

# Stable and Conservative High-Order Methods on Triangular Elements Using Tensor-Product Summation-by-Parts Operators

Tristan Montoya and David W. Zingg  
Corresponding author: tristan.montoya@mail.utoronto.ca

Institute for Aerospace Studies, University of Toronto  
4925 Dufferin St, Toronto, ON M3H 5T6, Canada

**Abstract:** We present provably stable and conservative discretizations of arbitrary order which exploit the geometric flexibility afforded by triangular elements alongside the computational efficiency of tensor-product operators through the use of a collapsed coordinate transformation. Collocated nodal approximations as well as modal formulations employing orthogonal polynomial expansions are presented, both of which allow for the efficient evaluation of the semi-discrete residual through the application of one-dimensional operations with respect to each component of the collapsed coordinate system. In order to develop such schemes, we first propose an approach for constructing summation-by-parts (SBP) operators of any order of accuracy which are amenable to such decompositions on triangular elements. These operators are then used to construct energy-stable discretizations of the linear advection equation in curvilinear coordinates by way of a skew-symmetric formulation local to each element, with adjacent elements coupled using numerical flux functions. We also prove that the proposed schemes are discretely conservative and free-stream preserving for linear as well as nonlinear problems. Numerical experiments are presented in which the linear advection equation is solved on a sequence of curvilinear meshes, confirming the discrete conservation, energy dissipation (for an upwind numerical flux), and energy conservation (for a central numerical flux) properties which were established theoretically, and demonstrating that the proposed nodal and modal schemes obtain convergence rates and levels of error comparable to those of a standard discontinuous Galerkin method of the same degree. Examining the spectra of the resulting semi-discrete operators, we also verify that the eigenvalues of the proposed schemes all have a non-positive real part for an upwind flux and lie on the imaginary axis for a central flux. This work represents the first application of the SBP methodology to tensor-product discretizations in collapsed coordinates and offers a promising approach to the development of efficient and robust high-order methods for simulations of fluid flow around complex geometries.

*Keywords:* High-order methods, unstructured grids, summation-by-parts, discontinuous Galerkin, spectral-element, tensor-product, energy stability, conservation laws, curvilinear coordinates.

## 1 Introduction

Scale-resolving simulations (i.e. direct numerical simulations and large eddy simulations) of turbulent flows largely rely on the computational efficiency achieved by spectral methods and high-order finite-difference methods when applied to smooth problems on relatively simple geometries. These methods conventionally exploit tensor-product formulations wherein multidimensional approximation procedures are decomposed so as to consist of individual one-dimensional operations, which are in part responsible for their efficiency benefits relative to inherently multidimensional finite-volume or finite-element approaches employing unstructured grids. The application of this approach to spectral methods in multiple space dimensions originates with the

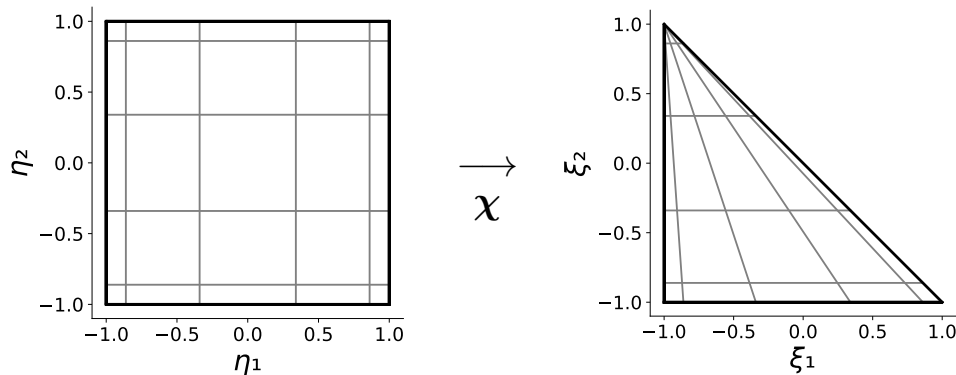


Figure 1: Illustration of the collapsed coordinate transformation (lines of constant  $\eta_1$  and  $\eta_2$  pictured)

work of Orszag [1] and is commonly referred to *sum factorization* in such contexts. Despite these advantages, the generation of curvilinear structured and block-structured grids, which are necessary for the application of traditional spectral and finite-difference methods to complex geometries, remains a major bottleneck for practical flow simulations, motivating the use of multidimensional discretization techniques which are more amenable to general element types including triangles and tetrahedra. Unfortunately, multidimensional discretizations not exploiting tensor-product decompositions generally result in algorithms that impose tighter coupling between numerical degrees of freedom, which, as discussed by Vos *et al.* [2] and Cantwell *et al.* [3], can lead to increased computational expense relative to algorithms employing such decompositions, particularly at higher orders of accuracy.

A promising approach to exploiting the geometric flexibility of simplicial elements while retaining the computational benefits of a tensorial operator structure involves the use of a collapsed coordinate transformation such as that illustrated in Figure 1. Although such transformations have been employed in the context of continuous Galerkin (CG) as well as discontinuous Galerkin (DG) methods for several decades (see, for example, Sherwin and Karniadakis [4] and Kirby *et al.* [5]) and have been shown to result in efficient algorithms employing single-instruction multiple-data (SIMD) vectorization on modern hardware (see, for example, Moxey *et al.* [6]), the theoretical stability and conservation properties of the resulting schemes have received relatively little attention prior to the present work. This particularly calls into question their applicability to curvilinear meshes and nonlinear problems, which are known to present significant challenges to the robustness of any high-order method for which an *a priori* proof of stability is not available.

In recent years, summation-by-parts (SBP) operators have been instrumental in providing a rigorous yet versatile approach to constructing provably stable and conservative high-order methods for linear and nonlinear problems (see, for example, the review papers by Del Rey Fernández *et al.* [7] and Svärd and Nordström [8]) and, as discussed by the authors in [9], have been recognized as providing a unifying framework for the analysis of a wide range of novel as well as existing numerical methods. Although linearly as well as nonlinearly stable SBP discretizations on triangular and tetrahedral elements have recently been developed (see, for example, the contributions in [10–18] as well as the review paper by Chen and Shu [19]), such methods rarely employ polynomial degrees greater than about four or five. Moreover, potential extensions of such schemes to substantially higher orders of accuracy are limited in efficiency relative to tensor-product discretizations as a result of the inherently multidimensional nature of their constituent operators.

With the aim of developing robust numerical methods which extend efficiently to arbitrarily high order on simplicial elements, the present work represents the first application (to the authors’ knowledge) of the SBP approach to discretizations in collapsed coordinates. Restricting our attention to methods based on discontinuous solution spaces due to the advantages of their diagonal or block-diagonal mass matrices in the context of explicit temporal integration, our contributions include the development of efficient nodal (i.e. evolving point values) as well as modal (i.e. evolving orthogonal polynomial expansion coefficients) tensor-product discretizations for conservation laws on triangular elements which are discretely conservative and free-stream preserving for general fluxes and energy stable for the linear advection equation in curvilinear coordinates by way of a skew-symmetric splitting. Notably, as the SBP stability theory does not rely on integral exactness, de-aliasing procedures based on over-integration (e.g. those described in [20–22]),

which are often employed for the *ad hoc* stabilization of conventional DG methods, are not necessary for the proposed schemes. Moreover, the proposed approach extends in a straightforward manner to provably entropy-stable discretizations of nonlinear systems such as the Euler and Navier-Stokes equations using techniques such as those developed in [11] and [16].

The structure of the remainder of this paper is as follows. In §2, we introduce some notational conventions and describe the model problem, mesh, and curvilinear coordinate transformation. In §3, we introduce the fundamental building blocks of the proposed tensor-product discretizations in collapsed coordinates, which we use to construct SBP operators on the triangle which are amenable to sum-factorization algorithms. In §4, we present skew-symmetric nodal and modal formulations employing such operators on curvilinear meshes, which we prove through matrix analysis to be conservative, free-stream preserving, and energy stable. In §5, we use the proposed schemes to solve the linear advection equation in curvilinear coordinates, assessing their accuracy properties relative to a standard DG scheme and confirming their theoretical conservation and energy stability properties. Concluding remarks and directions for future research are provided in §6.

## 2 Preliminaries

### 2.1 Notation

The notation in this paper closely follows that employed in [9]. Symbols bearing single underlines denote vectors (treated as column matrices), while symbols bearing double underlines denote matrices. Symbols in bold such as  $\mathbf{x}$  and  $\mathbf{y}$  are used specifically to denote Cartesian (i.e. spatial) vectors, for which we employ the usual dot product  $\mathbf{x} \cdot \mathbf{y} := x_1 y_1 + \dots + x_d y_d$ , Euclidean norm  $\|\mathbf{x}\|^2 := \mathbf{x} \cdot \mathbf{x}$ , and del operator  $\nabla_{\mathbf{x}} := [\partial/\partial x_1, \dots, \partial/\partial x_d]^T$ . The symbols  $\mathbb{R}$ ,  $\mathbb{R}^+$ ,  $\mathbb{R}_0^+$ ,  $\mathbb{N}$ ,  $\mathbb{N}_0$ , and  $\mathbb{S}^{d-1}$  denote the real numbers, the positive real numbers, the non-negative real numbers, the natural numbers (excluding zero), the natural numbers including zero, and the unit  $(d-1)$ -sphere, which is given by  $\{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\| = 1\}$ . The symbols  $\mathbf{0}^{(N)}$  and  $\mathbf{1}^{(N)}$  are reserved for vectors of length  $N \in \mathbb{N}$  containing all zeros and all ones, respectively, and we use  $\{1 : N\}$  as shorthand for an index set of the form  $\{1, 2, \dots, N\}$ . Furthermore, given any domain  $\mathcal{D} \subset \mathbb{R}^d$ , we use the symbol  $\partial\mathcal{D}$  to denote its boundary and  $\bar{\mathcal{D}} := \mathcal{D} \cup \partial\mathcal{D}$  to denote its closure; the interior of a closed domain  $\mathcal{D}$  is then given by  $\mathring{\mathcal{D}} := \mathcal{D} \setminus \partial\mathcal{D}$ . Other relevant notational conventions are introduced as they appear.

### 2.2 Scalar Conservation Law

We consider as a model problem a first-order conservation law in two space dimensions governing the evolution of a scalar quantity  $U(\mathbf{x}, t) \in \Upsilon \subseteq \mathbb{R}$ , where  $\Upsilon$  denotes the set of admissible solution states. Such partial differential equations (PDEs) take the general form

$$\begin{aligned} \frac{\partial U(\mathbf{x}, t)}{\partial t} + \nabla_{\mathbf{x}} \cdot \mathbf{F}(U(\mathbf{x}, t)) &= 0, & \forall (\mathbf{x}, t) \in \Omega \times (0, T), \\ U(\mathbf{x}, 0) &= U^0(\mathbf{x}), & \forall \mathbf{x} \in \Omega, \end{aligned} \quad (2.1)$$

subject to appropriate boundary conditions, where  $\Omega \subset \mathbb{R}^2$  denotes a fixed domain with a piecewise smooth boundary,  $T \in \mathbb{R}^+$  denotes the final time,  $\mathbf{F}(U(\mathbf{x}, t)) \in \mathbb{R}^2$  denotes the flux vector, and  $U^0(\mathbf{x}) \in \Upsilon$  denotes the initial data. Although the methods discussed in this work readily extend to discretizations of nonlinear problems and systems of equations, our analysis of energy stability is based on the linear advection equation, for which the flux is given by  $\mathbf{F}(U(\mathbf{x}, t)) := \mathbf{a}U(\mathbf{x}, t)$ , where  $\mathbf{a} \in \mathbb{R}^2$  denotes the (constant) advection velocity, and we assume that periodic boundary conditions are imposed in all directions.

### 2.3 Mesh and Curvilinear Coordinate Transformation

We begin our description of the discretization by introducing a mesh  $\mathcal{T}^h := \{\Omega^{(\kappa)}\}_{\kappa=1}^{N_e}$  consisting of  $N_e \in \mathbb{N}$  elements of characteristic size  $h \in \mathbb{R}^+$  with nonempty interiors, satisfying the following assumption.

**Assumption 2.1.** *The elements which constitute the mesh  $\mathcal{T}^h$  satisfy  $\bigcup_{\kappa=1}^{N_e} \Omega^{(\kappa)} = \bar{\Omega}$  and  $\mathring{\Omega}^{(\kappa)} \cap \mathring{\Omega}^{(\nu)} = \emptyset$  for  $\kappa \neq \nu$ . Each element is the image of the reference triangle  $\hat{\mathcal{T}}^2 := \{\boldsymbol{\xi} \in [-1, 1]^2 : \xi_1 + \xi_2 \leq 0\}$  under*

a smooth, time-invariant mapping  $\mathbf{X}^{(\kappa)} : \hat{\mathcal{T}}^2 \rightarrow \Omega^{(\kappa)}$ . The Jacobian of such a mapping is denoted by  $\nabla_{\xi} \mathbf{X}^{(\kappa)}(\xi) \in \mathbb{R}^{2 \times 2}$ , where the determinant  $J^{(\kappa)}(\xi) := \det(\nabla_{\xi} \mathbf{X}^{(\kappa)}(\xi))$  is strictly positive for all  $\xi \in \hat{\mathcal{T}}^2$ .

Using the transformation  $\mathbf{x} = \mathbf{X}^{(\kappa)}(\xi)$  to obtain a formulation of (2.1) in reference coordinates, we may apply the chain and product rules as well as the *metric identities* (see, for example, Thomas and Lombard [23] or Kopriva [24]), which are given for  $n \in \{1, 2\}$  by

$$\nabla_{\xi} \cdot \begin{bmatrix} J^{(\kappa)}(\xi) [(\nabla_{\xi} \mathbf{X}^{(\kappa)}(\xi))^{-1}]_{1,n} \\ J^{(\kappa)}(\xi) [(\nabla_{\xi} \mathbf{X}^{(\kappa)}(\xi))^{-1}]_{2,n} \end{bmatrix} = 0, \quad (2.2)$$

in order to obtain a conservative formulation on the reference element. The transformed PDE is then given by

$$\frac{\partial J^{(\kappa)}(\xi) U(\mathbf{X}^{(\kappa)}(\xi), t)}{\partial t} + \nabla_{\xi} \cdot (J^{(\kappa)}(\xi) (\nabla_{\xi} \mathbf{X}^{(\kappa)}(\xi))^{-1} \mathbf{F}(U(\mathbf{X}^{(\kappa)}(\xi), t))) = 0, \quad (2.3)$$

which has spatially varying coefficients under any mapping which is not affine, even when (2.1) is a constant-coefficient problem. Indexing the edges  $\hat{\mathcal{E}}^{(\zeta)} \subset \partial \hat{\mathcal{T}}^2$  of the reference triangle (shown on the right in Figure 1) in a counter-clockwise order, beginning with the bottom edge, the outward unit normal vectors are given by

$$\mathbf{n}^{(1)} := \begin{bmatrix} 0 \\ -1 \end{bmatrix}, \quad \mathbf{n}^{(2)} := \begin{bmatrix} 1/\sqrt{2} \\ 1/\sqrt{2} \end{bmatrix}, \quad \mathbf{n}^{(3)} := \begin{bmatrix} -1 \\ 0 \end{bmatrix}. \quad (2.4)$$

The outward unit normal to the curved physical facet  $\Gamma^{(\kappa, \zeta)} \subset \partial \Omega^{(\kappa)}$  which is the image of  $\hat{\mathcal{E}}^{(\zeta)}$  under  $\mathbf{X}^{(\kappa)}$  is then given by  $\mathbf{n}^{(\kappa, \zeta)} : \Gamma^{(\kappa, \zeta)} \rightarrow \mathbb{S}^1$  according to Nanson's formula (see, for example, Gurtin *et al.* [25, §8.1] for a derivation in the context of continuum mechanics),

$$J^{(\kappa, \zeta)}(\xi) \mathbf{n}^{(\kappa, \zeta)}(\mathbf{X}^{(\kappa)}(\xi)) = J^{(\kappa)}(\xi) (\nabla_{\xi} \mathbf{X}^{(\kappa)}(\xi))^{-T} \hat{\mathbf{n}}^{(\zeta)}, \quad (2.5)$$

where we define  $J^{(\kappa, \zeta)}(\xi) := \|J^{(\kappa)}(\xi) (\nabla_{\xi} \mathbf{X}^{(\kappa)}(\xi))^{-T} \hat{\mathbf{n}}^{(\zeta)}\|$ . Integrating (2.3) by parts against a continuously differentiable test function  $V(\xi)$  and using (2.5), we then obtain a weak formulation given by

$$\begin{aligned} \int_{\hat{\mathcal{T}}^2} \left( V(\xi) \frac{\partial J^{(\kappa)}(\xi) U(\mathbf{X}^{(\kappa)}(\xi), t)}{\partial t} - \nabla_{\xi} V(\xi) \cdot (J^{(\kappa)}(\xi) (\nabla_{\xi} \mathbf{X}^{(\kappa)}(\xi))^{-1} \mathbf{F}(U(\mathbf{X}^{(\kappa)}(\xi), t))) \right) d\xi \\ + \sum_{\zeta=1}^3 \int_{\hat{\mathcal{E}}^{(\zeta)}} V(\xi) J^{(\kappa, \zeta)}(\xi) (\mathbf{n}^{(\kappa, \zeta)}(\mathbf{X}^{(\kappa)}(\xi)) \cdot \mathbf{F}(U(\mathbf{X}^{(\kappa)}(\xi), t))) d\hat{s} = 0, \end{aligned} \quad (2.6)$$

which is the starting point for the construction of the schemes proposed in this paper.

## 3 Tensor-Product Approximations in Collapsed Coordinates

### 3.1 Collapsed Coordinates

Defining the *collapsed coordinate system*  $\boldsymbol{\eta} \in [-1, 1]^2$ , we may obtain any point  $\xi \in \hat{\mathcal{T}}^2$  by way of the transformation

$$\chi(\boldsymbol{\eta}) := \begin{bmatrix} \frac{1}{2}(1 + \eta_1)(1 - \eta_2) - 1 \\ \eta_2 \end{bmatrix}, \quad (3.1)$$

which is sometimes referred to as the *Duffy transformation* [26], and collapses the top edge of the square onto a single vertex of the triangle, as shown in Figure 1. The Jacobian of such a mapping is given by

$$\nabla_{\boldsymbol{\eta}} \chi(\boldsymbol{\eta}) = \begin{bmatrix} \frac{1}{2}(1 - \eta_2) & -\frac{1}{2}(1 + \eta_1) \\ 0 & 1 \end{bmatrix}, \quad (3.2)$$

where  $\det(\nabla_{\boldsymbol{\eta}} \chi(\boldsymbol{\eta})) = \frac{1}{2}(1 - \eta_2)$  is strictly positive for all  $\boldsymbol{\eta} \in [-1, 1] \times [-1, 1]$ , and is zero for  $\eta_2 = 1$ .

### 3.2 Approximation Spaces

Several finite-dimensional approximation spaces are referenced in this work. The first is the standard *total-degree polynomial space* on the reference element, which is given for  $p \in \mathbb{N}_0$  by

$$\mathbb{P}_p(\hat{\mathcal{T}}^2) := \text{span} \{ \hat{\mathcal{T}}^2 \ni \boldsymbol{\xi} \mapsto \xi_1^{\alpha_1} \xi_2^{\alpha_2} : \boldsymbol{\alpha} \in \mathcal{N}(p) \}, \quad (3.3)$$

where we define the multi-index set  $\mathcal{N}(p) := \{ \boldsymbol{\alpha} \in \mathbb{N}_0^2 : \alpha_1 + \alpha_2 \leq p \}$ . The dimension of the space  $\mathbb{P}_p(\hat{\mathcal{T}}^2)$ , or, equivalently, the cardinality of the set  $\mathcal{N}(p)$ , is given by  $N_p^* := \frac{1}{2}(p+1)(p+2)$ . We also define the (potentially anisotropic) *tensor-product polynomial space* on the square, which is given in terms of  $\mathcal{I}(\mathbf{q}) := \{0 : q_1\} \times \{0 : q_2\}$  for  $\mathbf{q} \in \mathbb{N}_0^2$  as

$$\mathbb{Q}_{\mathbf{q}}([-1, 1]^2) := \text{span} \{ [-1, 1]^2 \ni \boldsymbol{\eta} \mapsto \eta_1^{\alpha_1} \eta_2^{\alpha_2} : \boldsymbol{\alpha} \in \mathcal{I}(\mathbf{q}) \}, \quad (3.4)$$

which is of dimension  $N_{\mathbf{q}} := (q_1+1)(q_2+1)$ . We then have the space of functions which belong to  $\mathbb{Q}_{\mathbf{q}}([-1, 1]^2)$  when expressed under the mapping in (3.1), which is given by

$$\mathbb{R}_{\mathbf{q}}(\hat{\mathcal{T}}^2) := \{ V : \hat{\mathcal{T}}^2 \setminus \{[-1, 1]^T\} \rightarrow \mathbb{R} \mid V \circ \boldsymbol{\chi} \in \mathbb{Q}_{\mathbf{q}}([-1, 1]^2) \}, \quad (3.5)$$

where the domain is restricted in order to exclude the singularity of the transformation  $\boldsymbol{\chi} : [-1, 1]^2 \rightarrow \hat{\mathcal{T}}^2$ .

### 3.3 Quadrature Rules

For each  $m \in \{1, 2\}$ , we introduce  $q_m + 1$  distinct quadrature nodes  $\{\eta_m^{(i)}\}_{i=0}^{q_m}$  on the interval  $[-1, 1]$  and corresponding positive weights  $\{\omega_m^{(i)}\}_{i=0}^{q_m}$ , which allow one to construct an approximation of the form

$$\int_{-1}^1 V(\eta) d\eta \approx \sum_{i=0}^{q_m} V(\eta_m^{(i)}) \omega_m^{(i)}, \quad (3.6)$$

which is said to be of degree  $\tau \in \mathbb{N}_0$  if the above holds as an equality whenever  $V(\eta)$  is a polynomial of at most degree  $\tau$ . Important classes of one-dimensional quadrature rules with positive weights include the Legendre-Gauss (LG), Legendre-Gauss-Radau (LGR), and Legendre-Gauss-Lobatto (LGL) families, which employ nodal sets including zero, one, and two of the interval endpoints, respectively, and are of degree  $2q_m + 1$ ,  $2q_m$ , and  $2q_m - 1$ , respectively (see, for example, Abramowitz and Stegun [27, §25.4]).<sup>1</sup> In order to avoid the singularity of the mapping in (3.1), we make the following assumption, which, importantly, precludes the use of an LGL quadrature rule in the  $\eta_2$  direction.

**Assumption 3.1.** *The one-dimensional quadrature rules in the  $\eta_1$  and  $\eta_2$  directions employ nodal sets satisfying  $-1 \leq \eta_1^{(0)} < \dots < \eta_1^{(q_1)} \leq 1$  and  $-1 \leq \eta_2^{(0)} < \dots < \eta_2^{(q_2)} < 1$ , respectively.*

Using the approximations in (3.6) to integrate with respect to  $\eta_1$  and  $\eta_2$  then allows for the definition of a tensor-product quadrature rule on the triangle with nodes  $\{\boldsymbol{\xi}^{(i)}\}_{i=1}^{N_{\mathbf{q}}}$  and weights  $\{\omega^{(i)}\}_{i=1}^{N_{\mathbf{q}}}$  given by

$$\boldsymbol{\xi}^{(\sigma(\boldsymbol{\alpha}))} := \boldsymbol{\chi}(\eta_1^{(\alpha_1)}, \eta_2^{(\alpha_2)}) \quad \text{and} \quad \omega^{(\sigma(\boldsymbol{\alpha}))} := \frac{1 - \eta_2^{(\alpha_2)}}{2} \omega_1^{(\alpha_1)} \omega_2^{(\alpha_2)}, \quad (3.7)$$

where  $\sigma : \mathcal{I}(\mathbf{q}) \rightarrow \{1 : N_{\mathbf{q}}\}$  associates each multi-index with a unique scalar index. This allows for the approximation of integrals on the reference element as

$$\int_{\hat{\mathcal{T}}^2} V(\boldsymbol{\xi}) d\boldsymbol{\xi} \approx \sum_{i=1}^{N_{\mathbf{q}}} V(\boldsymbol{\xi}^{(i)}) \omega^{(i)}, \quad (3.8)$$

where the multidimensional quadrature weights are strictly positive under Assumption 3.1. Integrals over

<sup>1</sup>Legendre-Gauss-Radau quadrature rules may be defined to include nodes at either the left endpoint ( $\eta = -1$ ) or the right endpoint ( $\eta = 1$ ) of the interval  $\eta \in [-1, 1]$ ; in this paper, the term *LGR quadrature* refers to a rule employing the left endpoint.

each edge (i.e. facet) of the reference triangle may be similarly approximated as

$$\int_{\hat{\mathcal{E}}^{(\zeta)}} V(\boldsymbol{\xi}) d\hat{s} \approx \sum_{i=1}^{N_\zeta} V(\boldsymbol{\xi}^{(\zeta,i)}) \omega^{(\zeta,i)} \quad (3.9)$$

in terms of  $N_\zeta \in \mathbb{N}$  distinct quadrature nodes  $\{\boldsymbol{\xi}^{(\zeta,i)}\}_{i=1}^{N_\zeta}$  on  $\hat{\mathcal{E}}^{(\zeta)} \subset \partial\hat{\mathcal{T}}^2$  and non-negative weights  $\{\omega^{(\zeta,i)}\}_{i=1}^{N_\zeta}$ .

*Remark 3.1.* It is often desirable from an efficiency perspective to align the facet quadrature nodes along lines of volume quadrature nodes in order to make use of efficient one-dimensional interpolation/extrapolation procedures, which can be achieved by taking  $N_1 = q_1 + 1$  and  $N_2 = N_3 = q_2 + 1$ , and defining

$$\begin{aligned} \boldsymbol{\xi}^{(1,\sigma_1(\alpha_1))} &:= \boldsymbol{\chi}(\eta_1^{(\alpha_1)}, -1), & \boldsymbol{\xi}^{(2,\sigma_2(\alpha_2))} &:= \boldsymbol{\chi}(1, \eta_2^{(\alpha_2)}), & \boldsymbol{\xi}^{(3,\sigma_3(\alpha_2))} &:= \boldsymbol{\chi}(-1, \eta_2^{(\alpha_2)}), \\ \omega^{(1,\sigma_1(\alpha_1))} &:= \omega_1^{(\alpha_1)}, & \omega^{(2,\sigma_2(\alpha_2))} &:= \sqrt{2}\omega_2^{(\alpha_2)}, & \omega^{(3,\sigma_3(\alpha_2))} &:= \omega_2^{(\alpha_2)}, \end{aligned} \quad (3.10)$$

where  $\sigma_\zeta : \{1 : N_\zeta\} \rightarrow \{1 : N_\zeta\}$  denotes an arbitrary reordering of the facet quadrature nodes.

We also make the following assumption, which requires the facet quadrature nodes for any two elements sharing an interface to coincide (or in the case of periodic interfaces, to have opposite outward unit normals), and correspond to quadrature weights which are equal when scaled according to the mapping.

**Assumption 3.2.** *The mesh is conforming in the sense that for each pair of indices  $\kappa, \nu \in \{1 : N_e\}$  with  $\kappa \neq \nu$  such that  $\partial\Omega^{(\kappa)} \cap \partial\Omega^{(\nu)} \neq \emptyset$ , there exist  $\zeta, \eta \in \{1 : 3\}$  such that either  $\Gamma^{(\kappa,\zeta)} \subset \partial\Omega^{(\kappa)}$  and  $\Gamma^{(\nu,\eta)} \subset \partial\Omega^{(\nu)}$  are coincident but oppositely oriented such that  $\mathbf{n}^{(\kappa,\zeta)}(\mathbf{x}) = -\mathbf{n}^{(\nu,\eta)}(\mathbf{x})$ , or such facets are connected via periodic boundary conditions. Furthermore, any two abutting facets contain an equal number of quadrature nodes  $N_\zeta = N_\eta$ , where for every  $i \in \{1 : N_\zeta\}$  there exists a unique  $j \in \{1 : N_\zeta\}$  such that*

$$\mathbf{X}^{(\kappa)}(\boldsymbol{\xi}^{(\zeta,i)}) = \mathbf{X}^{(\nu)}(\boldsymbol{\xi}^{(\eta,j)}) \quad \text{and} \quad \omega^{(\zeta,i)} J^{(\kappa,\zeta)}(\boldsymbol{\xi}^{(\zeta,i)}) = \omega^{(\eta,j)} J^{(\nu,\eta)}(\boldsymbol{\xi}^{(\eta,j)}), \quad (3.11)$$

with the first condition replaced by  $\mathbf{n}^{(\kappa,\zeta)}(\mathbf{X}^{(\kappa)}(\boldsymbol{\xi}^{(\zeta,i)})) = -\mathbf{n}^{(\nu,\eta)}(\mathbf{X}^{(\nu)}(\boldsymbol{\xi}^{(\eta,j)}))$  for any periodic interface.

### 3.4 Nodal and Modal Bases

The discretizations described in this paper are constructed using basis functions which are separable in the sense that they consist of products of univariate functions when expressed in collapsed coordinates. The first of such bases which we consider is a nodal basis  $\{\ell^{(i)}\}_{i=1}^{N_q}$  for the space  $\mathbb{R}_q(\hat{\mathcal{T}}^2)$ , which is given by

$$\ell^{(\sigma(\alpha))}(\boldsymbol{\chi}(\boldsymbol{\eta})) := \ell_1^{(\alpha_1)}(\eta_1) \ell_2^{(\alpha_2)}(\eta_2) \quad (3.12)$$

in terms of the one-dimensional Lagrange polynomials  $\{\ell_m^{(i)}\}_{i=0}^{q_m}$  employing the quadrature nodes in (3.6), which satisfy the *cardinal property*  $\ell_m^{(i)}(\eta_m^{(j)}) = \delta_{ij}$  and are given explicitly for  $m \in \{1, 2\}$  as

$$\ell_m^{(i)}(\eta) := \prod_{j=0, j \neq i}^{q_m} \frac{\eta - \eta_m^{(j)}}{\eta_m^{(i)} - \eta_m^{(j)}}. \quad (3.13)$$

The basis in (3.12) then corresponds to a standard tensor-product polynomial approximation in collapsed coordinates, but contains rational functions when the mapping in (3.1) is inverted in order to express the basis functions in terms of coordinates on the triangle.

We also make use of an orthogonal modal basis  $\{\phi^{(i)}\}_{i=1}^{N_p^*}$  for the space  $\mathbb{P}_p(\hat{\mathcal{T}}^2)$  as described by Prorior [28], Koornwinder [29], and Dubiner [30], which is given by

$$\phi^{(\pi(\alpha))}(\boldsymbol{\chi}(\boldsymbol{\eta})) := \underbrace{\sqrt{2}P_{\alpha_1}^{(0,0)}(\eta_1)}_{=: \psi_1^{(\alpha_1)}(\eta_1)} \underbrace{(1 - \eta_2)^{\alpha_1} P_{\alpha_2}^{(2\alpha_1+1,0)}(\eta_2)}_{=: \psi_2^{(\alpha_1, \alpha_2)}(\eta_2)}, \quad (3.14)$$

with  $\pi : \mathcal{N}(p) \rightarrow \{1 : N_p^*\}$  denoting an arbitrary ordering of the multi-index values and  $P_i^{(a,b)}(\eta)$  denoting a

Jacobi polynomial (see, for example, [27, §22.2]), which we have normalized to satisfy

$$\int_{-1}^1 P_i^{(a,b)}(\eta) P_j^{(a,b)}(\eta) (1-\eta)^a (1+\eta)^b d\eta = \delta_{ij}, \quad (3.15)$$

such that the basis in (3.14) is orthonormal with respect to the standard  $L^2$  inner product on  $\hat{\mathcal{T}}^2$ . Similarly to the nodal basis in (3.12), each modal basis function in (3.14) may be decomposed as a product of one-dimensional polynomials  $\psi_1^{(\alpha_1)}(\eta_1)$  and  $\psi_2^{(\alpha_1, \alpha_2)}(\eta_2)$ , each of which is of at most degree  $p$ . Besides leading to efficient matrix-free algorithms for evaluating such basis functions at the nodes of tensor-product quadrature rules described as in §3.3, this allows for any given function  $V \in \mathbb{P}_p(\hat{\mathcal{T}}^2)$  to be represented equivalently in terms of either basis when  $p \leq \min(q_1, q_2)$ , a property which will be important throughout our analysis.

### 3.5 Summation-by-Parts Operators on the Reference Element

Whether or not explicit in their construction, existing energy-stable and entropy-stable high-order methods on simplicial elements typically rely on the multidimensional summation-by-parts property, which was first introduced in the context of nodal methods by Hicken *et al.* [10], who proposed the following definition.

**Definition 3.1** (Nodal SBP operator). *Let  $\mathcal{D} \subset \mathbb{R}^d$  denote a compact, connected domain with an outward unit normal vector given by  $\hat{\mathbf{n}} : \partial\mathcal{D} \rightarrow \mathbb{S}^{d-1}$ . A matrix  $\underline{\underline{D}}^{(m)} \in \mathbb{R}^{N \times N}$  approximating the partial derivative  $\partial/\partial\xi_m$  on  $N \in \mathbb{N}$  distinct nodes  $\{\boldsymbol{\xi}^{(i)}\}_{i=1}^N$  is a nodal SBP operator of (at least) degree  $p \in \mathbb{N}_0$  if it satisfies*

$$\sum_{j=1}^N \underline{\underline{D}}_{ij}^{(m)} V(\boldsymbol{\xi}^{(j)}) = \frac{\partial V}{\partial \xi_m}(\boldsymbol{\xi}^{(i)}), \quad \forall i \in \{1 : N\}, \quad \forall V \in \mathbb{P}_p(\mathcal{D}), \quad (3.16)$$

and may be decomposed in terms of  $\underline{\underline{M}}, \underline{\underline{Q}}^{(m)} \in \mathbb{R}^{N \times N}$  as  $\underline{\underline{D}}^{(m)} = \underline{\underline{M}}^{-1} \underline{\underline{Q}}^{(m)}$  such that the SBP property  $\underline{\underline{Q}}^{(m)} + (\underline{\underline{Q}}^{(m)})^T = \underline{\underline{E}}^{(m)}$  is satisfied, where the matrix  $\underline{\underline{E}}^{(m)} \in \mathbb{R}^{N \times N}$  approximates an integral over  $\partial\mathcal{D}$  as

$$\sum_{i=1}^N \sum_{j=1}^N U(\boldsymbol{\xi}^{(i)}) \underline{\underline{E}}_{ij}^{(m)} V(\boldsymbol{\xi}^{(j)}) = \int_{\partial\mathcal{D}} U(\boldsymbol{\xi}) V(\boldsymbol{\xi}) \hat{\mathbf{n}}_m(\boldsymbol{\xi}) d\hat{\mathbf{s}}, \quad \forall U, V \in \mathbb{P}_p(\mathcal{D}). \quad (3.17)$$

*Remark 3.2.* The degree of an SBP operator is typically defined uniquely as the *maximum* value of  $p$  for which the accuracy conditions in (3.16) are satisfied, where (3.17) holds for  $U, V \in \mathbb{P}_r(\mathcal{D})$  with  $r \geq p$ . We also refer to any SBP operator for which the associated  $\underline{\underline{M}}$  is diagonal as a *diagonal-norm* SBP operator, referring to the role of such a matrix in defining a discrete norm in which energy stability may be proven.

Noting that satisfying (3.16) for a given  $p$  requires at least  $N_p^*$  nodes, which in two dimensions scales as  $O(p^2)$ , the number of floating-point operations required for differentiating at all nodes using the matrix  $\underline{\underline{D}}^{(m)}$  scales as  $O(p^4)$  when the local degrees of freedom are fully coupled, which, to the authors' knowledge, is the case for all existing triangular SBP operators considered prior to the present work (see, for example, [10, 12–14, 18]). If a tensor-product structure is exploited, however, as is commonplace for quadrilateral elements, the number of floating-point operations instead scales as  $O(p^3)$ , resulting in a reduced computational cost for operators of sufficiently high order. The existence of such tensor-product operators with the SBP property on triangular elements for any arbitrary  $p \in \mathbb{N}$  is established with the following lemma, in which a spectral collocation approach is used to construct such operators.

**Lemma 3.1.** *Suppose that the quadrature rules in (3.6) are of at least degree  $2q_1$  and  $2q_2$  in the  $\eta_1$  and  $\eta_2$  directions, respectively,<sup>2</sup> and that volume and facet quadrature rules satisfying Assumption 3.1 are constructed as in (3.7) and (3.10), respectively. Furthermore, let  $\underline{\underline{D}}^{(m)} \in \mathbb{R}^{N_{\mathbf{q}} \times N_{\mathbf{q}}}$ ,  $\underline{\underline{M}} \in \mathbb{R}^{N_{\mathbf{q}} \times N_{\mathbf{q}}}$ ,  $\underline{\underline{R}}^{(\zeta)} \in \mathbb{R}^{N_{\zeta} \times N_{\zeta}}$ , and  $\underline{\underline{B}}^{(\zeta)} \in \mathbb{R}^{N_{\zeta} \times N_{\zeta}}$  denote the derivative, mass, interpolation/extrapolation, and facet quadrature matrices,*

<sup>2</sup>This can be achieved for arbitrary  $q_1, q_2 \in \mathbb{N}$  through the use of LG or LGR quadrature rules.



respectively, with entries given in terms of such quadrature rules and the nodal basis functions in (3.12) by

$$D_{ij}^{(m)} := \frac{\partial \ell^{(j)}}{\partial \xi_m}(\xi^{(i)}), \quad M_{ij} := \omega^{(i)} \delta_{ij}, \quad R_{ij}^{(\zeta)} := \ell^{(j)}(\xi^{(\zeta, i)}), \quad B_{ij}^{(\zeta)} := \omega^{(\zeta, i)} \delta_{ij}. \quad (3.18)$$

Then, for  $m \in \{1, 2\}$ , the matrix  $\underline{\underline{D}}^{(m)}$  defined above is a diagonal-norm SBP operator of at least degree  $p = \min(q_1, q_2)$  on  $\hat{\mathcal{T}}^2$  in the sense of Definition 3.1, where the boundary operators in (3.17) are given by

$$\underline{\underline{E}}^{(m)} := \sum_{\zeta=1}^3 \hat{n}_m^{(\zeta)} (\underline{\underline{R}}^{(\zeta)})^T \underline{\underline{B}}^{(\zeta)} \underline{\underline{R}}^{(\zeta)}. \quad (3.19)$$

*Proof.* Considering the accuracy conditions in (3.16), we recall from §3.4 that for  $p \leq \min(q_1, q_2)$ , any function  $V \in \mathbb{P}_p(\hat{\mathcal{T}}^2)$  may be represented exactly in terms of the nodal basis in (3.12). The partial derivatives of such an expansion are therefore given for  $m \in \{1, 2\}$  by

$$\frac{\partial V(\xi)}{\partial \xi_m} = \sum_{i=1}^{N_q} V(\xi^{(i)}) \frac{\partial \ell^{(i)}(\xi)}{\partial \xi_m}. \quad (3.20)$$

Evaluating the above at each volume quadrature node, we obtain (3.16) directly from the definition of  $\underline{\underline{D}}^{(m)}$  in (3.18). Next, examining the structure of the matrices  $\underline{\underline{Q}}^{(1)} := \underline{\underline{M}} \underline{\underline{D}}^{(1)}$  and  $\underline{\underline{Q}}^{(2)} := \underline{\underline{M}} \underline{\underline{D}}^{(2)}$ , we obtain

$$\begin{aligned} Q_{\sigma(\alpha), \sigma(\beta)}^{(1)} &= \left( \frac{1 - \eta_2^{(\alpha_2)}}{2} \omega_1^{(\alpha_1)} \omega_2^{(\alpha_2)} \right) \left( \frac{2}{1 - \eta_2^{(\alpha_2)}} \frac{d\ell_1^{(\beta_1)}}{d\eta_1} (\eta_1^{(\alpha_1)}) \ell_2^{(\beta_2)} (\eta_2^{(\alpha_2)}) \right) \\ &= \left( \sum_{i=0}^{q_1} \ell_1^{(\alpha_1)} (\eta_1^{(i)}) \frac{d\ell_1^{(\beta_1)}}{d\eta_1} (\eta_1^{(i)}) \omega_1^{(i)} \right) \left( \sum_{j=0}^{q_2} \ell_2^{(\alpha_2)} (\eta_2^{(j)}) \ell_2^{(\beta_2)} (\eta_2^{(j)}) \omega_2^{(j)} \right) \end{aligned} \quad (3.21)$$

and

$$\begin{aligned} Q_{\sigma(\alpha), \sigma(\beta)}^{(2)} &= \left( \frac{1 - \eta_2^{(\alpha_2)}}{2} \omega_1^{(\alpha_1)} \omega_2^{(\alpha_2)} \right) \left( \frac{1 + \eta_1^{(\alpha_1)}}{1 - \eta_2^{(\alpha_2)}} \frac{d\ell_1^{(\beta_1)}}{d\eta_1} (\eta_1^{(\alpha_1)}) \ell_2^{(\beta_2)} (\eta_2^{(\alpha_2)}) + \ell_1^{(\beta_1)} (\eta_1^{(\alpha_1)}) \frac{d\ell_2^{(\beta_2)}}{d\eta_2} (\eta_2^{(\alpha_2)}) \right) \\ &= \left( \sum_{i=0}^{q_1} \frac{1 + \eta_1^{(i)}}{2} \ell_1^{(\alpha_1)} (\eta_1^{(i)}) \frac{d\ell_1^{(\beta_1)}}{d\eta_1} (\eta_1^{(i)}) \omega_1^{(i)} \right) \left( \sum_{j=0}^{q_2} \ell_2^{(\alpha_2)} (\eta_2^{(j)}) \ell_2^{(\beta_2)} (\eta_2^{(j)}) \omega_2^{(j)} \right) \\ &\quad + \left( \sum_{i=0}^{q_1} \ell_1^{(\alpha_1)} (\eta_1^{(i)}) \ell_1^{(\beta_1)} (\eta_1^{(i)}) \omega_1^{(i)} \right) \left( \sum_{j=0}^{q_2} \frac{1 - \eta_2^{(j)}}{2} \ell_2^{(\alpha_2)} (\eta_2^{(j)}) \frac{d\ell_2^{(\beta_2)}}{d\eta_2} (\eta_2^{(j)}) \omega_2^{(j)} \right), \end{aligned} \quad (3.22)$$

where the first equality in each of (3.21) and (3.22) follows from expressing the quadrature weights and basis functions in terms of their one-dimensional factors and applying the chain rule, while the second equality results from the cardinal property of the Lagrange basis. The boundary operators in (3.19) may be decomposed similarly, resulting in

$$E_{\sigma(\alpha), \sigma(\beta)}^{(1)} = \left( \ell_1^{(\alpha_1)}(1) \ell_1^{(\beta_1)}(1) - \ell_1^{(\alpha_1)}(-1) \ell_1^{(\beta_1)}(-1) \right) \left( \sum_{j=0}^{q_2} \ell_2^{(\alpha_2)} (\eta_2^{(j)}) \ell_2^{(\beta_2)} (\eta_2^{(j)}) \omega_2^{(j)} \right) \quad (3.23)$$

and

$$\begin{aligned} E_{\sigma(\alpha), \sigma(\beta)}^{(2)} &= \left( \ell_1^{(\alpha_1)}(1) \ell_1^{(\beta_1)}(1) \right) \left( \sum_{j=0}^{q_2} \ell_2^{(\alpha_2)} (\eta_2^{(j)}) \ell_2^{(\beta_2)} (\eta_2^{(j)}) \omega_2^{(j)} \right) \\ &\quad - \left( \sum_{i=0}^{q_1} \ell_1^{(\alpha_1)} (\eta_1^{(i)}) \ell_1^{(\beta_1)} (\eta_1^{(i)}) \omega_1^{(i)} \right) \left( \ell_2^{(\alpha_2)}(-1) \ell_2^{(\beta_2)}(-1) \right). \end{aligned} \quad (3.24)$$



The SBP property in the  $\xi_1$  direction then follows from the fact that integration by parts may be applied to the first factor on the second line of (3.21) when using any one-dimensional quadrature rule of at least degree  $2q_1 - 1$  to integrate with respect to  $\eta_1$ , where the diagonal mass matrix  $\underline{M}$  is clearly SPD under Assumption 3.1. For the SBP property to hold in the  $\xi_2$  direction, however, the quadratures in (3.22) and (3.24) must be exact in *both* the  $\eta_1$  and  $\eta_2$  coordinates, requiring polynomials of degree  $2q_1$  and  $2q_2$ , respectively, to be integrated exactly. Under such conditions, the application of integration by parts and the product rule results in

$$\begin{aligned}
Q_{\sigma(\alpha), \sigma(\beta)}^{(2)} &= \left( \int_{-1}^1 \frac{1+\eta_1}{2} \ell_1^{(\alpha_1)}(\eta_1) \frac{d\ell_1^{(\beta_1)}(\eta_1)}{d\eta_1} d\eta_1 \right) \left( \int_{-1}^1 \ell_2^{(\alpha_2)}(\eta_2) \ell_2^{(\beta_2)}(\eta_2) d\eta_2 \right) \\
&\quad + \left( \int_{-1}^1 \ell_1^{(\alpha_1)}(\eta_1) \ell_1^{(\beta_1)}(\eta_1) d\eta_1 \right) \left( \int_{-1}^1 \frac{1-\eta_2}{2} \ell_2^{(\alpha_2)}(\eta_2) \frac{d\ell_2^{(\beta_2)}(\eta_2)}{d\eta_2} d\eta_2 \right) \\
&= \left[ \left( \ell_1^{(\alpha_1)}(1) \ell_1^{(\beta_1)}(1) \right) \left( \int_{-1}^1 \ell_2^{(\alpha_2)}(\eta_2) \ell_2^{(\beta_2)}(\eta_2) d\eta_2 \right) \right. \\
&\quad \left. - \left( \int_{-1}^1 \ell_1^{(\alpha_1)}(\eta_1) \ell_1^{(\beta_1)}(\eta_1) d\eta_1 \right) \left( \ell_2^{(\alpha_2)}(-1) \ell_2^{(\beta_2)}(-1) \right) \right] \\
&\quad - \left[ \left( \int_{-1}^1 \frac{1+\eta_1}{2} \frac{d\ell_1^{(\alpha_1)}(\eta_1)}{d\eta_1} \ell_1^{(\beta_1)}(\eta_1) d\eta_1 \right) \left( \int_{-1}^1 \ell_2^{(\alpha_2)}(\eta_2) \ell_2^{(\beta_2)}(\eta_2) d\eta_2 \right) \right. \\
&\quad \left. + \left( \int_{-1}^1 \ell_1^{(\alpha_1)}(\eta_1) \ell_1^{(\beta_1)}(\eta_1) d\eta_1 \right) \left( \int_{-1}^1 \frac{1-\eta_2}{2} \frac{d\ell_2^{(\alpha_2)}(\eta_2)}{d\eta_2} \ell_2^{(\beta_2)}(\eta_2) d\eta_2 \right) \right] \\
&= E_{\sigma(\alpha), \sigma(\beta)}^{(2)} - Q_{\sigma(\beta), \sigma(\alpha)}^{(2)}.
\end{aligned} \tag{3.25}$$

Finally, expanding any two polynomials  $U, V \in \mathbb{P}_p(\hat{\mathcal{T}}^2)$  in terms of the nodal basis in (3.12), the polynomial exactness of the quadrature approximations in (3.23) and (3.24) under the present assumptions results in

$$\begin{aligned}
\sum_{i=1}^{N_q} \sum_{j=1}^{N_q} U(\xi^{(i)}) E_{ij}^{(1)} V(\xi^{(j)}) &= \int_{-1}^1 U(\chi(1, \eta_2)) V(\chi(1, \eta_2)) d\eta_2 - \int_{-1}^1 U(\chi(-1, \eta_2)) V(\chi(-1, \eta_2)) d\eta_2 \\
&= \frac{1}{\sqrt{2}} \int_{\hat{\mathcal{E}}^{(2)}} U(\xi) V(\xi) d\hat{s} - \int_{\hat{\mathcal{E}}^{(3)}} U(\xi) V(\xi) d\hat{s}
\end{aligned} \tag{3.26}$$

and

$$\begin{aligned}
\sum_{i=1}^{N_q} \sum_{j=1}^{N_q} U(\xi^{(i)}) E_{ij}^{(2)} V(\xi^{(j)}) &= \int_{-1}^1 U(\chi(1, \eta_2)) V(\chi(1, \eta_2)) d\eta_2 - \int_{-1}^1 U(\chi(\eta_1, -1)) V(\chi(\eta_1, -1)) d\eta_1 \\
&= \frac{1}{\sqrt{2}} \int_{\hat{\mathcal{E}}^{(2)}} U(\xi) V(\xi) d\hat{s} - \int_{\hat{\mathcal{E}}^{(1)}} U(\xi) V(\xi) d\hat{s}.
\end{aligned} \tag{3.27}$$

Noting that the outward unit normal vectors to the each edge of the reference triangle are given by (2.4), we therefore obtain the accuracy conditions on the boundary operators given in (3.17).  $\square$

*Remark 3.3.* For operators based on spectral collocation in collapsed coordinates such as those constructed in Lemma 3.1, the required quadrature degree in each direction for integration by parts to hold discretely on the reference triangle is one higher than for the square (see, for example, Kopriva and Gassner [31, §3]).

*Remark 3.4.* While the matrix notation employed throughout this paper streamlines the analysis of the proposed schemes through algebraic techniques, our implementation employs a matrix-free strategy for evaluating tensor-product operators through SIMD vectorization in a similar manner to that employed in [6].

## 4 Stable and Conservative Schemes for Curvilinear Meshes

Although the approximations described in the previous section are directly applicable within *any* numerical framework which would otherwise employ non-tensor-product multidimensional SBP operators on triangular elements (including any of the nonlinearly stable formulations reviewed in [19]), we present as a concrete example a particular skew-symmetric formulation which is energy stable for the linear advection equation in curvilinear coordinates, as well as conservative and free-stream preserving for any linear or nonlinear conservation law in the form of (2.1).

### 4.1 Skew-Symmetric Nodal and Modal Formulations

Beginning with the weak formulation in (2.6), we approximate the solution  $U(\mathbf{X}^{(\kappa)}(\boldsymbol{\xi}), t)$  on the reference element by a function  $U^{(h, \kappa)}(\boldsymbol{\xi}, t)$  belonging to either  $\mathbb{P}_p(\hat{\mathcal{T}}^2)$  or  $\mathbb{R}_q(\hat{\mathcal{T}}^2)$  and take test functions  $V(\boldsymbol{\xi})$  belonging to the same space in order to obtain a Galerkin formulation. Furthermore, we approximate the flux vector within the volume integral by its nodal interpolant using the basis in (3.12) as

$$\mathbf{F}(U(\mathbf{X}^{(\kappa)}(\boldsymbol{\xi}), t)) \approx \underbrace{\sum_{i=1}^{N_q} \mathbf{F}(U^{(h, \kappa)}(\boldsymbol{\xi}^{(i)}, t)) \ell^{(i)}(\boldsymbol{\xi})}_{=: \mathbf{F}^{(h, \kappa)}(\boldsymbol{\xi}, t)}, \quad (4.1)$$

where the above holds as an equality for a linear, constant-coefficient problem. We resolve the discontinuity in the global approximation on the element boundary using a numerical flux function  $F^* : \Upsilon \times \Upsilon \times \mathbb{S}^1 \rightarrow \mathbb{R}$  in order to obtain an approximation of the normal trace of the flux on each facet  $\Gamma^{(\kappa, \zeta)} \subset \partial\Omega^{(\kappa)}$  as

$$\mathbf{n}^{(\kappa, \zeta)}(\mathbf{X}^{(\kappa)}(\boldsymbol{\xi})) \cdot \mathbf{F}(U(\mathbf{X}^{(\kappa)}(\boldsymbol{\xi}), t)) \approx \underbrace{F^*(U^{(h, \kappa)}(\boldsymbol{\xi}, t), U^{(+, \kappa, \zeta)}(\boldsymbol{\xi}, t), \mathbf{n}^{(\kappa, \zeta)}(\mathbf{X}^{(\kappa)}(\boldsymbol{\xi})))}_{=: F^{(*, \kappa, \zeta)}(\boldsymbol{\xi}, t)}, \quad (4.2)$$

where the external state (which may correspond either to prescribed boundary data or the numerical solution on another element across an interior or periodic interface) is given by  $U^{(+, \kappa, \zeta)}(\boldsymbol{\xi}, t) \in \mathbb{R}$  and the outward unit normal vector in physical space is obtained using (2.5). We typically require the numerical flux to be consistent and conservative in the following sense.

**Definition 4.1** (Consistent and conservative numerical flux). *A numerical flux  $F^* : \Upsilon \times \Upsilon \times \mathbb{S}^1 \rightarrow \mathbb{R}$  for the scalar conservation law in (2.1) is consistent if it satisfies  $F^*(U, U, \mathbf{n}) = \mathbf{F}(U) \cdot \mathbf{n}$  for all  $U \in \Upsilon$  and  $\mathbf{n} \in \mathbb{S}^1$  and conservative if it satisfies  $F^*(U^-, U^+, \mathbf{n}) = -F^*(U^+, U^-, -\mathbf{n})$  for all  $U^-, U^+ \in \Upsilon$  and  $\mathbf{n} \in \mathbb{S}^1$ .*

Inserting the approximations discussed above into (2.6) and applying the product rule, integration by parts, and the metric identities in (2.2) to half of the second term inside the volume integral, we obtain a skew-symmetric variational formulation given by

$$\begin{aligned} \int_{\hat{\mathcal{T}}^2} \left( V(\boldsymbol{\xi}) \frac{\partial J^{(\kappa)}(\boldsymbol{\xi}) U^{(h, \kappa)}(\boldsymbol{\xi}, t)}{\partial t} - \frac{1}{2} \sum_{m=1}^2 \sum_{n=1}^2 \frac{\partial V(\boldsymbol{\xi})}{\partial \xi_m} [J^{(\kappa)}(\boldsymbol{\xi}) (\nabla_{\boldsymbol{\xi}} \mathbf{X}^{(\kappa)}(\boldsymbol{\xi}))^{-1}]_{mn} F_n^{(h, \kappa)}(\boldsymbol{\xi}, t) \right. \\ \left. + \frac{1}{2} \sum_{m=1}^2 \sum_{n=1}^2 V(\boldsymbol{\xi}) [J^{(\kappa)}(\boldsymbol{\xi}) (\nabla_{\boldsymbol{\xi}} \mathbf{X}^{(\kappa)}(\boldsymbol{\xi}))^{-1}]_{mn} \frac{\partial F_n^{(h, \kappa)}(\boldsymbol{\xi}, t)}{\partial \xi_m} \right) d\boldsymbol{\xi} \\ + \sum_{\zeta=1}^3 \int_{\hat{\mathcal{E}}^{(\zeta)}} V(\boldsymbol{\xi}) \left( J^{(\kappa, \zeta)}(\boldsymbol{\xi}) F^{(*, \kappa, \zeta)}(\boldsymbol{\xi}, t) - \frac{1}{2} \sum_{n=1}^2 J^{(\kappa, \zeta)}(\boldsymbol{\xi}) n_n^{(\kappa, \zeta)}(\mathbf{X}^{(\kappa)}(\boldsymbol{\xi})) F_n^{(h, \kappa)}(\boldsymbol{\xi}, t) \right) d\hat{s} = 0, \end{aligned} \quad (4.3)$$

where the semi-discrete formulations which we propose are completed by specifying a basis for  $U(\mathbf{X}^{(\kappa)}(\boldsymbol{\xi}), t)$  and  $V(\boldsymbol{\xi})$  as well as quadrature rules to approximate the above volume and facet integrals.

#### 4.1.1 Nodal Formulation

In order to obtain a nodal approximation of (4.3), the numerical solution  $U^{(h,\kappa)}(\cdot, t) \in \mathbb{R}_{\mathbf{q}}(\hat{\mathcal{T}}^2)$  is expanded in terms of the basis in (3.12) and the nodal expansion coefficients  $\underline{u}^{(h,\kappa)}(t) \in \mathbb{R}^{N_{\mathbf{q}}}$  as

$$U^{(h,\kappa)}(\boldsymbol{\xi}, t) = \sum_{i=1}^{N_{\mathbf{q}}} u_i^{(h,\kappa)}(t) \ell^{(i)}(\boldsymbol{\xi}), \quad (4.4)$$

where  $u_i^{(h,\kappa)}(t) = U^{(h,\kappa)}(\boldsymbol{\xi}^{(i)}, t)$ . We then use the quadrature rules in (3.7) and (3.9) to approximate the integrals, resulting in a *skew-symmetric nodal formulation* given for  $\kappa \in \{1 : N_e\}$  and  $t \in (0, T)$  by

$$\begin{aligned} \underline{\underline{M}} J^{(\kappa)} \frac{d\underline{u}^{(h,\kappa)}(t)}{dt} = & \frac{1}{2} \sum_{m=1}^2 \sum_{n=1}^2 (\underline{\underline{D}}^{(m)})^T \underline{\underline{M}} \underline{\underline{A}}^{(\kappa,m,n)} \underline{f}^{(h,\kappa,n)}(t) - \frac{1}{2} \sum_{m=1}^2 \sum_{n=1}^2 \underline{\underline{A}}^{(\kappa,m,n)} \underline{\underline{M}} \underline{\underline{D}}^{(m)} \underline{f}^{(h,\kappa,n)} \\ & - \sum_{\zeta=1}^3 (\underline{\underline{R}}^{(\zeta)})^T \underline{\underline{B}}^{(\zeta)} \left( \underline{J}^{(\kappa,\zeta)} \underline{f}^{(*,\kappa,\zeta)}(t) - \frac{1}{2} \sum_{n=1}^2 \underline{\underline{N}}^{(\kappa,\zeta,n)} \underline{\underline{R}}^{(\zeta)} \underline{f}^{(h,\kappa,n)}(t) \right), \end{aligned} \quad (4.5)$$

where the diagonal matrices  $\underline{\underline{J}}^{(\kappa)} \in \mathbb{R}^{N_{\mathbf{q}} \times N_{\mathbf{q}}}$ ,  $\underline{\underline{A}}^{(\kappa,m,n)} \in \mathbb{R}^{N_{\mathbf{q}} \times N_{\mathbf{q}}}$ ,  $\underline{J}^{(\kappa,\zeta)} \in \mathbb{R}^{N_{\zeta} \times N_{\zeta}}$ , and  $\underline{\underline{N}}^{(\kappa,\zeta)} \in \mathbb{R}^{N_{\zeta} \times N_{\zeta}}$  have entries given by

$$\begin{aligned} J_{ij}^{(\kappa)} &:= J^{(\kappa)}(\boldsymbol{\xi}^{(i)}) \delta_{ij}, & A_{ij}^{(\kappa,m,n)} &:= [J^{(\kappa)}(\boldsymbol{\xi}^{(i)}) (\nabla_{\boldsymbol{\xi}} \mathbf{X}^{(\kappa)}(\boldsymbol{\xi}^{(i)}))^{-1}]_{mn} \delta_{ij} \\ J_{ij}^{(\kappa,\zeta)} &:= J^{(\kappa,\zeta)}(\boldsymbol{\xi}^{(\zeta,i)}) \delta_{ij}, & N_{ij}^{(\kappa,\zeta,n)} &:= J^{(\kappa,\zeta)}(\boldsymbol{\xi}^{(\zeta,i)}) \mathbf{n}^{(\kappa,\zeta)}(\mathbf{X}^{(\kappa)}(\boldsymbol{\xi}^{(\zeta,i)})) \delta_{ij}, \end{aligned} \quad (4.6)$$

and the vectors containing the nodal values of the physical and numerical flux functions are given by

$$\underline{f}^{(h,\kappa,n)}(t) := \begin{bmatrix} F_n^{(h,\kappa)}(\boldsymbol{\xi}^{(1)}, t) \\ \vdots \\ F_n^{(h,\kappa)}(\boldsymbol{\xi}^{(N_{\mathbf{q}})}, t) \end{bmatrix}, \quad \underline{f}^{(*,\kappa,\zeta)}(t) := \begin{bmatrix} F^{(*,\kappa,\zeta)}(\boldsymbol{\xi}^{(\zeta,1)}, t) \\ \vdots \\ F^{(*,\kappa,\zeta)}(\boldsymbol{\xi}^{(\zeta,N_{\zeta})}, t) \end{bmatrix}. \quad (4.7)$$

#### 4.1.2 Modal Formulation

As an alternative to the nodal approach described above, one may choose instead to seek a numerical solution  $U^{(h,\kappa)}(\cdot, t) \in \mathbb{P}_p(\hat{\mathcal{T}}^2)$  represented in terms of the modal basis in (3.14) as

$$U^{(h,\kappa)}(\boldsymbol{\xi}, t) = \sum_{i=1}^{N_p^*} \tilde{u}_i^{(h,\kappa)}(t) \phi^{(i)}(\boldsymbol{\xi}), \quad (4.8)$$

where the vector  $\tilde{\underline{u}}^{(h,\kappa)}(t) \in \mathbb{R}^{N_p^*}$  contains the corresponding expansion coefficients. Evaluating the nodal solution vector in terms of the modal expansion in (4.8) as  $\underline{u}^{(h,\kappa)}(t) = \underline{\underline{V}} \tilde{\underline{u}}^{(h,\kappa)}(t)$  and pre-multiplying both sides of (4.5) by  $\underline{\underline{V}}^T$  in order to restrict the test space to  $\mathbb{P}_p(\hat{\mathcal{T}}^2)$ , we obtain a *skew-symmetric modal formulation* given for  $\kappa \in \{1 : N_e\}$  and  $t \in (0, T)$  by

$$\underline{\underline{V}}^T \underline{\underline{M}} J^{(\kappa)} \underline{\underline{V}} \frac{d\tilde{\underline{u}}^{(h,\kappa)}(t)}{dt} = \underline{\underline{V}}^T \underline{r}^{(h,\kappa)}(t), \quad (4.9)$$

where  $\underline{r}^{(h,\kappa)}(t) \in \mathbb{R}^{N_{\mathbf{q}}}$  denotes the right-hand side of the nodal formulation in (4.5). In order to evaluate the initial condition, we approximate the  $L^2$  projection on the physical element for  $\kappa \in \{1 : N_e\}$  as

$$\underline{\underline{V}}^T \underline{\underline{M}} J^{(\kappa)} \underline{\underline{V}} \tilde{\underline{u}}^{(h,\kappa)}(0) = \underline{\underline{V}}^T \underline{\underline{W}} J^{(\kappa)} \begin{bmatrix} U^0(\mathbf{X}^{(\kappa)}(\boldsymbol{\xi}^{(1)})) \\ \vdots \\ U^0(\mathbf{X}^{(\kappa)}(\boldsymbol{\xi}^{(N_{\mathbf{q}})})) \end{bmatrix}. \quad (4.10)$$

The following lemma establishes that the local mass matrix is SPD, ensuring that (4.10) has a unique solution and allowing for the time derivative in (4.9) to be obtained, for example, using a preconditioned conjugate-gradient method (see, for example, Pazner and Persson [32, §3.3]).

**Lemma 4.1.** *Considering tensor-product quadrature rules and modal basis functions given as in (3.7) and (3.14), respectively, where it is assumed that  $p \leq \min(q_1, q_2)$ , the modal mass matrix  $\underline{\underline{V}}^T \underline{\underline{M}} \underline{\underline{J}}^{(\kappa)} \underline{\underline{V}}$  is SPD for all  $\kappa \in \{1 : N_e\}$  under Assumptions 2.1 and 3.1.*

*Proof.* First, we note that since any polynomial  $V \in \mathbb{P}_p(\hat{T}^2)$  admits a unique expansion in terms of the nodal basis in (3.12) when  $p \leq \min(q_1, q_2)$ , the matrix  $\underline{\underline{V}}$ , which defines an injective mapping from modal expansion coefficients to nodal expansion coefficients, is of full column rank and thus has a nullspace containing only the zero vector. The positive-definiteness of  $\underline{\underline{V}}^T \underline{\underline{M}} \underline{\underline{J}}^{(\kappa)} \underline{\underline{V}}$  then follows from the fact that  $\underline{\underline{J}}^{(\kappa)}$  and  $\underline{\underline{M}}$  are diagonal and positive-definite under Assumptions 2.1 and 3.1, respectively.  $\square$

*Remark 4.1.* Due to the separability of the basis functions in (3.14) and the fact that the quadrature nodes are based on a Cartesian product in collapsed coordinates, the action of  $\underline{\underline{V}}$  or its transpose on a vector can be computed using one-dimensional operations along lines of nodes, employing efficient matrix-free algorithms such as those described in [6]. As such, there is no need to store dense operator matrices for the reference element nor for the physical element, and  $O(p^4)$  procedures are entirely avoided.

## 4.2 Summation-by-Parts Operators on the Physical Element

Although not immediately obvious, the discretization of the skew-symmetric variational formulation in (4.3) is mathematically equivalent to the construction of an SBP operator on each physical element following the approach described by Crean *et al.* [11, §5]. To see this, we may group together all matrices pre-multiplying the vector  $\underline{f}^{(h, \kappa, n)}(t)$  in the nodal formulation given by (4.5), allowing for such a scheme to be expressed as

$$\underline{\underline{M}} \underline{\underline{J}}^{(\kappa)} \frac{d\underline{u}^{(h, \kappa)}(t)}{dt} = \sum_{n=1}^2 (\underline{\underline{Q}}^{(\kappa, n)})^T \underline{f}^{(h, \kappa, n)}(t) - \sum_{\zeta=1}^3 (\underline{\underline{R}}^{(\zeta)})^T \underline{\underline{B}}^{(\zeta)} \underline{\underline{J}}^{(\kappa, \zeta)} \underline{f}^{(*, \kappa, \zeta)}(t), \quad (4.11)$$

for  $\kappa \in \{1 : N_e\}$  and  $t \in (0, T)$ , where we define

$$\underline{\underline{Q}}^{(\kappa, n)} := \frac{1}{2} \sum_{m=1}^2 \left( \underline{\underline{A}}^{(\kappa, m, n)} \underline{\underline{M}} \underline{\underline{D}}^{(m)} - (\underline{\underline{D}}^{(m)})^T \underline{\underline{M}} \underline{\underline{A}}^{(\kappa, m, n)} \right) + \frac{1}{2} \sum_{\zeta=1}^3 (\underline{\underline{R}}^{(\zeta)})^T \underline{\underline{B}}^{(\zeta)} \underline{\underline{N}}^{(\kappa, \zeta, n)} \underline{\underline{R}}^{(\zeta)}, \quad (4.12)$$

which is identical to the operator proposed in [11, §5]. As a result of the skew-symmetry of the first term in (4.12) and the symmetry of the second term, the SBP property is satisfied on the physical element by construction for  $n \in \{1, 2\}$ , as given by

$$\underline{\underline{Q}}^{(\kappa, n)} + (\underline{\underline{Q}}^{(\kappa, n)})^T = \sum_{\zeta=1}^3 (\underline{\underline{R}}^{(\zeta)})^T \underline{\underline{B}}^{(\zeta)} \underline{\underline{N}}^{(\kappa, \zeta, n)} \underline{\underline{R}}^{(\zeta)}. \quad (4.13)$$

It is then straightforward to show that the weak formulation in (4.11) is equivalent to a semi-discrete strong formulation given for  $\kappa \in \{1 : N_e\}$  and  $t \in (0, T)$  by

$$\begin{aligned} \frac{d\underline{u}^{(h, \kappa)}(t)}{dt} = & - \sum_{n=1}^2 (\underline{\underline{M}} \underline{\underline{J}}^{(\kappa)})^{-1} \underline{\underline{Q}}^{(\kappa, n)} \underline{f}^{(h, \kappa, n)}(t) \\ & - \sum_{\zeta=1}^3 (\underline{\underline{M}} \underline{\underline{J}}^{(\kappa)})^{-1} (\underline{\underline{R}}^{(\zeta)})^T \underline{\underline{B}}^{(\zeta)} \left( \underline{\underline{J}}^{(\kappa, \zeta)} \underline{f}^{(*, \kappa, \zeta)}(t) - \sum_{n=1}^2 \underline{\underline{N}}^{(\kappa, \zeta, n)} \underline{\underline{R}}^{(\zeta)} \underline{f}^{(h, \kappa, n)}(t) \right), \end{aligned} \quad (4.14)$$

where the first term on the right-hand side is at least an  $O(h^{\min(q_1, q_2)})$  approximation to the divergence of the flux in (2.1) provided that the solution and flux are sufficiently smooth and that the SBP property is satisfied on the reference element [11, Theorem 9]. Due to the polynomial exactness of the extrapolation

operators, which follows from the fact that any function in  $\mathbb{P}_p(\hat{\mathcal{T}}^2)$  can be represented exactly in terms of the basis in (3.12), we also see that the second term on the right-hand side of (4.14) provides a design-order penalty (referred to as a *simultaneous approximation term* in the SBP literature, following Carpenter *et al.* [33]) if the numerical flux is consistent in the sense of Definition 4.1. Furthermore, we make the following assumption, which requires that constants are differentiated exactly in physical space.

**Assumption 4.1.** *The matrices in (4.12) satisfy  $\underline{\underline{Q}}^{(\kappa,n)} \underline{\underline{1}}^{(N_q)} = \underline{\underline{0}}^{(N_q)}$  for all  $n \in \{1, 2\}$  and  $\kappa \in \{1 : N_e\}$ .*

*Remark 4.2.* It was shown in [11, Theorem 6] that the above assumption is satisfied in two dimensions when the matrices  $\underline{\underline{D}}^{(m)}$  are SBP operators of degree  $p$  on the reference element and the degree of the mapping from reference to physical coordinates is less than or equal to  $p + 1$ . If this is not the case, the metric terms in (4.6) may be replaced by approximations which do satisfy Assumption 4.1, as in [11, §5.4] or [24, §7].

### 4.3 Theoretical Analysis

We now analyze the proposed nodal and modal schemes in (4.5) and (4.9), respectively, in terms of their conservation, free-stream preservation, and energy stability properties, where for simplicity we assume periodic boundary conditions throughout the analysis.

#### 4.3.1 Conservation

Beginning with conservation, the following theorem establishes that the proposed methods are locally conservative (in an element-wise sense) as well as globally conservative for a suitable choice of numerical flux.

**Theorem 4.1** (Conservation). *Under Assumption 4.1, the nodal and modal formulations given by (4.5) and (4.9), respectively, are locally conservative with respect to the schemes' quadrature rules, satisfying*

$$\frac{d}{dt} (\underline{\underline{1}}^{(N_q)})^T \underline{\underline{M}} \underline{\underline{J}}^{(\kappa)} \underline{\underline{u}}^{(h,\kappa)}(t) = - \sum_{\zeta=1}^3 (\underline{\underline{1}}^{(N_\zeta)})^T \underline{\underline{B}}^{(\zeta)} \underline{\underline{J}}^{(\kappa,\zeta)} \underline{\underline{f}}^{(*,\kappa,\zeta)}(t) \quad (4.15)$$

for all  $\kappa \in \{1 : N_e\}$  and  $t \in (0, T)$ . Moreover, for any numerical flux which is conservative in the sense of Definition 4.1, such schemes are globally conservative with the addition of Assumption 3.2, satisfying

$$\frac{d}{dt} \sum_{\kappa=1}^{N_e} (\underline{\underline{1}}^{(N_q)})^T \underline{\underline{M}} \underline{\underline{J}}^{(\kappa)} \underline{\underline{u}}^{(h,\kappa)}(t) = 0 \quad (4.16)$$

for all  $t \in (0, T)$  in the case of periodic boundary conditions.

*Proof.* To establish local conservation for the nodal formulation, we begin with the equivalent physical-operator formulation in (4.11) and pre-multiply both sides by  $(\underline{\underline{1}}^{(N_q)})^T$  in order to obtain

$$\frac{d}{dt} (\underline{\underline{1}}^{(N_q)})^T \underline{\underline{M}} \underline{\underline{J}}^{(\kappa)} \underline{\underline{u}}^{(h,\kappa)}(t) = \sum_{n=1}^2 (\underline{\underline{1}}^{(N_q)})^T (\underline{\underline{Q}}^{(\kappa,n)})^T \underline{\underline{f}}^{(h,\kappa,n)}(t) - \sum_{\zeta=1}^3 (\underline{\underline{1}}^{(N_\zeta)})^T \underline{\underline{B}}^{(\zeta)} \underline{\underline{J}}^{(\kappa,\zeta)} \underline{\underline{f}}^{(*,\kappa,\zeta)}(t), \quad (4.17)$$

where we have brought the term on the left-hand side inside of the time derivative and used the fact that  $\underline{\underline{R}}^{(\zeta)} \underline{\underline{1}}^{(N_q)} = \underline{\underline{1}}^{(N_\zeta)}$  holds for all  $\zeta \in \{1 : 3\}$  as the interpolation/extrapolation operators are exact for constant functions. The first term on the right-hand side then vanishes as a consequence of Assumption 4.1, resulting in the statement of local conservation in (4.15). To demonstrate that the methods are globally conservative, we sum (4.15) over all  $\kappa \in \{1 : N_e\}$ , and note that, as shown in the proof of global conservation in [9, Theorem 4.3], Assumption 3.2 leads to a cancellation of the interface terms due to the conservation property of the numerical flux, resulting in (4.16).

In the case of the modal formulation, we first note that since the approximation space in (3.3) includes constant functions, there exists a vector  $\underline{\underline{v}} \in \mathbb{R}^{N_p^*}$  of expansion coefficients such that  $\underline{\underline{V}} \underline{\underline{v}} = \underline{\underline{1}}^{(N_q)}$ . Pre-multiplying (4.9) by the transpose of such a vector and using  $\underline{\underline{u}}^{(h,\kappa)}(t) = \underline{\underline{V}} \underline{\underline{u}}^{(h,\kappa)}(t)$  to obtain the nodal

values of the numerical solution then results in (4.17). The remainder of the proof is therefore the same as for the nodal formulation, and hence the local and global conservation properties in (4.15) and (4.16) are satisfied for both discretization approaches.  $\square$

#### 4.3.2 Free-Stream Preservation

Particularly in the context of computational fluid dynamics, it is desirable for any constant solution to remain constant for all time, a property commonly referred to as *free-stream preservation*. This is established for the proposed discretizations with the following theorem.

**Theorem 4.2** (Free-stream preservation). *Suppose that the numerical flux function is consistent in the sense of Definition 4.1 and that Assumption 4.1 holds. Assuming periodic boundary conditions, the schemes in (4.5) and (4.9) are then free-stream preserving, satisfying*

$$\frac{d\mathbf{u}^{(h,\kappa)}(t)}{dt} = \mathbf{0}^{(N_q)} \quad (4.18)$$

for any solution given by  $\mathbf{u}^{(h,\kappa)}(t) = U\mathbf{1}^{(N_q)}$ , where  $U \in \Upsilon$  is constant for all  $\kappa \in \{1 : N_e\}$ .

*Proof.* Substituting a constant solution into the strong formulation in (4.14), which may be obtained from (4.5) by applying the SBP property in (4.13) to the physical-operator formulation in (4.13), we obtain

$$\begin{aligned} \frac{d\mathbf{u}^{(h,\kappa)}(t)}{dt} = & - \sum_{n=1}^2 (\underline{\underline{M}}J^{(\kappa)})^{-1} \underline{\underline{Q}}^{(\kappa,n)} \mathbf{1}^{(N_q)} F_n(U) \\ & - \sum_{\zeta=1}^3 (\underline{\underline{M}}J^{(\kappa)})^{-1} (\underline{\underline{R}}^{(\zeta)})^T \underline{\underline{B}}^{(\zeta)} \left( \underline{\underline{J}}^{(\kappa,\zeta)} \underline{\underline{f}}^{(*,\kappa,\zeta)}(t) - \sum_{n=1}^2 \underline{\underline{N}}^{(\kappa,\zeta,n)} \underline{\underline{R}}^{(\zeta)} \mathbf{1}^{(N_q)} F_n(U) \right). \end{aligned} \quad (4.19)$$

The first term on the right-hand side then vanishes under Assumption 4.1. Using the consistency of the numerical flux as well as  $\underline{\underline{R}}^{(\zeta)} \mathbf{1}^{(N_q)} = \mathbf{1}^{(N_\zeta)}$ , the second term on the right-hand side of (4.19) vanishes as well, resulting in free-stream preservation for the nodal scheme in (4.5). Since the nodal values of the semi-discrete residual for the modal scheme in (4.9) are related to those of the nodal scheme through a time-independent linear mapping (i.e. a discrete  $L^2$  projection), which is given by

$$\left. \frac{d\mathbf{u}^{(h,\kappa)}(t)}{dt} \right|_{\text{nodal}} = \underline{\underline{V}} (\underline{\underline{V}}^T \underline{\underline{M}}J^{(\kappa)} \underline{\underline{V}})^{-1} \underline{\underline{V}}^T \underline{\underline{M}}J^{(\kappa)} \left. \frac{d\mathbf{u}^{(h,\kappa)}(t)}{dt} \right|_{\text{nodal}} \quad (4.20)$$

for all  $t \in (0, T)$ , the time derivative on the left-hand side is zero whenever the nodal solution is constant in time, implying that the modal formulation is also free-stream preserving in the sense of (4.18).  $\square$

#### 4.3.3 Energy Stability

We now demonstrate that the proposed discretizations are energy stable for the constant-coefficient linear advection equation in curvilinear coordinates with respect to the discrete  $L^2$  norm associated with the volume quadrature of the scheme. While the following theorem establishes such a result for periodic boundary conditions, the extension to inflow/outflow boundary conditions is straightforward.

**Theorem 4.3** (Energy stability). *The nodal and modal formulations in (4.5) and (4.9), respectively, are energy stable under Assumptions 2.1, 3.1, and 3.2 when applied to a periodic constant-coefficient linear advection problem, provided that the numerical flux takes the form*

$$F^*(U^-, U^+, \mathbf{n}) := \frac{1}{2} \mathbf{a} \cdot \mathbf{n} (U^- + U^+) - \frac{\alpha}{2} |\mathbf{a} \cdot \mathbf{n}| (U^+ - U^-), \quad (4.21)$$

where  $\alpha \in \mathbb{R}_0^+$  takes the same value on either side of any interface. This is to say that such schemes satisfy

$$\frac{1}{2} \frac{d}{dt} \sum_{\kappa=1}^{N_e} (\underline{u}^{(h,\kappa)}(t))^T \underline{\underline{M}} \underline{\underline{J}}^{(\kappa)} \underline{u}^{(h,\kappa)}(t) \leq 0, \quad (4.22)$$

for all  $t \in (0, T)$ , which becomes an equality when a central numerical flux (i.e.  $\alpha = 0$ ) is used at all interfaces.

*Proof.* Pre-multiplying both sides of (4.5) by  $(\underline{u}^{(h,\kappa)}(t))^T$  or both sides of (4.9) by  $(\underline{\underline{u}}^{(h,\kappa)}(t))^T$  and applying the chain rule to the time derivative on the left-hand side, we note that the volume terms vanish due to the skew-symmetry of the matrix  $(\underline{\underline{D}}^{(m)})^T \underline{\underline{M}} \underline{\underline{A}}^{(\kappa,m,n)} - \underline{\underline{A}}^{(\kappa,m,n)} \underline{\underline{M}} \underline{\underline{D}}^{(m)}$  and therefore obtain

$$\begin{aligned} \frac{1}{2} \frac{d}{dt} (\underline{u}^{(h,\kappa)}(t))^T \underline{\underline{M}} \underline{\underline{J}}^{(\kappa)} \underline{u}^{(h,\kappa)}(t) = & - \sum_{\zeta=1}^3 \left( (\underline{u}^{(h,\kappa)}(t))^T (\underline{\underline{R}}^{(\zeta)})^T \underline{\underline{B}}^{(\zeta)} \underline{\underline{J}}^{(\kappa,\zeta)} \underline{\underline{f}}^{(*,\kappa,\zeta)}(t) \right. \\ & \left. - \frac{1}{2} (\underline{u}^{(h,\kappa)}(t))^T (\underline{\underline{R}}^{(\zeta)})^T \underline{\underline{B}}^{(\zeta)} \underline{\underline{J}}^{(\kappa,\zeta)} \underline{\underline{A}}^{(\kappa,\zeta)} \underline{\underline{R}}^{(\zeta)} \underline{u}^{(h,\kappa)}(t) \right), \end{aligned} \quad (4.23)$$

where we define the diagonal matrices  $\underline{\underline{A}}^{(\kappa,\zeta)} \in \mathbb{R}^{N_\zeta \times N_\zeta}$  with entries  $A_{ij}^{(\kappa,\zeta)} := \mathbf{a} \cdot \mathbf{n}^{(\kappa,\zeta)}(\mathbf{X}^{(\kappa)}(\boldsymbol{\xi}^{(\zeta,i)})) \delta_{ij}$ . Considering any two elements  $\Omega^{(\kappa)}, \Omega^{(\nu)} \in \mathcal{T}^h$  for which the facets  $\Gamma^{(\kappa,\zeta)} \subset \partial\Omega^{(\kappa)}$  and  $\Gamma^{(\nu,\eta)} \subset \partial\Omega^{(\nu)}$  are co-incident but oppositely oriented, or otherwise connected via periodic boundary conditions, the corresponding contributions to the energy balance resulting from a numerical flux in the form of (4.21) are given by

$$\begin{aligned} (\underline{u}^{(h,\kappa)}(t))^T (\underline{\underline{R}}^{(\zeta)})^T \underline{\underline{B}}^{(\zeta)} \underline{\underline{J}}^{(\kappa,\zeta)} \underline{\underline{f}}^{(*,\kappa,\zeta)}(t) = & \frac{1}{2} (\underline{u}^{(h,\kappa)}(t))^T (\underline{\underline{R}}^{(\zeta)})^T \underline{\underline{B}}^{(\zeta)} \underline{\underline{J}}^{(\kappa,\zeta)} \underline{\underline{A}}^{(\kappa,\zeta)} \underline{\underline{R}}^{(\zeta)} \underline{u}^{(h,\kappa)}(t) \\ & + \frac{1}{2} (\underline{u}^{(h,\kappa)}(t))^T (\underline{\underline{R}}^{(\zeta)})^T \underline{\underline{B}}^{(\zeta)} \underline{\underline{J}}^{(\kappa,\zeta)} \underline{\underline{A}}^{(\kappa,\zeta)} \underline{\underline{T}}^{(\kappa,\zeta)} \underline{\underline{R}}^{(\eta)} \underline{u}^{(h,\nu)}(t) \\ & - \frac{\alpha}{2} (\underline{u}^{(h,\kappa)}(t))^T (\underline{\underline{R}}^{(\zeta)})^T \underline{\underline{B}}^{(\zeta)} \underline{\underline{J}}^{(\kappa,\zeta)} |\underline{\underline{A}}^{(\kappa,\zeta)}| \underline{\underline{T}}^{(\kappa,\zeta)} \underline{\underline{R}}^{(\eta)} \underline{u}^{(h,\nu)}(t) \\ & + \frac{\alpha}{2} (\underline{u}^{(h,\kappa)}(t))^T (\underline{\underline{R}}^{(\zeta)})^T \underline{\underline{B}}^{(\zeta)} \underline{\underline{J}}^{(\kappa,\zeta)} |\underline{\underline{A}}^{(\kappa,\zeta)}| \underline{\underline{R}}^{(\zeta)} \underline{u}^{(h,\kappa)}(t) \end{aligned} \quad (4.24)$$

and

$$\begin{aligned} (\underline{u}^{(h,\nu)}(t))^T (\underline{\underline{R}}^{(\eta)})^T \underline{\underline{B}}^{(\eta)} \underline{\underline{J}}^{(\nu,\eta)} \underline{\underline{f}}^{(*,\kappa,\zeta)}(t) = & - \frac{1}{2} (\underline{u}^{(h,\nu)}(t))^T (\underline{\underline{R}}^{(\eta)})^T \underline{\underline{B}}^{(\eta)} \underline{\underline{J}}^{(\nu,\eta)} \underline{\underline{A}}^{(\nu,\eta)} \underline{\underline{R}}^{(\eta)} \underline{u}^{(h,\nu)}(t) \\ & - \frac{1}{2} (\underline{u}^{(h,\nu)}(t))^T (\underline{\underline{R}}^{(\eta)})^T (\underline{\underline{T}}^{(\kappa,\zeta)})^T \underline{\underline{B}}^{(\zeta)} \underline{\underline{J}}^{(\kappa,\zeta)} \underline{\underline{A}}^{(\kappa,\zeta)} \underline{\underline{R}}^{(\zeta)} \underline{u}^{(h,\kappa)}(t) \\ & - \frac{\alpha}{2} (\underline{u}^{(h,\nu)}(t))^T (\underline{\underline{R}}^{(\eta)})^T (\underline{\underline{T}}^{(\kappa,\zeta)})^T \underline{\underline{B}}^{(\zeta)} \underline{\underline{J}}^{(\kappa,\zeta)} |\underline{\underline{A}}^{(\kappa,\zeta)}| \underline{\underline{R}}^{(\zeta)} \underline{u}^{(h,\kappa)}(t) \\ & + \frac{\alpha}{2} (\underline{u}^{(h,\nu)}(t))^T (\underline{\underline{R}}^{(\zeta)})^T (\underline{\underline{T}}^{(\kappa,\zeta)})^T \underline{\underline{B}}^{(\zeta)} \underline{\underline{J}}^{(\kappa,\zeta)} |\underline{\underline{A}}^{(\kappa,\zeta)}| \underline{\underline{T}}^{(\kappa,\zeta)} \underline{\underline{R}}^{(\eta)} \underline{u}^{(h,\nu)}(t), \end{aligned} \quad (4.25)$$

where  $|\underline{\underline{A}}^{(\kappa,\zeta)}| \in \mathbb{R}^{N_\zeta \times N_\zeta}$  contains the absolute values of the entries of  $\underline{\underline{A}}^{(\kappa,\zeta)}$ , and  $\underline{\underline{T}}^{(\kappa,\zeta)} \in \mathbb{R}^{N_\zeta \times N_\zeta}$  denotes the permutation matrix corresponding to the reordering of the facet nodes described in Assumption 3.2. We therefore obtain the final energy balance in (4.22) by summing (4.23) over all  $\kappa \in \{1 : N_e\}$  and noting that (4.24) and (4.25) result in a net contribution due to the interface  $\Omega^{(\kappa)} \cap \Omega^{(\nu)}$  given by the quadratic form

$$- \frac{\alpha}{2} (\underline{\underline{R}}^{(\zeta)} \underline{u}^{(h,\kappa)}(t) - \underline{\underline{T}}^{(\kappa,\zeta)} \underline{\underline{R}}^{(\eta)} \underline{u}^{(h,\nu)}(t))^T \underline{\underline{B}}^{(\zeta)} \underline{\underline{J}}^{(\kappa,\zeta)} |\underline{\underline{A}}^{(\kappa,\zeta)}| (\underline{\underline{R}}^{(\zeta)} \underline{u}^{(h,\kappa)}(t) - \underline{\underline{T}}^{(\kappa,\zeta)} \underline{\underline{R}}^{(\eta)} \underline{u}^{(h,\nu)}(t)), \quad (4.26)$$

which is non-positive for  $\alpha \geq 0$  due to the fact that the diagonal matrix  $\underline{\underline{B}}^{(\zeta)} \underline{\underline{J}}^{(\kappa,\zeta)} |\underline{\underline{A}}^{(\kappa,\zeta)}|$  is positive semidefinite, and is equal to zero for  $\alpha = 0$ , thus resulting in energy conservation in such a case.  $\square$

*Remark 4.3.* It is straightforward to show that the numerical flux in (4.21), which represents a standard choice for the linear advection equation, is consistent and conservative, as needed for Theorems 4.1 and 4.2.



## 5 Numerical Experiments

### 5.1 Problem Definition

We solve the two-dimensional linear advection equation on a square domain  $\Omega := (0, L)^2$  of side length  $L \in \mathbb{R}^+$  with periodic boundary conditions in both directions, where the velocity is given by  $\mathbf{a} := a[\sin(\theta), \cos(\theta)]^T$ , with a magnitude of  $a \in \mathbb{R}^+$  and a direction of  $\theta \in [0, 2\pi)$ . The initial condition is given by

$$U^0(\mathbf{x}) := \sin(2\pi x_1/L) \sin(2\pi x_2/L), \quad (5.1)$$

and the solution is evolved forward in time for one period, until  $T = L/(a \max(|\cos(\theta)|, |\sin(\theta)|))$ . For the numerical experiments in this section, the solver is run with values of  $a = \sqrt{2}$ ,  $\theta = \pi/4$ , and  $L = 1$ .

### 5.2 Numerical Methods and Implementation

We now describe the specifics of the methods described in §4 as implemented in `CLOUD.jl` (Conservation Laws on Unstructured Domains),<sup>3</sup> an unstructured high-order solver for PDEs developed by the first author.

#### 5.2.1 Spatial Discretization

For the numerical experiments in this paper, we implement the skew-symmetric nodal and modal formulations in (4.5) and (4.9), where, taking  $q_1 = q_2 = p$ , we use LG and LGR quadrature rules in the  $\eta_1$  and  $\eta_2$  directions, respectively. This ensures that the conditions of Lemmas 3.1 and 4.1 are met, resulting in valid SBP operators on the reference element as well as the physical element. We employ a direct solver for the mass matrix as opposed to an iterative solver, as our goal is to confirm the theoretical results, which assume that such a system is solved exactly. For comparison, we also implement standard quadrature-based modal DG methods, in which the numerical solution  $U^{(h,\kappa)}(\cdot, t) \in \mathbb{P}_p(\hat{\mathcal{T}}^2)$  satisfies a weak formulation on the reference element given for all  $t \in (0, T)$  and  $\kappa \in \{1 : N_e\}$  by

$$\begin{aligned} \int_{\hat{\mathcal{T}}^2} \left( V(\boldsymbol{\xi}) \frac{\partial J^{(\kappa)}(\boldsymbol{\xi}) U^{(h,\kappa)}(\boldsymbol{\xi}, t)}{\partial t} - \nabla_{\boldsymbol{\xi}} V(\boldsymbol{\xi}) \cdot (J^{(\kappa)}(\boldsymbol{\xi}) (\nabla_{\boldsymbol{\xi}} \mathbf{X}^{(\kappa)}(\boldsymbol{\xi}))^{-1} \mathbf{F}(U^{(h,\kappa)}(\boldsymbol{\xi}, t))) \right) d\boldsymbol{\xi} \\ + \sum_{\zeta=1}^3 \int_{\hat{\mathcal{E}}(\zeta)} V(\boldsymbol{\xi}) J^{(\kappa,\zeta)}(\boldsymbol{\xi}) F^{(*,\kappa,\zeta)}(\boldsymbol{\xi}, t) d\hat{s} = 0, \quad \forall V \in \mathbb{P}_p(\hat{\mathcal{T}}^2), \end{aligned} \quad (5.2)$$

where we approximate the volume integrals using symmetric quadrature rules of degree  $2p$  from Xiao and Gimbutas [34] and approximate the facet integrals using LG quadrature rules of degree  $2p+1$ .

*Remark 5.1.* As discussed by the authors in [9], the standard quadrature-based DG method satisfies the SBP property on the reference element for volume and facet quadrature rules of at least degree  $2p-1$  and  $2p$ , respectively, resulting in energy stability for an affine mapping. However, such a scheme does not satisfy the SBP property on curved physical elements unless the integrals in (5.2) are evaluated exactly. For DG methods employing inexact integration, we cannot therefore guarantee *a priori* that the discrete energy will decay monotonically when solving the linear advection equation on a curvilinear mesh. While this could be remedied by using a skew-symmetric formulation analogous to (4.9), we have chosen to use the standard weak-form DG method as a baseline scheme for comparison due to its widespread use among practitioners.

#### 5.2.2 Curvilinear Mesh Generation

To generate the meshes used for the computations in this section, we begin with a regular Cartesian grid consisting of  $M \in \mathbb{N}$  equal intervals in each direction and subdivide each quadrilateral into two triangles in order to obtain a total of  $N_e = 2M^2$  elements, each of which is the image of  $\hat{\mathcal{T}}^2$  under an affine mapping  $\mathbf{X}_{\text{affine}}^{(\kappa)} \in [\mathbb{P}_1(\hat{\mathcal{T}}^2)]^2$ . As the nodal sets considered in this work are not symmetric, special care must be taken regarding the element orientation in order to ensure that the facet quadrature nodes in (3.10) satisfy

<sup>3</sup>`CLOUD.jl` is available under the GNU General Public License at <https://github.com/tristanmontoya/CLOUD.jl>.

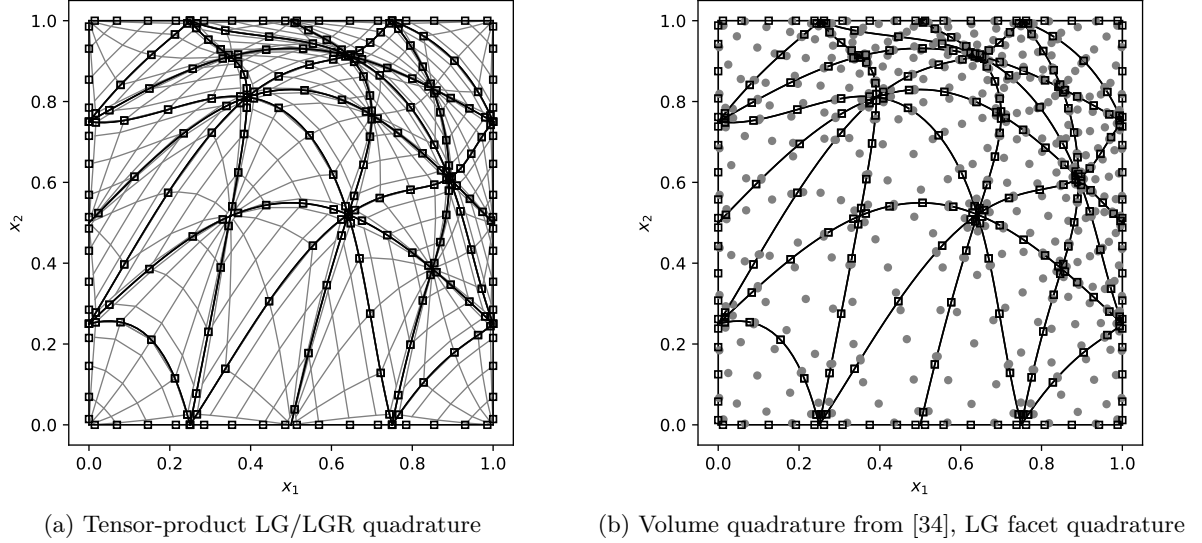


Figure 2: Isoparametric warped meshes for  $p = 4$  and  $M = 4$

Assumption 3.2. Next, in order to mimic the highly stretched curvilinear elements which could result from the meshing of complex geometries, we consider the warping function  $\mathbf{W} : [0, L]^2 \rightarrow [0, L]^2$  given by

$$\mathbf{W}(\mathbf{x}) := \begin{bmatrix} x_1 + \frac{1}{5}L \sin(\pi x_1/L) \sin(\pi x_2/L) \\ x_2 + \frac{1}{5}L \exp(1 - x_2/L) \sin(\pi x_1/L) \sin(\pi x_2/L) \end{bmatrix}, \quad (5.3)$$

which was employed by Del Rey Fernández *et al.* in [35], and construct a polynomial mapping  $\mathbf{X}^{(\kappa)} \in [\mathbb{P}_l(\hat{\mathcal{T}}^2)]^2$  by interpolating such a function on a set of nodes  $\{\boldsymbol{\xi}_{\text{map}}^{(i)}\}_{i=1}^{N_l^*}$  supporting a Lagrange basis  $\{\ell_{\text{map}}^{(i)}\}_{i=1}^{N_l^*}$  of degree  $l \in \mathbb{N}$ , which is additionally required to include  $l + 1$  LGL nodes on each edge of the triangle.<sup>4</sup> The resulting coordinate transformation is therefore given for  $\kappa \in \{1 : N_e\}$  by

$$\mathbf{X}^{(\kappa)}(\boldsymbol{\xi}) := \sum_{i=1}^{N_l^*} \mathbf{W}(\mathbf{X}_{\text{affine}}^{(\kappa)}(\boldsymbol{\xi}_{\text{map}}^{(i)})) \ell_{\text{map}}^{(i)}(\boldsymbol{\xi}). \quad (5.4)$$

For the computations in this section, we consider isoparametric mappings (i.e.  $l = p$ ), which, as discussed in Remark 4.2, result in physical operators which satisfy Assumption 4.1. Examples of meshes constructed based on the above procedure are shown in Figure 2, where the volume quadrature nodes and facet quadrature nodes for the proposed tensor-product approach as well as the standard DG method are pictured.

### 5.2.3 Temporal Discretization

The five-stage, fourth-order explicit Runge-Kutta method proposed by Carpenter and Kennedy [37] is used to advance the solution in time, with a time step given by  $\Delta t = C_t h/a$ , where the characteristic element size is taken to be  $h = L/M$ . Motivated by the Courant-Friedrichs-Lewy (CFL) condition described for the standard DG method by Cockburn and Shu [38, §2.2], we choose  $C_t = \beta/(2p+1)$  with  $\beta = 2.5 \times 10^{-3}$  in order to ensure that the error due to the temporal discretization is dominated by that of the spatial discretization.

## 5.3 Results

### 5.3.1 Refinement Studies

We solve the linear advection problem described in §5.1 using the skew-symmetric nodal and modal tensor-product formulations as well as a standard DG method, where we present results for polynomial degrees

<sup>4</sup>We use the nodes provided in `NodesAndModes.jl`, which are based on the “interpolatory warp and blend” procedure in [36].

$p = 4$  and  $p = 9$  using the upwind and central numerical fluxes obtained from (4.21) with  $\alpha = 1$  and  $\alpha = 0$ , respectively. The meshes are constructed following the procedure described in §5.2, beginning with  $M = 2$  and doubling the number of edges in each direction at each refinement. Suitable measures of conservation and energy stability for periodic problems may be defined as the net change in the quantities within the time derivatives on the left-hand sides of (4.16) and (4.22), respectively, as given by

$$\text{Conservation Metric} := \sum_{\kappa=1}^{N_e} (\mathbf{1}^{(N_q)})^T \underline{\underline{M}} \underline{J}^{(\kappa)} \underline{u}^{(h,\kappa)}(T) - \sum_{\kappa=1}^{N_e} (\mathbf{1}^{(N_q)})^T \underline{\underline{M}} \underline{J}^{(\kappa)} \underline{u}^{(h,\kappa)}(0), \quad (5.5)$$

and

$$\text{Energy Metric} := \frac{1}{2} \sum_{\kappa=1}^{N_e} (\underline{u}^{(h,\kappa)}(T))^T \underline{\underline{M}} \underline{J}^{(\kappa)} \underline{u}^{(h,\kappa)}(T) - \frac{1}{2} \sum_{\kappa=1}^{N_e} (\underline{u}^{(h,\kappa)}(0))^T \underline{\underline{M}} \underline{J}^{(\kappa)} \underline{u}^{(h,\kappa)}(0), \quad (5.6)$$

where the former is expected to be zero for the proposed nodal and modal tensor-product schemes as a consequence of Theorem 4.1, whereas the latter is expected to be non-positive for  $\alpha = 1$  and zero for  $\alpha = 0$  due to Theorem 4.3. The corresponding conservation and energy metrics for the standard DG scheme are defined as in [9, §5.1.9]. We also evaluate the accuracy of the proposed schemes using the error metric

$$\text{Error Metric} := \sqrt{\sum_{\kappa=1}^{N_e} (\underline{e}^{(h,\kappa)}(T))^T \underline{\underline{M}} \underline{J}^{(\kappa)} \underline{e}^{(h,\kappa)}(T)}, \quad (5.7)$$

defining an approximation to the  $L^2$  norm of the solution error, which is evaluated at the volume quadrature nodes for each element as  $\underline{e}^{(h,\kappa)}(T) := [U(\mathbf{X}^{(\kappa)}(\boldsymbol{\xi}^{(1)}), T) - u_1^{(h,\kappa)}(T), \dots, U(\mathbf{X}^{(\kappa)}(\boldsymbol{\xi}^{(N_q)}), T) - u_{N_q}^{(h,\kappa)}(T)]$ , and we present estimates of the order of convergence between successive grids with respect to  $h = L/M$ , or equivalently (for fixed  $p$ ), the square root of the total number of degrees of freedom.

The results for the skew-symmetric nodal and modal tensor-product formulations are shown in Tables 1 and 2, respectively, with those for the standard DG formulation provided for comparison in Table 3. As expected from the analysis in §4.3, both the nodal and modal discretizations are discretely conservative, while the upwind and central schemes conserve and dissipate energy, respectively, up to small discrepancies due to the influence of roundoff error as well as the temporal discretization, neither of which were accounted for in the analysis. Moreover, as discussed in Remark 5.1, the standard DG method does not satisfy the SBP property on the physical element, and finite energy growth is accordingly observed in some cases with  $\alpha = 0$ . As such, we cannot guarantee that the energy will remain bounded for all time, relying instead on the dissipation of the upwind numerical flux to stabilize the scheme in practice. The three discretization approaches are seen to be fairly similar in accuracy, both in terms of the level of error and convergence rate, where the upwind schemes achieve a nearly optimal (i.e.  $p+1$ ) rate of convergence, indicating that exploiting the potential benefits of the proposed tensor-product SBP formulations does not incur any loss of accuracy relative to the baseline standard DG method.

### 5.3.2 Semi-Discrete Operator Spectra

Using the implicitly restarted Arnoldi method provided in `Arpack.jl`, a Julia wrapper for the ARPACK library, to numerically compute the eigenvalues of the semi-discrete formulations resulting from the proposed spatial discretizations in a non-intrusive manner, we plot the spectra of the proposed nodal and modal skew-symmetric tensor-product discretizations on triangles in Figure 3, where schemes of degree  $p = 4$  are constructed on the isoparametric curvilinear meshes with  $M = 4$  pictured in Figure 2. In accordance with the theoretical analysis, the eigenvalues for discretizations employing an upwind numerical flux all have a non-positive real part, whereas those for the discretizations employing a central numerical flux are purely imaginary, with discrepancies on the order of at most  $10^{-12}$  likely arising as a result of the inexact solution of the eigenvalue problem.

Additionally, the clustering of nodes near the singularity of the mapping in (3.1) results in some eigenvalues for the nodal tensor-product formulation in (4.5) being very large in magnitude, which has the effect

Table 1: Refinement studies for the skew-symmetric tensor-product nodal formulation

$p$	$N_e$	Conservation Metric		Energy Metric		Error Metric		Order	
		Upwind	Central	Upwind	Central	Upwind	Central	Upwind	Central
4	8	9.454e-16	4.528e-16	-2.073e-02	-3.410e-13	9.955e-02	2.531e-01		
	32	-1.619e-15	-8.964e-16	-1.296e-03	-3.886e-16	1.536e-02	4.280e-02	2.70	2.56
	128	-3.486e-16	2.030e-16	-7.058e-06	1.166e-15	5.764e-04	4.651e-03	4.74	3.20
	512	-2.929e-15	-3.324e-15	-1.886e-08	2.082e-15	1.750e-05	1.751e-04	5.04	4.73
	2048	-1.746e-14	-1.759e-14	-3.987e-11	1.846e-14	5.660e-07	7.911e-06	4.95	4.47
9	8	-4.318e-15	-4.800e-15	-1.254e-05	-5.412e-15	1.308e-03	5.353e-03		
	32	-9.622e-15	-9.331e-15	-2.111e-08	-5.496e-15	5.013e-05	1.962e-04	4.71	4.77
	128	3.035e-14	2.972e-14	-1.404e-13	-1.957e-14	5.795e-08	5.013e-07	9.76	8.61
	512	9.695e-14	9.687e-14	-8.613e-14	-8.599e-14	6.604e-11	8.130e-10	9.78	9.27

Table 2: Refinement studies for the skew-symmetric tensor-product modal formulation

$p$	$N_e$	Conservation Metric		Energy Metric		Error Metric		Order	
		Upwind	Central	Upwind	Central	Upwind	Central	Upwind	Central
4	8	-2.984e-16	9.957e-16	-2.283e-02	-7.327e-15	1.098e-01	1.815e-01		
	32	-2.334e-15	-6.301e-16	-1.668e-03	1.180e-15	1.415e-02	4.044e-02	2.96	2.17
	128	5.539e-16	-2.626e-16	-8.943e-06	1.818e-15	5.349e-04	4.170e-03	4.73	3.28
	512	-3.107e-15	-3.065e-15	-2.362e-08	2.692e-15	1.581e-05	2.062e-04	5.08	4.34
	2048	-1.703e-14	-1.851e-14	-4.991e-11	1.776e-14	4.973e-07	3.566e-06	4.99	5.85
9	8	-5.492e-15	-6.687e-15	-1.761e-05	-6.939e-16	1.138e-03	6.186e-03		
	32	-1.194e-14	-8.830e-15	-2.742e-08	-3.830e-15	4.104e-05	1.780e-04	4.79	5.12
	128	2.936e-14	2.964e-14	-1.670e-13	-1.700e-14	4.118e-08	5.288e-07	9.96	8.40
	512	9.791e-14	9.849e-14	-8.665e-14	-8.667e-14	4.831e-11	7.207e-10	9.74	9.52

Table 3: Refinement studies for the standard DG formulation

$p$	$N_e$	Conservation Metric		Energy Metric		Error Metric		Order	
		Upwind	Central	Upwind	Central	Upwind	Central	Upwind	Central
4	8	-3.652e-15	-7.060e-16	-1.838e-02	3.665e-03	1.150e-01	2.110e-01		
	32	-6.618e-16	-9.793e-16	-1.539e-03	-3.245e-05	1.489e-02	4.131e-02	2.95	2.35
	128	-1.081e-15	-4.471e-16	-8.251e-06	-1.319e-07	5.461e-04	3.587e-03	4.77	3.53
	512	-5.347e-15	-6.036e-15	-2.192e-08	8.231e-10	1.541e-05	2.128e-04	5.15	4.08
	2048	-3.426e-14	-3.431e-14	-4.633e-11	2.168e-13	4.864e-07	3.370e-06	4.99	5.98
9	8	-8.859e-15	-1.245e-14	-1.498e-05	5.288e-06	1.441e-03	6.297e-03		
	32	-1.615e-14	-1.894e-14	-2.497e-08	9.624e-09	4.746e-05	2.032e-04	4.92	4.95
	128	5.784e-14	5.878e-14	-1.835e-13	-2.266e-14	4.915e-08	5.612e-07	9.92	8.50
	512	1.929e-13	1.940e-13	-1.723e-13	-1.704e-13	5.912e-11	7.283e-10	9.70	9.59

of severely limiting the maximum allowable time step when used with explicit time-marching methods in a stability-limited context. This limitation was recognized as early as [30], wherein the modal basis in (3.14) was proposed as a remedy for such an issue, effectively redistributing the local degrees of freedom to avoid such clustering. We find that the proposed modal formulation in (4.9) is indeed superior in this regard, resulting in operators with much smaller spectral radii than those of the nodal schemes. Whether or not this justifies the additional cost for evaluating the semi-discrete residual (i.e. resulting from the mass matrix and generalized Vandermonde matrix) is dependent on the particular time-marching scheme being used as well as the problem being solved and is the subject of further investigation. However, we emphasize that both the nodal and modal approaches have been shown theoretically and numerically to result in energy stable

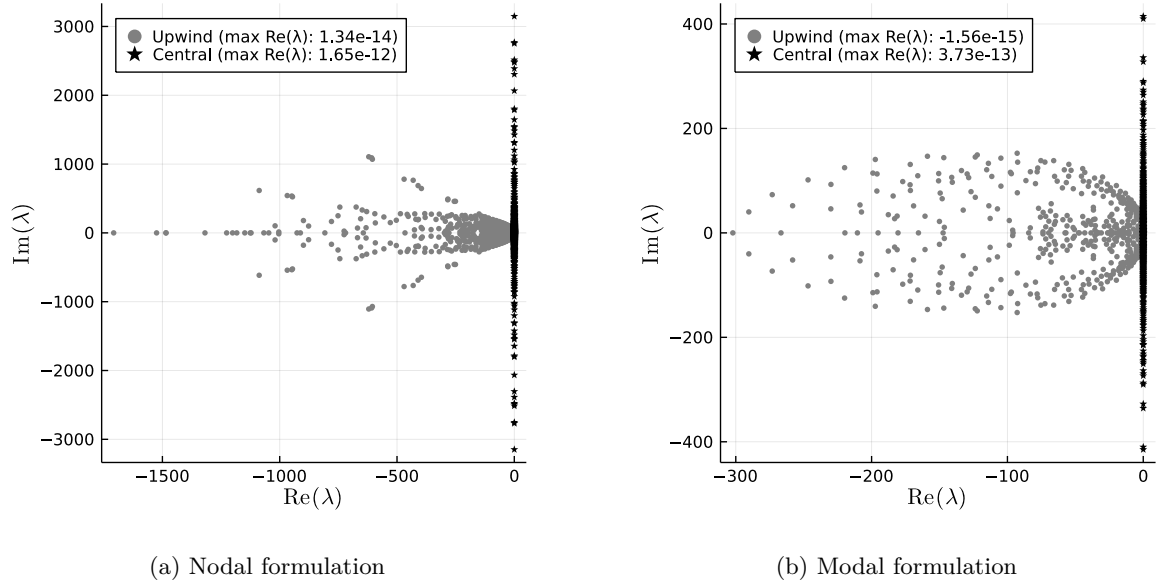


Figure 3: Semi-discrete operator spectra for proposed tensor-product discretizations with  $p = 4$  and  $M = 4$

and conservative schemes which are amenable to efficient matrix-free algorithms due to the tensor-product structure of their constituent operators.

## 6 Conclusions

By extending the SBP approach to tensor-product discretizations in collapsed coordinates, we have developed a methodology for constructing provably stable and conservative discretizations of any order on triangular elements which enable efficient matrix-free algorithms for evaluating local operators to be exploited. Although the focus of this paper is on skew-symmetric formulations for curvilinear coordinates, the approach is applicable to nonlinearly stable discretizations through the use of entropy-conservative two-point flux functions, and is expected to extend in a straightforward manner to three-dimensional discretizations on tetrahedral, prismatic, and pyramidal elements. While both the nodal and modal approaches proposed in this paper result in schemes which are similar in accuracy to a standard DG method of the same polynomial degree, the cost of residual evaluation for the nodal approach is lower, as the mass matrix for the curved physical element is diagonal, and the solution is available directly at the volume quadrature nodes without the need for interpolation. However, unlike the modal schemes, the nodal methods result in semi-discrete operator spectra containing eigenvalues which are very large in magnitude, thereby requiring smaller time steps to maintain stability when used with explicit time-marching methods, and potentially resulting in linear systems with poorer conditioning in the context of implicit methods. The study of such competing objectives along with a thorough analysis of the efficiency of the proposed schemes relative to existing methods is the subject of future work, as is the development of optimized nodal schemes with reduced spectral radii.

## Acknowledgements

The authors are grateful to Gianmarco Mengaldo for the discussions which initiated this project. A full list of dependencies for `CLOUD.jl` can be found on the solver's GitHub page; we are especially grateful for the core functionality made possible through the use of the `LinearMaps.jl`, `StartUpDG.jl`, and `DifferentialEquations.jl` packages. Computations were performed on the Niagara supercomputer at the SciNet HPC Consortium [39], which is funded by the Canada Foundation for Innovation, the Government of Ontario, the Ontario Research Fund – Research Excellence, and the University of Toronto.

## References

- [1] S. A. Orszag, “Spectral methods for problems in complex geometries,” *Journal of Computational Physics*, vol. 37, no. 1, pp. 70–92, Aug. 1980.
- [2] P. E. J. Vos, S. J. Sherwin, and R. M. Kirby, “From  $h$  to  $p$  efficiently: Implementing finite and spectral/ $hp$  element methods to achieve optimal performance for low- and high-order discretisations,” *Journal of Computational Physics*, vol. 229, no. 13, pp. 5161–5181, Jul. 2010.
- [3] C. D. Cantwell, S. J. Sherwin, R. M. Kirby, and P. H. J. Kelly, “From  $h$  to  $p$  efficiently: Strategy selection for operator evaluation on hexahedral and tetrahedral elements,” *Computers & Fluids*, vol. 43, no. 1, pp. 23–28, Apr. 2011.
- [4] S. J. Sherwin and G. E. Karniadakis, “A triangular spectral element method: Applications to the incompressible Navier-Stokes equations,” *Computer Methods in Applied Mechanics and Engineering*, vol. 123, no. 1-4, pp. 189–229, Jun. 1995.
- [5] R. M. Kirby, T. C. Warburton, I. Lomtev, and G. E. Karniadakis, “A discontinuous Galerkin spectral/ $hp$  method on hybrid grids,” *Applied Numerical Mathematics*, vol. 33, no. 1-4, pp. 393–405, May 2000.
- [6] D. Moxey, R. Amici, and R. M. Kirby, “Efficient matrix-free high-order finite element evaluation for simplicial elements,” *SIAM Journal on Scientific Computing*, vol. 42, no. 3, pp. C97–C123, May 2020.
- [7] D. C. Del Rey Fernández, J. E. Hicken, and D. W. Zingg, “Review of summation-by-parts operators with simultaneous approximation terms for the numerical solution of partial differential equations,” *Computers & Fluids*, vol. 95, pp. 171–196, May 2014.
- [8] M. Svärd and J. Nordström, “Review of summation-by-parts schemes for initial-boundary-value problems,” *Journal of Computational Physics*, vol. 268, pp. 17–38, Jul. 2014.
- [9] T. Montoya and D. W. Zingg, “A unifying algebraic framework for discontinuous Galerkin and flux reconstruction methods based on the summation-by-parts property,” *Journal of Scientific Computing*, vol. 92, no. 3, Sep. 2022.
- [10] J. E. Hicken, D. C. Del Rey Fernández, and D. W. Zingg, “Multidimensional summation-by-parts operators: General theory and application to simplex elements,” *SIAM Journal on Scientific Computing*, vol. 38, no. 4, A1935–A1958, Jul. 2016.
- [11] J. Crean, J. E. Hicken, D. C. Del Rey Fernández, D. W. Zingg, and M. H. Carpenter, “Entropy-stable summation-by-parts discretization of the Euler equations on general curved elements,” *Journal of Computational Physics*, vol. 356, pp. 410–438, Mar. 2018.
- [12] D. C. Del Rey Fernández, J. E. Hicken, and D. W. Zingg, “Simultaneous approximation terms for multi-dimensional summation-by-parts operators,” *Journal of Scientific Computing*, vol. 75, no. 1, pp. 83–110, Apr. 2018.
- [13] D. C. Del Rey Fernández, J. Crean, M. H. Carpenter, and J. E. Hicken, “Staggered-grid entropy-stable multidimensional summation-by-parts discretizations on curvilinear coordinates,” *Journal of Computational Physics*, vol. 392, pp. 161–186, Sep. 2019.
- [14] A. L. Marchildon and D. W. Zingg, “Optimization of multidimensional diagonal-norm summation-by-parts operators on simplices,” *Journal of Computational Physics*, vol. 411, Jun. 2020.
- [15] S. Shadpey and D. W. Zingg, “Entropy-stable multidimensional summation-by-parts discretizations on  $hp$ -adaptive curvilinear grids for hyperbolic conservation laws,” *Journal of Scientific Computing*, vol. 82, no. 3, Mar. 2020.
- [16] J. Chan, “On discretely entropy conservative and entropy stable discontinuous Galerkin methods,” *Journal of Computational Physics*, vol. 362, pp. 346–374, Jun. 2018.
- [17] J. Chan and L. C. Wilcox, “On discretely entropy stable weight-adjusted discontinuous Galerkin methods: Curvilinear meshes,” *Journal of Computational Physics*, vol. 378, pp. 366–393, Feb. 2019.
- [18] T. Chen and C.-W. Shu, “Entropy stable high order discontinuous Galerkin methods with suitable quadrature rules for hyperbolic conservation laws,” *Journal of Computational Physics*, vol. 345, pp. 427–461, Sep. 2017.



- [19] T. Chen and C.-W. Shu, “Review of entropy stable discontinuous Galerkin methods for systems of conservation laws on unstructured simplex meshes,” *CSIAM Transactions on Applied Mathematics*, vol. 1, no. 1, pp. 1–52, Jun. 2020.
- [20] R. M. Kirby and G. E. Karniadakis, “De-aliasing on non-uniform grids: Algorithms and applications,” *Journal of Computational Physics*, vol. 191, no. 1, pp. 249–264, Oct. 2003.
- [21] R. M. Kirby and S. J. Sherwin, “Aliasing errors due to quadratic nonlinearities on triangular spectral/*hp* element discretisations,” *Journal of Engineering Mathematics*, vol. 56, no. 3, pp. 273–288, Oct. 2006.
- [22] G. Mengaldo, D. De Grazia, D. Moxey, P. E. Vincent, and S. J. Sherwin, “Dealiasing techniques for high-order spectral element methods on regular and irregular grids,” *Journal of Computational Physics*, vol. 299, pp. 56–81, Oct. 2015.
- [23] P. D. Thomas and C. K. Lombard, “Geometric conservation law and its application to flow computations on moving grids,” *AIAA Journal*, vol. 17, no. 10, pp. 1030–1037, Oct. 1979.
- [24] D. A. Kopriva, “Metric identities and the discontinuous spectral element method on curvilinear meshes,” *Journal of Scientific Computing*, vol. 26, no. 3, pp. 301–327, Mar. 2006.
- [25] M. E. Gurtin, E. Fried, and L. Anand, *The Mechanics and Thermodynamics of Continua*. Cambridge University Press, 2010.
- [26] M. G. Duffy, “Quadrature over a pyramid or cube of integrands with a singularity at a vertex,” *SIAM Journal on Numerical Analysis*, vol. 19, no. 6, pp. 1260–1262, Dec. 1982.
- [27] M. Abramowitz and I. A. Stegun, Eds., *Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables*. New York: Dover Publications, 1965.
- [28] J. Proriol, “Sur une famille de polynomes à deux variables orthogonaux dans un triangle,” *Comptes Rendus Hebdomadaires des Séances de l’Académie des Sciences*, vol. 245, pp. 2459–2461, Dec. 1957.
- [29] T. Koornwinder, “Two-variable analogues of the classical orthogonal polynomials,” in *Theory and Application of Special Functions*, R. Askey, Ed., Academic Press, 1975, pp. 435–495.
- [30] M. Dubiner, “Spectral methods on triangles and other domains,” *Journal of Scientific Computing*, vol. 6, no. 4, pp. 345–390, Dec. 1991.
- [31] D. A. Kopriva and G. J. Gassner, “On the quadrature and weak form choices in collocation type discontinuous Galerkin spectral element methods,” *Journal of Scientific Computing*, vol. 44, pp. 136–155, Aug. 2010.
- [32] W. Pazner and P.-O. Persson, “Approximate tensor-product preconditioners for very high order discontinuous Galerkin methods,” *Journal of Computational Physics*, vol. 354, pp. 344–369, Feb. 2018.
- [33] M. H. Carpenter, D. Gottlieb, and S. Abarbanel, “Time-stable boundary conditions for finite-difference schemes solving hyperbolic systems: Methodology and application to high-order compact schemes,” *Journal of Computational Physics*, vol. 111, no. 2, pp. 220–236, Apr. 1994.
- [34] H. Xiao and Z. Gimbutas, “A numerical algorithm for the construction of efficient quadrature rules in two and higher dimensions,” *Computers & Mathematics with Applications*, vol. 59, no. 2, pp. 663–676, Jan. 2010.
- [35] D. C. Del Rey Fernández, P. D. Boom, M. Shademan, and D. W. Zingg, “Numerical investigation of tensor-product summation-by-parts discretization strategies and operators,” in *55th AIAA Aerospace Sciences Meeting*, American Institute of Aeronautics and Astronautics, Jan. 2017.
- [36] J. Chan and T. Warburton, “A comparison of high order interpolation nodes for the pyramid,” *SIAM Journal on Scientific Computing*, vol. 37, no. 5, A2151–A2170, Sep. 2015.
- [37] M. H. Carpenter and C. A. Kennedy, “Fourth-order  $2N$ -storage Runge-Kutta schemes,” NASA Technical Memorandum 109112, Tech. Rep., Jun. 1994.
- [38] B. Cockburn and C.-W. Shu, “Runge-Kutta discontinuous Galerkin methods for convection-dominated problems,” *Journal of Scientific Computing*, vol. 16, no. 3, Sep. 2001.
- [39] M. Ponce, R. van Zon, S. Northrup, D. Gruner, J. Chen, F. Ertinaz, A. Fedoseev, L. Groer, F. Mao, B. C. Mundim, M. Nolta, J. Pinto, M. Saldarriaga, V. Slavnic, E. Spence, C.-H. Yu, and W. R. Peltier, “Deploying a top-100 supercomputer for large parallel workloads,” in *Proceedings of the Practice and Experience in Advanced Research Computing on Rise of the Machines (Learning)*, Association for Computing Machinery, Jul. 2019.