.0 \cdot 10^{-03}$	$8.2 \cdot 10^{-17}$
$10^{-14}$	$8.0 \cdot 10^{+25}$	$3.0 \cdot 10^{-17}$	$5.5 \cdot 10^{+11}$	$0.0 \cdot 10^{+00}$	$8.0 \cdot 10^{-03}$	$8.2 \cdot 10^{-17}$
$10^{-15}$	$8.0 \cdot 10^{+28}$	$4.2 \cdot 10^{-17}$	$5.5 \cdot 10^{+13}$	$0.0 \cdot 10^{+00}$	$6.1 \cdot 10^{-02}$	$1.1 \cdot 10^{-16}$

is  $O(\varepsilon)$  times smaller than  $(\|\mathbf{R}^u - \mathbf{R}\|/\|\mathbf{R}\|)$ . The following results explain this behaviour.

$$(6.3) \quad \mathbf{Z}\mathbf{A} = 4\varepsilon^2\mathbf{Z}, \quad \mathbf{Z}\mathbf{A}^2 = 16\varepsilon^3\mathbf{Z}$$

$$(6.4) \quad \mathbf{R}^u - \mathbf{R} \approx (\Phi^u(\star, 2\varepsilon)[- \varepsilon; - \varepsilon; 3\varepsilon] - \Phi(\star, 2\varepsilon)[- \varepsilon; - \varepsilon; 3\varepsilon])\mathbf{Z}$$

$$(6.5) \quad \mathbf{Q}^u - \mathbf{Q} = (\mathbf{R}^u - \mathbf{R})\mathbf{A}, \quad \mathbf{P}^u - \mathbf{P} = (\mathbf{R}^u - \mathbf{R})\mathbf{A}^2$$

Equation (6.4) holds because  $\Phi(\star, 2\varepsilon)[- \varepsilon; - \varepsilon; 3\varepsilon]$  is the highest-order divided difference term in Eq. 4.25 and its computation error dominates over the rest. The relative error in the usual (naïve) computation of  $\Phi(\star, 2\varepsilon)[- \varepsilon; - \varepsilon; 3\varepsilon]$  is approximately  $10^{40}$  when  $\varepsilon = 10^{-15}$ . Recall that the matrices  $\mathbf{P}(\tau, \mathbf{A})$  and  $\mathbf{Q}(\tau, \mathbf{A})$  govern the evolution of the particle positions. Likewise, the matrices  $\mathbf{Q}(\tau, \mathbf{A})$  and  $\mathbf{R}(\tau, \mathbf{A})$  govern the evolution of the particle velocities.

**6.2. Example 2.** Here the details differ from the previous example only in the definitions of the matrices  $\mathbf{A}$  and  $\mathbf{Z}$ .

$$(6.6) \quad \mathbf{A} := \begin{bmatrix} a + \varepsilon & b & c \\ 0 & a + 2\varepsilon & d \\ 0 & 0 & a + 3\varepsilon \end{bmatrix} \Rightarrow \text{eigs}(\mathbf{A}) := \begin{bmatrix} \lambda_1 \\ \lambda_2 \\ \lambda_3 \end{bmatrix} = \begin{bmatrix} a + \varepsilon \\ a + 2\varepsilon \\ a + 3\varepsilon \end{bmatrix}$$

$$(6.7) \quad \mathbf{Z} := (\mathbf{A} - \lambda_1\mathbf{I})(\mathbf{A} - \lambda_2\mathbf{I}) = \begin{bmatrix} 0 & 0 & bd + ch \\ 0 & 0 & 2dh \\ 0 & 0 & 2h^2 \end{bmatrix} \Rightarrow \begin{aligned} \mathbf{Z}\mathbf{A} &= (a + 3\varepsilon)\mathbf{Z} \\ \mathbf{Z}\mathbf{A}^2 &= (a + 3\varepsilon)^2\mathbf{Z} \end{aligned}$$

By construction, all the eigenvalues and the gap between them can be driven to zero with decreasing values of  $\varepsilon$  for appropriate choice of the parameter  $a$ .

Table 4 illustrates the relative errors in the computations of the matrices  $\mathbf{P}(\tau, \mathbf{A})$ ,  $\mathbf{Q}(\tau, \mathbf{A})$  and  $\mathbf{R}(\tau, \mathbf{A})$  choosing  $\tau = 1$ ,  $a = 0$  and  $b = c = d = 1$ . The behaviour of the usual and stable computations are similar to what is observed in the previous example. The maximum relative error in the usual computations of  $\mathbf{P}^u$ ,  $\mathbf{Q}^u$  and  $\mathbf{R}^u$  are of the order of  $10^{-2}$ ,  $10^{13}$  and  $10^{28}$ , respectively. As before  $(\|\mathbf{P}^u - \mathbf{P}\|/\|\mathbf{P}\|)$  is  $O(\varepsilon)$  times smaller than  $(\|\mathbf{Q}^u - \mathbf{Q}\|/\|\mathbf{Q}\|)$

which in turn is  $O(\varepsilon)$  times smaller than  $(\|\mathbf{R}^u - \mathbf{R}\|/\|\mathbf{R}\|)$ . The following results explain this behaviour.

$$(6.8) \quad \mathbf{Z}\mathbf{A} = 3\varepsilon\mathbf{Z}, \quad \mathbf{Z}\mathbf{A}^2 = 9\varepsilon^2\mathbf{Z}$$

$$(6.9) \quad \mathbf{R}^u - \mathbf{R} \approx (\mathbf{r}^u[\varepsilon; 2\varepsilon; 3\varepsilon] - \mathbf{r}[\varepsilon; 2\varepsilon; 3\varepsilon])\mathbf{Z}$$

$$(6.10) \quad \mathbf{Q}^u - \mathbf{Q} = (\mathbf{R}^u - \mathbf{R})\mathbf{A}, \quad \mathbf{P}^u - \mathbf{P} = (\mathbf{R}^u - \mathbf{R})\mathbf{A}^2$$

Equation (6.9) holds because  $\mathbf{r}[\varepsilon; 2\varepsilon; 3\varepsilon]$  is the highest-order divided difference term in Eq. (4.14) and its computation error dominates over the rest. The relative error in the usual (naïve) computation of  $\mathbf{r}[\varepsilon; 2\varepsilon; 3\varepsilon]$  is approximately  $10^{28}$  when  $\varepsilon = 10^{-15}$ .

**7. Conclusions.** Formula computations in the neighbourhood of removable singularities suffer loss of significance when they are done using finite precision arithmetic. Formulas for the solution of the X-IVAS scheme involve many removable singularities. Hence, the use of numerically stable formulas for the same is a criteria for robustness.

We have proposed numerically stable formulas for the closed-form analytical solution of the X-IVAS scheme. Therein, the Newton form of the polynomial interpolation definition is used for the functions of matrices which appear in the formulas. In this definition, removable singularities and the terms/expressions that participate to yield a finite limit at these points are grouped together as divided differences. In other algebraically equivalent forms, these terms/expressions get dispersed. The poor reputation of nearly confluent divided differences with respect to the loss of significance in floating point computations is a blessing in disguise. We get *a priori* warning about possible numerical instabilities in formula computations. To control the loss of significance, we have presented piecewise definitions for these divided differences. The piecewise definitions switch the computations to the respective series approximations of the divided differences should the gap between the independent variables be less than a specified threshold. These divided differences are expressible as the divided difference of the exponential function of an appropriate order less than or equal to four. For the terms involving the second, third and fourth order divided differences, the double precision floating-point computation of their piecewise definitions guarantee at least 14, 12 and 10 significant decimal digits to be exact, respectively. The implementation of these piecewise definitions is simple and the computations are stable.

**8. Acknowledgement.** I thank Mr. Guillermo Casas-González for reading the manuscript in draft form and suggesting improvements.

#### REFERENCES

- [1] MARCO CALIARI, *Accurate evaluation of divided differences for polynomial interpolation of exponential propagators*, Computing, 80 (2007), pp. 189–201.
- [2] MARCO CALIARI, ALEXANDER OSTERMANN, AND STEFAN RAINER, *Meshfree Exponential Integrators*, SIAM Journal on Scientific Computing, 35 (2013), pp. A431–A452.
- [3] DARIN P. DIACHIN AND JAMES A. HERZOG, *Analytic streamline calculations on linear tetrahedra*, in 13th Computational Fluid Dynamics Conference, Reston, Virginia, June 1997, American Institute of Aeronautics and Astronautics, pp. 733–742.
- [4] FREEMAN GILBERT AND GEORGE E. BACKUS, *Propagator matrices in elastic wave and vibration problems*, Geophysics, 31 (1966), pp. 326–332.
- [5] OPENCFD LTD (ESI GROUP), *OpenFOAM—The open source CFD toolbox*. <http://www.openfoam.org/>.
- [6] NICHOLAS J. HIGHAM, *Accuracy and Stability of Numerical Algorithms*, Society for Industrial and Applied Mathematics, Jan. 2002.
- [7] ———, *Functions of Matrices: Theory and Computation*, Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 2008.



- [8] MARLIS HOCHBRUCK, CHRISTIAN LUBICH, AND HUBERT SELHOFER, *Exponential Integrators for Large Systems of Differential Equations*, SIAM Journal on Scientific Computing, 19 (1998), pp. 1552–1574.
- [9] SERGIO R. IDELSOHN, JULIO MARTI, PABLO BECKER, AND EUGENIO OÑATE, *Analysis of multifluid flows with large time steps using the particle finite element method*, International Journal for Numerical Methods in Fluids, 75 (2014), pp. 621–644.
- [10] SERGIO R. IDELSOHN, NORBERTO NIGRO, ALEJANDRO LIMACHE, AND EUGENIO OÑATE, *Large time-step explicit integration method for solving problems with dominant convection*, Computer Methods in Applied Mechanics and Engineering, 217-220 (2012), pp. 168–185.
- [11] SERGIO R. IDELSOHN, EUGENIO OÑATE, AND FACUNDO DEL PIN, *The particle finite element method: a powerful tool to solve incompressible flows with free-surfaces and breaking waves*, International Journal for Numerical Methods in Engineering, 61 (2004), pp. 964–989.
- [12] WILLAM KAHAN AND JOSEPH D. DARCY, *How Java’s Floating-Point Hurts Everyone Everywhere*. ACM 1998 Workshop on Java for High-Performance Network Computing. <http://www.cs.berkeley.edu/~wkahan/JAVAhurt.pdf>, 1998. Online; accessed: 05-02-2014.
- [13] PETER KIPFER, FRANK RECK, AND GUNTHER GREINER, *Local Exact Particle Tracing on Unstructured Grids*, Computer Graphics Forum, 22 (2003), pp. 133–142.
- [14] ALLAN CHARLES MCCURDY, KWOK CHOI NG, AND BERESFORD NEILL PARLETT, *Accurate computation of divided differences of the exponential function*, Mathematics of Computation, 43 (1984), pp. 501–501.
- [15] GREGORY M. NIELSON AND IL-HONG JUNG, *Tools for computing tangent curves for linearly varying vector fields over tetrahedral domains*, IEEE Transactions on Visualization and Computer Graphics, 5 (1999), pp. 360–372.
- [16] EUGENIO OÑATE, SERGIO R. IDELSOHN, FACUNDO DEL PIN, AND ROMAIN AUBRY, *The particle finite element method. An overview*, International Journal of Computational Methods, 1 (2004), pp. 267–307.
- [17] ALEXANDER OSTERMANN, MECHTHILD THALHAMMER, AND WILLIAM M. WRIGHT, *A Class of Explicit Exponential General Linear Methods*, BIT Numerical Mathematics, 46 (2006), pp. 409–431.
- [18] BERESFORD NEILL PARLETT, *A recurrence among the elements of functions of triangular matrices*, Linear Algebra and its Applications, 14 (1976), pp. 117–121.
- [19] JAMES F. PRICE, *Lagrangian and Eulerian Representations of Fluid Flow: Kinematics and the Equations of Motion*. Woods Hole Oceanographic Institution. <http://www.whoi.edu/science/P0/people/jprice/class/ELreps.pdf>, 2006. Online; accessed: 10-12-2012.
- [20] ERIC W. WEISSTEIN, *Cubic Formula*. From MathWorld—A Wolfram Web Resource. <http://mathworld.wolfram.com/CubicFormula.html>. Online; accessed: 23-10-2013.