



# Generalized Summation-By-Parts Methods: Coordinate Transformations, Quadrature Accuracy, and Functional Superconvergence

David A. Craig Penner\* and David W. Zingg†

University of Toronto Institute for Aerospace Studies, Toronto, Ontario, M3H 5T6

**We investigate coordinate transformations, quadrature accuracy, and functional superconvergence for diagonal-norm tensor-product generalized summation-by-parts operators. We show that projection operators of degree  $r \geq 2p$  are required to preserve quadrature accuracy, and therefore functional superconvergence, in curvilinear coordinates when: (1) the Jacobian of the coordinate transformation is approximated by the same generalized summation-by-parts operator that is used to approximate the flux terms and (2) the degree of the generalized summation-by-parts operator is lower than the degree of the polynomial used to represent the geometry of interest. Legendre-Gauss-Lobatto and Legendre-Gauss element-type operators are considered. When the aforementioned condition (2) is violated for the Legendre-Gauss operators, there is an even-odd quadrature convergence pattern that is explained by the cancellation of the leading truncation error terms for the projection operators that correspond to the odd-degree Legendre-Gauss operators.**

## I. Introduction

THIS paper examines some of the practical issues associated with using generalized summation-by-parts (SBP) schemes to solve aerodynamic flows around complex geometries on curvilinear domains. In combination with simultaneous approximation terms (SATs) [1–3], which allow boundary conditions to be enforced in a stable manner, SBP methods [4, 5] provide a provably stable, conservative, and consistent way to numerically solve a wide class of linear and nonlinear partial differential equations (PDEs) [6–8]. For classical diagonal-norm SBP operators, if the discretization is dual consistent, then the underlying quadrature and solution functionals converge at the same rate as the order of the interior operator [9, 10]. The objective of the present work is to develop the conditions necessary to preserve quadrature accuracy and superconvergence on curvilinear domains discretized using generalized SBP operators that do not include one or both boundary nodes.

## II. Notation and Definitions

The notation is adapted from [6, 9, 11]. Upper case script letters, e.g.,  $\mathcal{U}$ , denote continuous functions, while lower case bold letters, e.g.,  $\mathbf{u}$ , indicate the restriction of these continuous functions onto a set of nodes. A sans-serif capital letter, e.g.,  $H$ , represents a matrix. Let  $\Omega \subset \mathbb{R}^d$  be a  $d$ -dimensional Lipschitz domain. The inner product and norm are defined for two square-integrable real-valued functions,  $\mathcal{U}, \mathcal{V} \in L^2(\Omega)$ , as

$$(\mathcal{U}, \mathcal{V}) \equiv \int_{\Omega} \mathcal{U}\mathcal{V} \, d\Omega, \quad \text{and} \quad \|\mathcal{U}\|_{L^2}^2 \equiv \int_{\Omega} \mathcal{U}^2 \, d\Omega, \quad (1)$$

which are approximated by the discrete inner product,  $(\mathbf{u}, \mathbf{v})_H \equiv \mathbf{u}^T H \mathbf{v}$ , and norm,  $\|\mathbf{u}\|_H^2 \equiv \mathbf{u}^T H \mathbf{u}$ , respectively. We use the following definition of a generalized SBP operator [11].

**Definition 1.** *Generalized summation-by-parts operator for the first derivative: A matrix operator,  $D \in \mathbb{R}^{(n+1) \times (n+1)}$ , is an SBP operator that approximates the derivative  $\frac{\partial}{\partial x}$ , on the nodal distribution  $\Omega_x \in [a, b]$  having  $n + 1$  nodes, of degree  $p$  if*

\*Ph.D. Candidate, Institute for Aerospace Studies, 4925 Dufferin Street, and AIAA Student Member (david.craigpenner@mail.utoronto.ca).

†University of Toronto Distinguished Professor of Computational Aerodynamics and Sustainable Aviation, Director, Centre for Research in Sustainable Aviation, Director, Centre for Computational Science and Engineering, Institute for Aerospace Studies, 4925 Dufferin Street, and AIAA Associate Fellow (dwz@oddjob.utias.utoronto.ca).

- 1)  $D\mathbf{x}^k = \mathbf{H}^{-1}\mathbf{Q}\mathbf{x}^{k-1} = k\mathbf{x}^{k-1}$ ,  $k = 0, 1, \dots, p$ ;
- 2)  $\mathbf{H}$ , the norm matrix, is symmetric and positive definite; and
- 3)  $\mathbf{Q} + \mathbf{Q}^T = \mathbf{E}$ , where  $(\mathbf{x}^i)^T \mathbf{E} \mathbf{x}^j = b^{i+j} - a^{i+j}$ ,  $i, j = 0, 1, \dots, r, r \geq p$ .

From Definition 1, we see that the accuracy of an SBP operator is expressed in terms of the maximum degree of monomial for which it is exact. For operators constructed according to Definition 1 on tensor-product domains, it is common to decompose  $\mathbf{E}$  as

$$\mathbf{E} = \mathbf{t}_R \mathbf{t}_R^T - \mathbf{t}_L \mathbf{t}_L^T, \quad \text{where} \quad \mathbf{t}_L^T \mathbf{x}^k = a^k, \quad \mathbf{t}_R^T \mathbf{x}^k = b^k, \quad k = 0, 1, \dots, r. \quad (2)$$

Throughout this work we refer to  $\mathbf{t}_R$  and  $\mathbf{t}_L$  as projection operators.

We define two classes of SBP operators: classical SBP operators and generalized SBP operators [11]. Classical SBP operators are constructed with a repeated interior stencil on a uniform nodal distribution that includes both boundary nodes. Generalized SBP operators can be constructed on nonuniform nodal distributions that do not include one or both boundary nodes. A classical diagonal-norm SBP operator of degree  $p$  is associated with a degree  $\tau = 2p - 1$  quadrature rule, while a generalized diagonal-norm SBP operator of degree  $p$  is associated with a degree  $\tau \geq 2p - 1$  quadrature rule [11]. The order of a quadrature is equal to  $\tau + 1$ . Note that classical SBP operators are a subset of generalized SBP operators.

### III. Analysis

Solving PDEs on complex geometries normally requires the use of a curvilinear coordinate transformation which relates points in the physical domain to points in a reference space. For classical SBP operators, the impact that the geometric terms introduced by the coordinate transformation have on the accuracy of diagonal-norm SBP quadrature was studied in [12]. It was found that classical diagonal-norm SBP quadrature retains its order  $2p$  theoretical accuracy in curvilinear coordinates when the Jacobian of the transformation is constructed using the SBP derivative operator associated with the quadrature (i.e., the norm,  $\mathbf{H}$ ). For tensor-product domains, besides retaining quadrature accuracy for classical SBP operators, the motivation for constructing the Jacobian using the SBP derivative operator associated with the norm arises from studies by, for example, [13], that imply that the derivative operators used to approximate the fluxes should also be used to compute the metrics to satisfy the metric invariants.

To see the effect of curvilinear transformations on the quadrature accuracy of generalized SBP schemes that do not include one or both boundary nodes, we first formalize the coordinate transformation in a one-dimensional setting, similar to [12]. We decompose the physical domain,  $\Omega_x \in [a, b]$ , into  $n_e$  non-overlapping elements  $K_i$ , such that  $\Omega_x = \cup_{i=1}^{n_e} K_i$ ,  $K_i \cap K_j = \emptyset$ ,  $i \neq j$ . Suppose there exists an invertible transformation that maps the projection of the physical domain onto the element  $K_i$ ,  $(\Omega_x)_{K_i}$ , to the reference domain,  $\Omega_\xi \in [0, 1]$ . The change of variable theorem gives

$$\int_a^b \mathbf{U} dx = \sum_{i=1}^{n_e} \int_0^1 \mathbf{U}_{K_i} \mathcal{J} d\xi, \quad (3)$$

where  $\mathcal{J} = \frac{dx_{K_i}}{d\xi}$  is the Jacobian of the inverse transformation, and  $\mathbf{U}_{K_i}$  is  $\mathbf{U}$  on the element  $K_i$ . The right-hand side of Eq. (3) can be approximated as

$$\sum_{i=1}^{n_e} \mathbf{u}_{K_i}^T \mathbf{H} D_\xi \mathbf{x}_{K_i} = \sum_{i=1}^{n_e} \mathbf{u}_{K_i}^T \mathbf{Q}_\xi \mathbf{x}_{K_i}. \quad (4)$$

Here,  $D_\xi$  and  $\mathbf{Q}_\xi$  are defined on the reference domain. For classical SBP operators, to show that Eq. (4) is a  $2p$ -order accurate approximation to the right-hand side of Eq. (3), [12] introduces and proves the following theorem.

**Theorem 1.** *Let  $\mathbf{D} = \mathbf{H}^{-1}\mathbf{Q}$  be an SBP operator of degree  $p$  approximating the first derivative. Then*

$$(\mathbf{z}, \mathbf{D}\mathbf{u})_H = \mathbf{z}^T \mathbf{Q}\mathbf{u}$$

*is a  $2p$ -order accurate approximation to the integral*

$$\int_0^1 \mathcal{Z} \frac{d\mathbf{U}}{dx} dx,$$

*where  $\mathcal{Z} \frac{d\mathbf{U}}{dx} \in C^{2p-1}[0, 1]$ .*

**Remark 1.** The first derivative SBP operator in Theorem 1 is a classical SBP operator constructed on a uniform domain that includes both boundary nodes.

**Remark 2.** The accuracy of the discrete quadrature,  $(z, \mathbf{D}\mathbf{u})_H$ , in Theorem 1 is defined in terms of order rather than degree. In this case, the quadrature rule associated with the norm matrix,  $H$ , is a degree  $\tau = 2p - 1$  quadrature rule. For classical diagonal-norm SBP operators that include both boundary nodes, if  $(z, \mathbf{D}\mathbf{u})_H$  is a  $2p$ -order approximation to the integral

$$\int_0^1 \mathcal{Z} \frac{d\mathcal{U}}{dx} dx,$$

this means that

$$\int_0^1 \mathcal{Z} \frac{d\mathcal{U}}{dx} dx = (z, \mathbf{D}\mathbf{u})_H + O(h^{2p}),$$

where  $h$  is the uniform mesh spacing. Thus, the order of the quadrature is equal to the degree of the quadrature plus one.

For generalized SBP operators, we can prove an analogous theorem, which is essentially the same as Theorem 3.4 in [14], where the proof is similar to that given in [12] for Theorem 1 above. The theorem and proof follow.

**Theorem 2.** Let  $D = H^{-1}Q$  be a generalized SBP operator of degree  $p$  approximating the first derivative operator, as in Definition 1. Then

$$(z, \mathbf{D}\mathbf{u})_H = z^T Q \mathbf{u} = \int_0^1 \mathcal{Z} \frac{d\mathcal{U}}{dx} dx, \quad i, j \leq r, i + j \leq 2p,$$

where  $\mathcal{Z} \frac{d\mathcal{U}}{dx} \in C^{2p-1}[0, 1]$ ,  $z = \mathbf{x}^i$ , and  $\mathbf{u} = \mathbf{x}^j$ .

*Proof.* Let  $\mathbf{u}'$  and  $z'$  be the exact derivatives of  $\mathcal{U}$  and  $\mathcal{Z}$  evaluated at the nodes, respectively. Due to the accuracy of  $H$ , the result will follow if we can show that

$$(z, \mathbf{u}')_H = (z, \mathbf{D}\mathbf{u})_H, \quad i, j \leq r, i + j \leq 2p. \quad (5)$$

First take  $j \leq p$ , this gives

$$\mathbf{D}\mathbf{u} = j\mathbf{x}^{j-1} = \mathbf{u}', \quad (6)$$

which means  $(z, \mathbf{u}')_H = (z, \mathbf{D}\mathbf{u})_H$  for  $j \leq p$ . Next, we consider  $j \geq p$ , which means that  $i \leq p$  and  $Dz = i\mathbf{x}^{i-1} = z'$ . Using  $Dz = z'$  along with the SBP property,  $Q + Q^T = E$ , gives

$$\begin{aligned} (z, \mathbf{D}\mathbf{u})_H &= z^T H \mathbf{D}\mathbf{u} \\ &= z^T (E - Q^T) \mathbf{u} \\ &= z^T E \mathbf{u} - z^T Q^T \mathbf{u} \\ &= z^T E \mathbf{u} - (\mathbf{u}, D\mathbf{z})_H \\ &= z^T E \mathbf{u} - (\mathbf{u}, z')_H \\ &= z^T E \mathbf{u} - \int_0^1 \mathcal{U} \frac{d\mathcal{Z}}{dx} dx, \end{aligned}$$

or, taking  $i, j \leq r$  and using the accuracy condition on  $E$  gives,

$$\begin{aligned} (z, \mathbf{D}\mathbf{u})_H &= z^T E \mathbf{u} - \int_0^1 \mathcal{U} \frac{d\mathcal{Z}}{dx} dx \\ &= \int_0^1 \frac{d(\mathcal{U}\mathcal{Z})}{dx} dx - \int_0^1 \mathcal{U} \frac{d\mathcal{Z}}{dx} dx \\ &= \int_0^1 \mathcal{Z} \frac{d\mathcal{U}}{dx} dx. \end{aligned}$$

Therefore,  $(z, \mathbf{u}')_H = (z, \mathbf{D}\mathbf{u})_H$  for  $j \geq p$  and  $i, j \leq r$ , with  $i + j \leq 2p$ . We have thus shown the desired result.  $\square$

From Theorems 1 and 2, we can see that, for generalized SBP operators,  $\mathbf{z}^T \mathbf{Q} \mathbf{u}$  is at least a  $2p$ -order approximation of  $\int_0^1 \mathcal{Z} \frac{d\mathcal{U}}{dx} dx$  if and only if  $r \geq 2p$ . This is a more stringent condition compared to Definition 1 where the accuracy requirement on  $r$  is  $r \geq p$ . The implication of Theorem 2 when considered along with the preceding one-dimensional example is that quadrature accuracy in curvilinear coordinates is decreased for generalized SBP operators if: (1)  $r < 2p$ , (2) the Jacobian of the transformation is approximated by the same SBP operator that is associated with the norm, and (3) an element uses a higher degree representation of the geometry compared to the degree of the SBP operator associated with that element.

The practical implication of this theoretical decrease in quadrature accuracy under the preceding conditions is the loss of superconvergent functionals in curvilinear coordinates under the above conditions. To appreciate this, we take a small detour and examine functional superconvergence *without* a curvilinear transformation for tensor-product SBP operators. We build upon [9], which explains functional superconvergence for classical tensor-product SBP discretizations, and [15], which extends the theory of superconvergent linear functionals to generalized SBP time-marching methods.

Consider the one-dimensional problem

$$\begin{aligned} \frac{d\mathcal{U}(x)}{dx} &= \mathcal{F}(x) \quad \forall x \in \Omega = [0, 1], \\ \mathcal{U}(x=0) &= \mathcal{U}_L, \end{aligned} \quad (7)$$

and the associated discretization

$$\mathbf{D} \mathbf{u}_h = \mathbf{f} - \overbrace{\mathbf{H}^{-1} \mathbf{t}_L (\mathbf{t}_L^T \mathbf{u}_h - \mathcal{U}_L)}^{\text{SAT}}. \quad (8)$$

We introduce the linear functional,

$$\mathcal{I}(\mathcal{U}) = \int_{\Omega} \mathcal{G} \mathcal{U} dx + \alpha(\mathcal{U})|_{x=1}, \quad (9)$$

and the discrete approximation of the preceding functional,

$$\mathcal{I}_h(\mathbf{u}_h) = \mathbf{g}^T \mathbf{H} \mathbf{u}_h + \alpha \mathbf{t}_R^T \mathbf{u}_h. \quad (10)$$

Here,  $\alpha$  is a scalar constant. The dual problem associated with this PDE and linear functional is [9]

$$\begin{aligned} -\frac{d\psi(x)}{dx} &= \mathcal{G}(x) \quad \forall x \in \Omega = [0, 1], \\ \psi(x=1) &= \alpha. \end{aligned} \quad (11)$$

The question we are interested in is: how well does  $\mathcal{I}_h(\mathbf{u}_h)$  approximate  $\mathcal{I}(\mathcal{U})$ ? To answer this question we require the following two lemmas.

**Lemma 1.** *The term  $\mathbf{t}_R^T \mathbf{u}_h$  is a degree  $\tau$  and order  $\tau + 1$  approximation of  $\mathcal{U}|_{x=1}$ .*

*Proof.* We have

$$\begin{aligned} \mathcal{U}|_{x=1} - \mathbf{t}_R^T \mathbf{u}_h &= \int_{\Omega} \frac{d\mathcal{U}}{dx} dx + \mathcal{U}|_{x=0} - \mathbf{t}_R^T \mathbf{u}_h && \text{(expand } \mathcal{U}|_{x=1}) \\ &= \int_{\Omega} \mathcal{F} dx + \mathcal{U}_L - \mathbf{t}_R^T \mathbf{u}_h && \text{(definition of PDE)} \\ &= \mathbf{1}^T \mathbf{H} \mathbf{f} + \mathbf{1}^T \mathbf{t}_L \mathcal{U}_L - \mathbf{t}_R^T \mathbf{u}_h + \mathcal{O}(h^{\tau+1}), && \text{(insert H)} \end{aligned}$$

since  $\mathbf{1}^T \mathbf{t}_L = 1$  and, using the accuracy of  $\mathbf{H}$ ,

$$\mathbf{1}^T \mathbf{H} \mathbf{f} = \int_{\Omega} \mathcal{F} dx + \mathcal{O}(h^{\tau+1}).$$

Also, from the definition of the discretization of the PDE we have

$$\mathbf{H} \mathbf{f} + \mathbf{t}_L \mathcal{U}_L = (\mathbf{Q} + \mathbf{t}_L \mathbf{t}_L^T) \mathbf{u}_h$$

and using the summation-by-parts (SBP) property gives

$$\mathbf{H}\mathbf{f} + \mathbf{t}_L \mathcal{U}_L = (-\mathbf{Q}^T + \mathbf{t}_R \mathbf{t}_R^T) \mathbf{u}_h.$$

Therefore,

$$\begin{aligned} \mathcal{U}|_{x=1} - \mathbf{t}_R^T \mathbf{u}_h &= \mathbf{1}^T (\mathbf{H}\mathbf{f} + \mathbf{t}_L \mathcal{U}_L) - \mathbf{t}_R^T \mathbf{u}_h + O(h^{\tau+1}) && \text{(collect terms)} \\ &= \mathbf{1}^T (\mathbf{Q} + \mathbf{t}_L \mathbf{t}_L^T) \mathbf{u}_h - \mathbf{t}_R^T \mathbf{u}_h + O(h^{\tau+1}) && \text{(definition of discrete PDE)} \\ &= \mathbf{1}^T (-\mathbf{Q}^T + \mathbf{t}_R \mathbf{t}_R^T) \mathbf{u}_h - \mathbf{t}_R^T \mathbf{u}_h + O(h^{\tau+1}), && \text{(SBP property)} \end{aligned}$$

or, since  $\mathbf{1}^T \mathbf{t}_R = 1$  and  $\mathbf{Q}\mathbf{1} = \mathbf{0}$ ,

$$\mathcal{U}|_{x=1} - \mathbf{t}_R^T \mathbf{u}_h = O(h^{\tau+1}).$$

This concludes the proof.  $\square$

**Remark 3.** Lemma 1 is essentially a variant of Theorem 4 in [15].

**Lemma 2.** Suppose we have  $\psi \frac{d\mathcal{U}}{dx} \in C^{\min(2p, \tau)}$ , where  $\psi$  and  $\mathcal{U}$  are the dual and primal solutions, respectively, associated with the PDE and functional defined by Eq. (7) and Eq. (9), respectively. Then,

$$\boldsymbol{\psi}^T \mathbf{H}(\mathbf{D}\mathbf{u} + \mathbf{H}^{-1} \mathbf{t}_L (\mathbf{t}_L^T \mathbf{u} - \mathcal{U}_L)) - \alpha \mathbf{t}_R^T \mathbf{u} + \alpha \mathcal{U}_R$$

is a degree  $\min(2p, \tau)$  and order  $\min(2p + 1, \tau + 1)$  approximation of

$$\int_{\Omega} \psi \frac{d\mathcal{U}}{dx} dx,$$

i.e.,

$$\boldsymbol{\psi}^T \mathbf{H}(\mathbf{D}\mathbf{u} + \mathbf{H}^{-1} \mathbf{t}_L (\mathbf{t}_L^T \mathbf{u} - \mathcal{U}_L)) - \alpha \mathbf{t}_R^T \mathbf{u} + \alpha \mathcal{U}_R = \int_{\Omega} \psi \frac{d\mathcal{U}}{dx} dx + O(h^{\min(2p+1, \tau+1)}).$$

*Proof.* We approach this proof in the same manner as the proof of Theorem 2. From the accuracy of  $\mathbf{H}$  we have

$$\int_{\Omega} \psi \frac{d\mathcal{U}}{dx} dx = (\boldsymbol{\psi}, \mathbf{u}')_{\mathbf{H}} + O(h^{\tau+1}). \quad (12)$$

Here,  $\mathbf{u}'$  is  $\frac{d\mathcal{U}}{dx}$  at the nodes. Based on the accuracy of  $\mathbf{H}$ , we would like to show

$$(\boldsymbol{\psi}, \mathbf{u}')_{\mathbf{H}} = \boldsymbol{\psi}^T \mathbf{H}(\mathbf{D}\mathbf{u} + \mathbf{H}^{-1} \mathbf{t}_L (\mathbf{t}_L^T \mathbf{u} - \mathcal{U}_L)) - \alpha \mathbf{t}_R^T \mathbf{u} + \alpha \mathcal{U}_R + O(h^{\min(2p+1, \tau+1)}). \quad (13)$$

Using a similar argument as in [12], it is sufficient to show that the preceding equation is exact for polynomial integrands of degree less than  $2p + 1$ . To this end, we consider  $\boldsymbol{\psi} = \mathbf{p}_k$  and  $\mathbf{u} = \mathbf{p}_m$  as degree  $k$  and  $m$  polynomials, respectively, where  $k + m \leq 2p + 1$  defines the highest permissible degree of the combined integrand. We begin by taking  $m \leq p$ , which gives

$$\boldsymbol{\psi}^T \mathbf{H}(\mathbf{D}\mathbf{u} + \mathbf{H}^{-1} \mathbf{t}_L (\mathbf{t}_L^T \mathbf{u} - \mathcal{U}_L)) - \alpha \mathbf{t}_R^T \mathbf{u} + \alpha \mathcal{U}_R = \boldsymbol{\psi}^T \mathbf{H}\mathbf{u}' = (\boldsymbol{\psi}, \mathbf{u}')_{\mathbf{H}}.$$

Next, we reverse the situation and take  $m > p$ , which means we must have  $k < p + 1$  due to the condition that  $k + m \leq 2p + 1$ . For this case we have

$$\begin{aligned} \boldsymbol{\psi}^T \mathbf{H}(\mathbf{D}\mathbf{u} + \mathbf{H}^{-1} \mathbf{t}_L (\mathbf{t}_L^T \mathbf{u} - \mathcal{U}_L)) - \alpha \mathbf{t}_R^T \mathbf{u} + \alpha \mathcal{U}_R &= \boldsymbol{\psi}^T (\mathbf{E} - \mathbf{Q}^T) \mathbf{u} + \boldsymbol{\psi}^T \mathbf{t}_L (\mathbf{t}_L^T \mathbf{u} - \mathcal{U}_L) - \alpha \mathbf{t}_R^T \mathbf{u} + \alpha \mathcal{U}_R && \text{(SBP property)} \\ &= \boldsymbol{\psi}^T (\mathbf{t}_R \mathbf{t}_R^T - \mathbf{Q}^T) \mathbf{u} - \boldsymbol{\psi}^T \mathbf{t}_L \mathcal{U}_L - \alpha \mathbf{t}_R^T \mathbf{u} + \alpha \mathcal{U}_R && \text{(definition of E)} \\ &= \alpha \mathbf{t}_R^T \mathbf{u} - \boldsymbol{\psi}^T \mathbf{Q}^T \mathbf{u} - \boldsymbol{\psi}_L \mathcal{U}_L - \alpha \mathbf{t}_R^T \mathbf{u} + \boldsymbol{\psi}_R \mathcal{U}_R && (\mathbf{t}_R^T \boldsymbol{\psi} = \boldsymbol{\psi}_R = \alpha) \\ &= -(\mathbf{D}\boldsymbol{\psi})^T \mathbf{H}\mathbf{u} + \boldsymbol{\psi}_R \mathcal{U}_R - \boldsymbol{\psi}_L \mathcal{U}_L \\ &= -(\boldsymbol{\psi}')^T \mathbf{H}\mathbf{u} + \boldsymbol{\psi} \mathcal{U}|_{x=0}^{x=1} \\ &= - \int_{\Omega} \frac{d\boldsymbol{\psi}}{dx} \mathcal{U} dx + \boldsymbol{\psi} \mathcal{U}|_{x=0}^{x=1} + O(h^{\tau+1}) && \text{(accuracy of H)} \\ &= \int_{\Omega} \boldsymbol{\psi} \frac{d\mathcal{U}}{dx} dx + O(h^{\tau+1}) && \text{(integ. by parts)} \end{aligned}$$

Since  $k + m \leq 2p + 1$ , we see that  $\boldsymbol{\psi}^T \mathbf{H}(\mathbf{D}\mathbf{u} + \mathbf{H}^{-1} \mathbf{t}_L (\mathbf{t}_L^T \mathbf{u} - \mathcal{U}_L)) - \alpha \mathbf{t}_R^T \mathbf{u} + \alpha \mathcal{U}_R$  is an order  $\min(2p + 1, \tau + 1)$  approximation of  $\int_{\Omega} \boldsymbol{\psi} \frac{d\mathcal{U}}{dx} dx$ , which concludes the proof.  $\square$

**Remark 4.** Lemma 2 is similar to Lemma 12 in [15], however here we consider general  $\alpha$  rather than  $\alpha = 0$ .

We now return to our initial question and state the following theorem as an answer.

**Theorem 3.** Let  $\mathbf{u}_h$  be the discrete solution that satisfies Eq. (8). Then  $\mathcal{I}_h(\mathbf{u}_h)$  (defined by Eq. (10)) is an order  $\min(2p + 1, \tau + 1)$  approximation of  $\mathcal{I}(\mathcal{U})$  (defined by Eq. (9)).

*Proof.* Consider

$$\begin{aligned}
\mathcal{I}(\mathcal{U}) - \mathcal{I}_h(\mathbf{u}_h) &= \int_{\Omega} \mathcal{G}\mathcal{U} \, dx + \alpha(\mathcal{U})|_{x=1} - \mathbf{g}^T \mathbf{H} \mathbf{u}_h - \alpha \mathbf{t}_R^T \mathbf{u}_h \\
&= \int_{\Omega} \mathcal{G}\mathcal{U} \, dx - \mathbf{g}^T \mathbf{H} \mathbf{u}_h + \alpha((\mathcal{U})|_{x=1} - \mathbf{t}_R^T \mathbf{u}_h) \\
&= \int_{\Omega} \mathcal{G}\mathcal{U} \, dx - \mathbf{g}^T \mathbf{H} \mathbf{u}_h + O(h^{\tau+1}) && \text{(using Lemma 1)} \\
&= \mathbf{g}^T \mathbf{H}(\mathbf{u} - \mathbf{u}_h) + O(h^{\tau+1}). && \text{(accuracy of H)}
\end{aligned}$$

Next, we must determine the order of the term  $\mathbf{g}^T \mathbf{H}(\mathbf{u} - \mathbf{u}_h)$ . To do this, we first note that the discretization of the dual problem associated with the primal PDE is

$$-\mathbf{D}\boldsymbol{\psi}_h = \mathbf{g} - \mathbf{H}^{-1} \mathbf{t}_R (\mathbf{t}_R^T \boldsymbol{\psi}_h - \alpha). \quad (14)$$

Based on our PDEs, we can define the truncation error associated with the primal and dual problems, respectively, as

$$\begin{aligned}
e_u &= \mathbf{D}\mathbf{u} - \mathbf{u}' + \mathbf{H}^{-1} \mathbf{t}_L (\mathbf{t}_L^T \mathbf{u} - \mathcal{U}_L) \\
e_\psi &= -\mathbf{D}\boldsymbol{\psi} + \boldsymbol{\psi}' + \mathbf{H}^{-1} \mathbf{t}_R (\mathbf{t}_R^T \boldsymbol{\psi} - \alpha).
\end{aligned}$$

Multiplying  $e_u$  by  $\mathbf{H}$  and rearranging gives

$$\begin{aligned}
\mathbf{H}e_u - (\mathbf{Q} + \mathbf{t}_L \mathbf{t}_L^T) \mathbf{u} + (\mathbf{H}\mathbf{f} + \mathbf{t}_L \mathcal{U}_L) &= \mathbf{0} \\
\mathbf{H}e_u - \mathbf{A}(\mathbf{u} - \mathbf{u}_h) &= \mathbf{0},
\end{aligned}$$

where we have introduced  $\mathbf{A} \equiv \mathbf{Q} + \mathbf{t}_L \mathbf{t}_L^T$  (note that from the definition of the discretization of the primal problem we can write  $\mathbf{A}\mathbf{u}_h = \mathbf{H}\mathbf{f} + \mathbf{t}_L \mathcal{U}_L$ ). Adding  $\boldsymbol{\psi}_h^T \mathbf{0} = \boldsymbol{\psi}_h^T \mathbf{H}e_u - \boldsymbol{\psi}_h^T \mathbf{A}(\mathbf{u} - \mathbf{u}_h)$  to the discrete integral equation,  $\mathbf{g}^T \mathbf{H}(\mathbf{u} - \mathbf{u}_h)$ , gives

$$\begin{aligned}
\mathbf{g}^T \mathbf{H}(\mathbf{u} - \mathbf{u}_h) &= \mathbf{g}^T \mathbf{H}(\mathbf{u} - \mathbf{u}_h) + \boldsymbol{\psi}_h^T \mathbf{H}e_u - \boldsymbol{\psi}_h^T \mathbf{A}(\mathbf{u} - \mathbf{u}_h) \\
&= \mathbf{g}^T \mathbf{H}(\mathbf{u} - \mathbf{u}_h) - \boldsymbol{\psi}_h^T \mathbf{A} \mathbf{H}^{-1} \mathbf{H}(\mathbf{u} - \mathbf{u}_h) + \boldsymbol{\psi}_h^T \mathbf{H}e_u && \text{(insert } \mathbf{I} = \mathbf{H}^{-1} \mathbf{H}) \\
&= (\mathbf{g}^T - \boldsymbol{\psi}_h^T \mathbf{A} \mathbf{H}^{-1}) \mathbf{H}(\mathbf{u} - \mathbf{u}_h) + \boldsymbol{\psi}_h^T \mathbf{H}e_u \\
&= (\mathbf{g} - \mathbf{H}^{-1} \mathbf{A}^T \boldsymbol{\psi}_h)^T \mathbf{H}(\mathbf{u} - \mathbf{u}_h) + \boldsymbol{\psi}_h^T \mathbf{H}e_u \\
&= (\mathbf{g} - \mathbf{H}^{-1} (-\mathbf{Q} + \mathbf{t}_R \mathbf{t}_R^T) \boldsymbol{\psi}_h)^T \mathbf{H}(\mathbf{u} - \mathbf{u}_h) + \boldsymbol{\psi}_h^T \mathbf{H}e_u && (\mathbf{A}^T = -\mathbf{Q} + \mathbf{t}_R \mathbf{t}_R^T) \\
&= (\mathbf{D}\boldsymbol{\psi}_h + \mathbf{g} - \mathbf{H}^{-1} \mathbf{t}_R \mathbf{t}_R^T \boldsymbol{\psi}_h)^T \mathbf{H}(\mathbf{u} - \mathbf{u}_h) + \boldsymbol{\psi}_h^T \mathbf{H}e_u \\
&= (\mathbf{D}\boldsymbol{\psi}_h + \mathbf{g} - \mathbf{H}^{-1} \mathbf{t}_R (\mathbf{t}_R^T \boldsymbol{\psi}_h - \alpha))^T \mathbf{H}(\mathbf{u} - \mathbf{u}_h) - \alpha (\mathbf{H}^{-1} \mathbf{t}_R)^T \mathbf{H}(\mathbf{u} - \mathbf{u}_h) + \boldsymbol{\psi}_h^T \mathbf{H}e_u \\
&= \boldsymbol{\psi}_h^T \mathbf{H}e_u - \alpha (\mathbf{H}^{-1} \mathbf{t}_R)^T \mathbf{H}(\mathbf{u} - \mathbf{u}_h),
\end{aligned}$$

where the first term in the final line above is zero due to the definition of the discretization of the dual problem. Continuing by adding and subtracting  $\boldsymbol{\psi}^T \mathbf{H}e_u$ , we have

$$\begin{aligned}
\mathbf{g}^T \mathbf{H}(\mathbf{u} - \mathbf{u}_h) &= -\alpha (\mathbf{H}^{-1} \mathbf{t}_R)^T \mathbf{H}(\mathbf{u} - \mathbf{u}_h) + \boldsymbol{\psi}_h^T \mathbf{H}e_u + \boldsymbol{\psi}^T \mathbf{H}e_u - \boldsymbol{\psi}^T \mathbf{H}e_u \\
&= (\boldsymbol{\psi}_h - \boldsymbol{\psi})^T \mathbf{H}e_u + \boldsymbol{\psi}^T \mathbf{H}e_u - \alpha \mathbf{t}_R^T (\mathbf{u} - \mathbf{u}_h).
\end{aligned}$$

But, note that  $\mathbf{A}^T (\boldsymbol{\psi} - \boldsymbol{\psi}_h) = \mathbf{H}e_\psi$ , which means

$$\mathbf{g}^T \mathbf{H}(\mathbf{u} - \mathbf{u}_h) = -(\mathbf{A}^{-T} \mathbf{H}e_\psi)^T \mathbf{H}e_u + \boldsymbol{\psi}^T \mathbf{H}e_u - \alpha \mathbf{t}_R^T (\mathbf{u} - \mathbf{u}_h).$$

Furthermore, making the assumption that  $\|A^{-T}H\|_\infty \leq C$ , where  $C$  is a constant, gives

$$\mathbf{g}^T H(\mathbf{u} - \mathbf{u}_h) = \boldsymbol{\psi}^T H e_u - \alpha \mathbf{t}_R^T (\mathbf{u} - \mathbf{u}_h) + O(h^{\tau+1}).$$

Alternatively, we can substitute  $e_u$  into the above and write,

$$\begin{aligned} \mathbf{g}^T H(\mathbf{u} - \mathbf{u}_h) &= \boldsymbol{\psi}^T H e_u - \alpha \mathbf{t}_R^T (\mathbf{u} - \mathbf{u}_h) + O(h^{\tau+1}) \\ &= \boldsymbol{\psi}^T H (D\mathbf{u} + H^{-1} \mathbf{t}_L (\mathbf{t}_L^T \mathbf{u} - \mathcal{U}_L)) - \alpha \mathbf{t}_R^T (\mathbf{u} - \mathbf{u}_h) - \boldsymbol{\psi}^T H \mathbf{u}' + O(h^{\tau+1}) \\ &= \boldsymbol{\psi}^T H (D\mathbf{u} + H^{-1} \mathbf{t}_L (\mathbf{t}_L^T \mathbf{u} - \mathcal{U}_L)) - \alpha \mathbf{t}_R^T (\mathbf{u} - \mathbf{u}_h) - \int_\Omega \boldsymbol{\psi} \frac{d\mathcal{U}}{dx} dx + O(h^{\tau+1}). \end{aligned}$$

Furthermore, using Lemma 1, we can substitute  $\mathbf{t}_R^T \mathbf{u}_h = \mathcal{U}_R + O(h^{\tau+1})$  into the above and write

$$\begin{aligned} \mathbf{g}^T H(\mathbf{u} - \mathbf{u}_h) &= \boldsymbol{\psi}^T H (D\mathbf{u} + H^{-1} \mathbf{t}_L (\mathbf{t}_L^T \mathbf{u} - \mathcal{U}_L)) - \alpha \mathbf{t}_R^T \mathbf{u} + \alpha \mathcal{U}_R - \int_\Omega \boldsymbol{\psi} \frac{d\mathcal{U}}{dx} dx + O(h^{\tau+1}) \\ &= O(h^{\min(2p+1, \tau+1)}) + O(h^{\tau+1}) \\ &= O(h^{\min(2p+1, \tau+1)}) \end{aligned} \quad (\text{Lemma 2})$$

Therefore,

$$\mathcal{I}(\mathcal{U}) - \mathcal{I}_h(\mathcal{U}) = O(h^{\min(2p+1, \tau+1)}), \quad (15)$$

which concludes the proof.  $\square$

To summarize, thus far we have seen that we can prove the superconvergence of linear functionals *without* a curvilinear transformation for the PDE defined by Eq. (7) without invoking Theorem 2. This changes in Section 4 of [9] when Theorem 1 of the present paper (Lemma 3 of [9] and Theorem 2 of [12]) is invoked to prove Theorem 6 (involving a curvilinear transformation) of [9], which concerns the accuracy of quadratures of the form

$$\mathcal{I}_h(\mathbf{u}) = \mathbf{u}^T (H \otimes H) \mathbf{J} \approx \iint_{\Omega_x} \mathcal{U} dx dy, \quad (16)$$

based on the uniform discretization of the computational domain. The proof of Theorem 6 is in Appendix A of [9], and it involves the repeated application of Theorem 1 to show that

$$\mathcal{I}_h(\mathbf{u}) = \mathbf{u}^T (H \otimes H) \mathbf{J} \quad (17)$$

is an order  $2p$  approximation of

$$\mathcal{I}(\mathcal{U}) = \iint_{\Omega_x} \mathcal{U} dx dy. \quad (18)$$

Because Theorem 1 does not apply to generalized SBP operators with  $r < 2p$ , we do not expect superconvergent functionals when: (1)  $r < 2p$ , (2) the Jacobian of the transformation is approximated by the same SBP operator that is associated with the norm, and (3) an element uses a higher degree representation of the geometry compared to the degree of the SBP operator associated with that element.

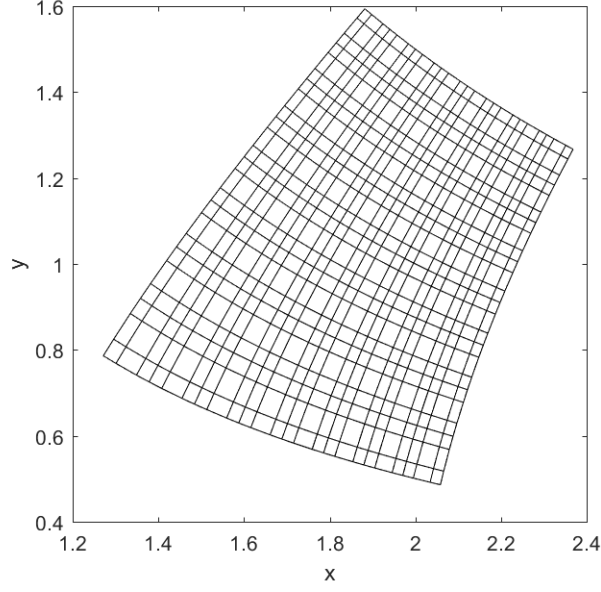
## IV. Results

We can confirm the quadrature accuracy of different generalized SBP operators numerically by examining a two-dimensional quadrature on a curvilinear domain. We take the test problem from Section 4.2 of [12] and restate it here for completeness. Consider the domain

$$\Omega = \{(x, y) \in \mathbb{R}^2 \mid 1 \leq xy \leq 3, 1 \leq x^2 - y^2 \leq 4\},$$

and the integral

$$\begin{aligned} \ell &= \iint_\Omega (x^2 + y^2) e^{\frac{1-x^2+y^2}{3}} \sin\left(\frac{xy-1}{2}\right) dx dy \\ &= 3(1 - e^{-1})(1 - \cos(1)). \end{aligned} \quad (19)$$



**Fig. 1** Example grid for  $\Omega$  with  $n_e = 8$  non-overlapping elements in the  $\tilde{x}$  and  $\tilde{y}$  directions. The  $p = 3$  LGL nodes are used in each element.

We can compute this integral numerically by introducing the global mapping coordinates

$$\tilde{x} = \frac{x^2 - y^2 - 1}{3} \quad \text{and} \quad \tilde{y} = \frac{xy - 1}{2}.$$

For a given number of non-overlapping elements,  $n_e$ , with  $n + 1$  nodes in the  $\tilde{x}$  and  $\tilde{y}$  directions in each element, we partition the square  $[0, 1]^2$  with  $n_e$  elements in the  $\tilde{x}$  and  $\tilde{y}$  directions. Figure 1 shows the grid for  $n_e = 8$  using the  $p = 3$  Legendre-Gauss-Lobatto (LGL) nodes in each element. We note here that although  $\tilde{x} = \tilde{x}(x, y)$  and  $\tilde{y} = \tilde{y}(x, y)$  are polynomial functions,  $x = x(\tilde{x}, \tilde{y})$  and  $y = y(\tilde{x}, \tilde{y})$  are not. This means that the geometry representation in each element is not a polynomial; however the geometry in each element corresponds to the analytical geometry.

On the  $i^{\text{th}}$  element, the Jacobian of the transformation is constructed as

$$\mathbf{J} = [(\mathbf{l} \otimes \mathbf{D})\mathbf{x}_i] \circ [(\mathbf{D} \otimes \mathbf{l})\mathbf{y}_i] - [(\mathbf{l} \otimes \mathbf{D})\mathbf{y}_i] \circ [(\mathbf{D} \otimes \mathbf{l})\mathbf{x}_i], \quad (20)$$

where  $\otimes$  denotes the Kronecker product and  $\circ$  denotes the Hadamard product. For a given  $n_e$ , the approximation of Eq. (19) is summed over all elements to obtain

$$\ell_{n_e} = \sum_{i=1}^{n_e} \sum_{j=1}^{n_e} \mathbf{J}^T (\mathbf{H} \otimes \mathbf{H}) \mathbf{f}, \quad (21)$$

where  $\mathbf{f}$  is the integrand of Eq. (19) computed using the  $\mathbf{x}$  and  $\mathbf{y}$  coordinates for each element. The error associated with the quadrature approximation is computed as  $E_{n_e} = |\ell - \ell_{n_e}|$ . Table 1 lists the generalized SBP operators that are used throughout this paper. Note that the LGL nodal distributions include the boundary nodes; and therefore the projection operators associated with this class of nodal distributions are exact (i.e.,  $r = \infty$ ). Details regarding the construction of the operators listed in Table 1 can be found in [11].

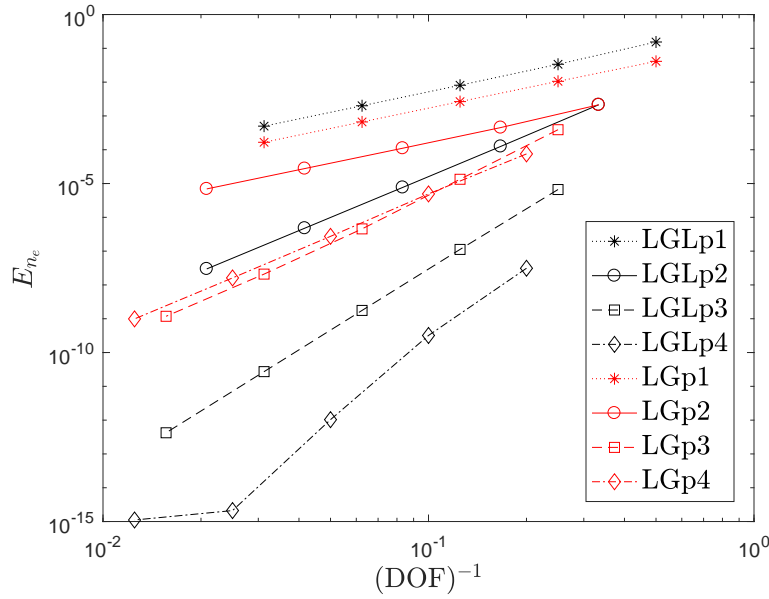
Figure 2 plots  $E_{n_e}$  as a function of  $(\text{DOF})^{-1}$  for the operators listed in Table 1, and Table 2 gives the associated convergence rates. Here,  $\text{DOF} = n_e(n + 1)$ , i.e., DOF is the square root of the total number of degrees of freedom used to discretize the domain  $\Omega_{\mathbf{x}}$ . The convergence rates were computed using data points from the three finest grid levels, except for the LGLp4 operator. The data point on the finest grid level for the LGLp4 operator was ignored due to round-off error and the slope was instead computed using the three preceding data points.

As expected based on Theorem 2, the quadratures computed using the LGL operators, which have  $r = \infty > 2p$ , all converge at a rate of approximately order  $\tau + 1$ , where  $\tau \geq 2p - 1$  denotes the degree of the quadrature rule associated



Operator	Nodal distribution	Operator degree, $p$	Projection degree, $r$	Quadrature degree, $\tau$
LGLp1	Legendre-Gauss-Lobatto	1	$\infty$	1
LGLp2	Legendre-Gauss-Lobatto	2	$\infty$	3
LGLp3	Legendre-Gauss-Lobatto	3	$\infty$	5
LGLp4	Legendre-Gauss-Lobatto	4	$\infty$	7
LGp1	Legendre-Gauss	1	1	3
LGp2	Legendre-Gauss	2	2	5
LGp3	Legendre-Gauss	3	3	7
LGp4	Legendre-Gauss	4	4	9

**Table 1** Some generalized summation-by-parts operators that satisfy Definition 1.



**Fig. 2** Convergence of the error when using the generalized SBP operators in Table 1 to approximate Eq. (19).

with each operator. These results agree with [11]. In contrast, the LG operators all converge at rates less than  $2p$ , which is also expected since for these operators  $r = p < 2p$ , the Jacobian is computed using the same SBP operator associated with the norm, and the geometry representation is non-polynomial. Specifically, there is an even-odd convergence pattern associated with the LG operators. The even-degree LG operators converge at a rate of  $p$  while the odd-degree LG operators converge at a rate of approximately  $p + 1$ . This even-odd quadrature convergence behaviour can be explained by considering the interactions between the leading truncation error terms associated with the respective even- and odd-degree LG projection operators. Consider decomposing  $E$  in terms of the projection operators  $t_L$  and  $tr$

$$E = t_R t_R^T - t_L t_L^T. \quad (22)$$

Next, recall the accuracy condition on  $E$  from Definition 1. Namely, for some  $\Omega_x \in [a, b]$ , we have

$$\left(x^i\right)^T E x^j = b^{i+j} - a^{i+j}, \quad i, j = 0, 1, \dots, r, r \geq p. \quad (23)$$

With Eq. (22) and Eq. (23) established, we take  $\Omega_x \in [-1, 1]$ . Substituting  $\Omega_x$ , i.e.,  $a = -1$  and  $b = 1$ , into Eq. (23) and moving all the terms to the right-hand side gives us an expression for the error in  $E$ ,  $E_E$ , for different values of  $i$  and  $j$ .

$$E_E = \left(x^i\right)^T E x^j - (1)^{i+j} + (-1)^{i+j}. \quad (24)$$

Operator	Convergence Rate
LGLp1	2.00994
LGLp2	4.00502
LGLp3	6.01197
LGLp4	8.60466
LGp1	1.99889
LGp2	1.99984
LGp3	4.29535
LGp4	4.05604

**Table 2** Convergence rates when using the generalized SBP operators in Table 1 to approximate Eq. (19).

LGp1	i	0	1	2						
	j	3	2	1						
	i+j	3	3	3						
LGp2	i	0	1	2	3	4				
	j	3	4	3	0	1				
	i+j	3	5	5	3	5				
LGp3	i	0	1	2	3	4	5	6		
	j	5	4	5	4	1	0	1		
	i+j	5	5	7	7	5	5	7		
LGp4	i	0	1	2	3	4	5	6	7	8
	j	5	6	5	6	5	0	1	0	1
	i+j	5	7	7	9	9	5	7	7	9

**Table 3** Values of  $i$ ,  $j$ , and  $i + j$  when the first nonzero value of  $E_E$  occurs. The values of  $i$  and  $j$  in Eq. (24) are increased with  $j$  running first.

For LG operators, based on Eq. (23), Eq. (24) will be equal to zero for  $i, j \leq r$ . Therefore, we are interested in the behaviour of Eq. (24) when  $i, j > r$ . Table 3 gives the values of  $i$ ,  $j$ , and  $i + j$  when the first nonzero value of  $E_E$  (i.e., Eq. (24)) occurs. For each operator, the minimum value of  $i + j$  is boxed. For the odd-degree LG operators, the minimum value of  $i + j$  is equal to  $p + 2$ . For the even-degree LG operators, the minimum value of  $i + j$  is equal to  $p + 1$ . To understand why this pattern occurs, we can examine Eq. (24) for the minimum values of  $i + j$  for each LG operator, as reported in Table 3.

From Table 3, the minimum value of  $i + j$  occurs for each operator when  $i = 0$ . Therefore, we substitute  $i = 0$  and Eq. (22) into Eq. (24), this gives

$$\begin{aligned}
 E_E|_{i=0} &= \left(x^0\right)^T \left(t_R t_R^T - t_L t_L^T\right) x^j - (1)^{0+j} + (-1)^{0+j} \\
 &= \mathbf{1}^T \left(t_R t_R^T - t_L t_L^T\right) x^j - (1)^j + (-1)^j \\
 &= t_R^T x^j - t_L^T x^j - 1 + (-1)^j,
 \end{aligned}$$

where we have used  $\mathbf{1}$  to denote a vector of ones. Note also that  $\mathbf{1}^T t_R = \mathbf{1}^T t_L = 1$  since the projection operators interpolate the constant function exactly. We now introduce two additional error metrics. Let

$$E_{t_R} = t_R^T x^j - 1, \quad (25)$$

$$E_{t_L} = t_L^T x^j - (-1)^j, \quad (26)$$

Operator	$j$	$E_E _{i=0}$	$E_{t_R}$	$E_{t_L}$	
LGp1	0	0.00000e + 00	-2.22045e - 16	-2.22045e - 16	
	$j = r$	1	-1.11022e - 15	-2.22045e - 16	2.22045e - 16
	$j = 2p$	2	0.00000e + 00	-6.66667e - 01	-6.66667e - 01
		3	-1.33333e + 00	-6.66667e - 01	6.66667e - 01
LGp2	0	0.00000e + 00	0.00000e + 00	0.00000e + 00	
	1	8.88178e - 16	0.00000e + 00	0.00000e + 00	
	$j = r$	2	0.00000e + 00	2.22045e - 16	2.22045e - 16
	3	-8.00000e - 01	-4.00000e - 01	4.00000e - 01	
	$j = 2p$	4	0.00000e + 00	-4.00000e - 01	-4.00000e - 01
LGp3	0	0.00000e + 00	0.00000e + 00	0.00000e + 00	
	1	0.00000e + 00	-5.55112e - 16	5.55112e - 16	
	2	0.00000e + 00	0.00000e + 00	0.00000e + 00	
	$j = r$	3	0.00000e + 00	-4.44089e - 16	4.44089e - 16
	4	0.00000e + 00	-2.28571e - 01	-2.28571e - 01	
	5	-4.57143e - 01	-2.28571e - 01	2.28571e - 01	
	$j = 2p$	6	0.00000e + 00	-4.24490e - 01	-4.24490e - 01
LGp4	0	2.22045e - 16	0.00000e + 00	2.22045e - 16	
	1	0.00000e + 00	0.00000e + 00	2.22045e - 16	
	2	1.11022e - 16	-2.22045e - 16	-2.22045e - 16	
	3	-4.44089e - 16	-2.22045e - 16	4.44089e - 16	
	$j = r$	4	1.11022e - 16	-1.11022e - 16	-1.11022e - 16
	5	-2.53968e - 01	-1.26984e - 01	1.26984e - 01	
	6	1.11022e - 16	-1.26984e - 01	-1.26984e - 01	
	7	-5.36155e - 01	-2.68078e - 01	2.68078e - 01	
$j = 2p$	8	1.11022e - 16	-2.68078e - 01	-2.68078e - 01	

**Table 4** Projection error associated with LG operators.

be the error associated with the projection operators  $t_R$  and  $t_L$ , respectively. Note that we can recast  $E_E|_{i=0}$  in terms of  $E_{t_R}$  and  $E_{t_L}$  as

$$E_E|_{i=0} = E_{t_R} - E_{t_L}. \quad (27)$$

Table 4 numerically tabulates these error terms for the LG operators listed in Table 1 for different  $j$ .

First note that all error terms are zero (or machine zero) up to  $j = r$ . This is expected from the accuracy condition on E associated with Definition 1. Therefore, we expect that  $E_{t_R}$  and  $E_{t_L}$  will be non-zero for  $j > r$ . The first non-zero (or non-machine-zero) values of  $E_{t_R}$  and  $E_{t_L}$  that appear when increasing  $j$  from 0 to  $2p$  are boxed, and this occurs for each operator when  $j = r + 1$ . Similarly, the first non-zero value of  $E_E|_{i=0}$  for each operator is boxed. For the even-degree operators, the first non-zero values of  $E_{t_R}$  and  $E_{t_L}$  are equal and opposite in sign, which results in the first non-zero value of  $E_E|_{i=0}$  occurring at  $j = r + 1$ . In contrast, the first non-zero values of  $E_{t_R}$  and  $E_{t_L}$  are equal and of the same sign, which causes them to cancel when  $j = r + 1$ . As a result of this cancellation, the first non-zero value of  $E_E|_{i=0}$  for the odd-degree operators occurs at  $j = r + 2$ , i.e., one value of  $j$  higher than for the even-degree operators. This explains the even-odd quadrature convergence behaviour observed in Figure 2 and Table 2 with the even-degree LG operators converging at a rate of  $p$  and the odd-degree LG operators converging at a rate of  $p + 1$ . Essentially, in curvilinear coordinates, when the Jacobian is approximated using the same SBP derivative operator associated with the norm and the geometry representation is either non-polynomial or of a degree higher than that of the SBP operator, the

<i>Polynomial Degree</i>	<i>Mesh Function</i>	<i>Abbreviation</i>
1	$\xi$	MFD1
2	$\frac{1}{2}\xi^2 + \frac{1}{2}\xi$	MFD2
3	$\frac{1}{3}\xi^3 + \frac{1}{3}\xi^2 + \frac{1}{3}\xi$	MFD3
4	$\frac{1}{4}\xi^4 + \frac{1}{4}\xi^3 + \frac{1}{4}\xi^2 + \frac{1}{4}\xi$	MFD4
5	$\frac{1}{5}\xi^5 + \frac{1}{5}\xi^4 + \frac{1}{5}\xi^3 + \frac{1}{5}\xi^2 + \frac{1}{5}\xi$	MFD5
non-polynomial	$(\exp(4\xi) - 1)/(\exp(4) - 1)$	MFNP

**Table 5 Mesh functions.**

LG quadrature accuracy is limited by the accuracy of the projection operators, and the even-odd quadrature convergence behaviour is associated with the leading truncation error cancellation of the projection operators that occurs for the odd-degree LG operators.

To examine functional superconvergence, consider the steady one-dimensional linear convection equation with unit wave speed on the domain  $\Omega_x \in [0, 1]$

$$\frac{\partial \mathcal{U}}{\partial x} = \mathcal{S}(x), \quad (28)$$

where  $\mathcal{S}(x)$  is a source term. The source term

$$\mathcal{S}(x) = \frac{\pi e^x}{e-1} \cos\left(\frac{\pi e^x - \pi + e - 1}{e-1}\right) \quad (29)$$

gives the steady-state solution

$$\mathcal{U}(x) = \sin\left(\frac{\pi(e^x - 1)}{e-1} + 1\right). \quad (30)$$

We use the mesh functions,  $x(\xi)$ , in Table 5 to introduce a non-constant metric Jacobian into integral functionals on  $\Omega_x$ , which simulates the effect of a curvilinear coordinate transformation.

Consider the integral functional

$$\mathcal{I}(\mathcal{U}) = \int_0^1 \mathcal{U} dx = \int_0^1 \mathcal{U} \frac{dx}{d\xi} d\xi, \quad (31)$$

discretized as

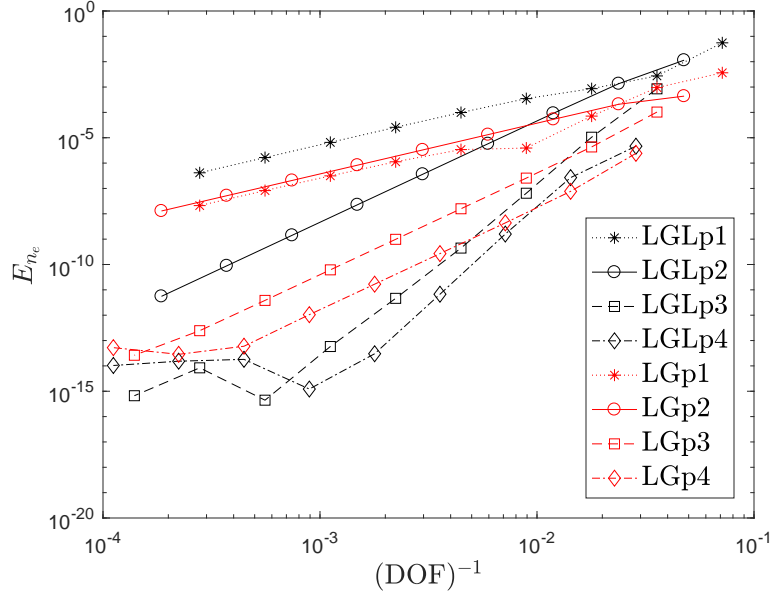
$$\mathcal{I}_{n_e}(\mathbf{u}) = \sum_{i=1}^{n_e} \mathbf{J}^T \mathbf{H} \mathbf{u}_i. \quad (32)$$

Table 6 gives the convergence results when using Eq. (32) to discretize Eq. (31) with the mesh functions in Table 5. Figure 3 shows the visual results when using the non-polynomial mesh function. The convergence rates were all computed using data from the asymptotic region prior to round-off error. Upwind SATs were used.

From Table 6, as for the two-dimensional quadrature example, the LGL operators converge at a rate of approximately the quadrature degree,  $\tau$ , plus one, i.e.,  $\tau + 1$ . In contrast, the LG operators converge at a rate of  $\tau + 1$  only when the mesh function is a polynomial whose degree is less than or equal to the degree of the SBP operator used to compute the Jacobian. When a mesh function is used that does not satisfy this condition, the odd- and even-degree LG operators converge at rates of  $p + 1$  and  $p$ , respectively. These convergence rates are boxed in Table 6. As before, the even-odd convergence behaviour can be explained by considering the error cancellation associated with the leading truncation error terms of the odd-degree LG projection operators.

## V. Conclusion

We have shown that, for tensor-product generalized SBP operators, projection operators of degree  $r \geq 2p$  are required to preserve quadrature accuracy and therefore superconvergent functionals in curvilinear coordinates when (1)



**Fig. 3** Convergence of the error when using the generalized SBP operators in Table 1 to approximate Eq. (31) via Eq. (32) using the non-polynomial mesh function.

<i>Operator</i>	MFD1	MFD2	MFD3	MFD4	MFD5	MFNP
LGLp1	2.00306	2.00912	1.97688	1.98138	1.98822	1.99188
LGLp2	4.00483	4.19932	4.00088	4.00030	4.00005	3.99982
LGLp3	6.00458	5.98032	6.05277	5.79955	5.65710	6.88729
LGLp4	8.40074	8.63669	8.01160	7.91618	7.82547	7.84233
LGp1	4.00515	2.00105	1.98979	1.97327	1.97019	1.94877
LGp2	6.00316	5.95041	1.99924	2.00003	2.00003	2.00004
LGp3	8.69240	7.96493	8.22565	3.96046	4.01958	4.01082
LGp4	9.48193	9.38194	10.39338	10.74542	3.99982	4.01261

**Table 6** Convergence rates when using the generalized SBP operators in Table 1 to approximate Eq. (31) via Eq. (32) with the mesh functions defined in Table 5.

the Jacobian of the transformation is approximated by the same SBP operator that is associated with the norm and (2) when a higher degree representation of the geometry is used compared to the degree of the SBP operator. Furthermore, when the geometry condition is violated for the LG SBP operators, which have  $r = p < 2p$ , there is an even-odd quadrature convergence behaviour that can be explained by considering the cancellation of the leading truncation error terms for the LG projection operators associated with the odd-degree LG operators.

### Acknowledgments

The authors gratefully acknowledge the financial support provided by the Natural Sciences and Engineering Research Council of Canada and the University of Toronto. D. A. Craig Penner thanks Dr. Pieter D. Boom for helpful discussions regarding some of the proofs in his Ph.D. thesis ([15]).

## References

- [1] Funaro, D., and Gottlieb, D., “A New Method of Imposing Boundary Conditions in Pseudospectral Approximations of Hyperbolic Equations,” *Mathematics of Computation*, Vol. 51, No. 184, 1988, pp. 599–613.
- [2] Carpenter, M. H., Gottlieb, D., and Abarbanel, S., “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, 1994, pp. 220–236.
- [3] Svärd, M., Carpenter, M. H., and Nordström, J., “A stable high-order finite difference scheme for the compressible Navier-Stokes equations, far-field boundary conditions,” *Journal of Computational Physics*, Vol. 225, No. 1, 2007, pp. 1020–1038.
- [4] Kreiss, H.-O., and Scherer, G., “Finite Element and Finite Difference Methods for Hyperbolic Partial Differential Equations,” *Mathematical Aspects of Finite Elements in Partial Differential Equations*, Academic Press, New York/London, 1974, pp. 195–212.
- [5] Strand, B., “Summation by Parts for Finite Difference Approximations for  $d/dx$ ,” *Journal of Scientific Computing*, Vol. 20, No. 1, 1994, pp. 47–67.
- [6] Del Rey Fernández, D. C., Hicken, J. E., and Zingg, D. W., “Review of summation-by-parts operators with simultaneous approximation terms for the numerical solution of partial differential equations,” *Computers & Fluids*, Vol. 95, 2014, pp. 171–196.
- [7] Svärd, M., and Nordström, J., “Review of summation-by-parts schemes for initial–boundary-value problems,” *Journal of Computational Physics*, Vol. 268, 2014, pp. 17–38.
- [8] Crean, J., Hicken, J. E., Del Rey Fernández, D. C., Zingg, D. W., and Carpenter, M. H., “Entropy-Stable Summation-By-Parts Discretization of the Euler Equations on General Curved Elements,” *Journal of Computational Physics*, Vol. 356, 2018, pp. 410–438.
- [9] Hicken, J. E., and Zingg, D. W., “Superconvergent functional estimates from summation-by-parts finite-difference discretizations,” *SIAM Journal on Scientific Computing*, Vol. 33, No. 2, 2011, pp. 893–922.
- [10] Hicken, J. E., and Zingg, D. W., “Dual consistency and functional accuracy: a finite-difference perspective,” *Journal of Computational Physics*, Vol. 256, 2014, pp. 161–182.
- [11] Del Rey Fernández, D. C., Boom, P. D., and Zingg, D. W., “A generalized framework for nodal first derivative summation-by-parts operators,” *Journal of Computational Physics*, Vol. 266, 2014, pp. 214–239.
- [12] Hicken, J. E., and Zingg, D. W., “Summation-by-parts operators and high-order quadrature,” *Journal of Computational and Applied Mathematics*, Vol. 237, No. 1, 2013, pp. 111–125.
- [13] Thomas, P. D., and Lombard, C. K., “Geometric conservation law and its application to flow computations on moving grids,” *AIAA Journal*, Vol. 17, No. 10, 1979, pp. 1030–1037.
- [14] Hicken, J. E., Del Rey Fernández, D. C., and Zingg, D. W., “Multidimensional summation-by-parts operators: General theory and application to simplex elements,” *SIAM Journal on Scientific Computing*, Vol. 38, No. 4, 2016, pp. A1935–A1958.
- [15] Boom, P. D., “High-order implicit time-marching methods for unsteady fluid flow simulation,” Ph.D. thesis, University of Toronto, 2015.