CrossMark

# On the exponent of exponential convergence of $p$-version FEM spaces

**Zhaonan Dong[1]** (ORCID)

## Abstract

We study the exponent of the exponential rate of convergence in terms of the number of degrees of freedom for various non-standard $p$-version finite element spaces employing reduced cardinality basis. More specifically, we show that serendipity finite element methods and discontinuous Galerkin finite element methods with total degree $\mathcal{P}_p$ basis have a faster exponential convergence with respect to the number of degrees of freedom than their counterparts employing the tensor product $\mathcal{Q}_p$ basis for quadrilateral/hexahedral elements, for piecewise analytic problems under $p$-refinement. The above results are proven by using a new $p$-optimal error bound for the $L^2$-orthogonal projection onto the total degree $\mathcal{P}_p$ basis, and for the $H^1$-projection onto the serendipity finite element space over tensor product elements with dimension $d \geq 2$. These new $p$-optimal error bounds lead to a larger exponent of the exponential rate of convergence with respect to the number of degrees of freedom. Moreover, these results show that part of the basis functions in $\mathcal{Q}_p$ basis plays no roles in achieving the $hp$-optimal error bound in the Sobolev space. The sharpness of theoretical results is also verified by a series of numerical examples.

**Keywords** $hp$-finite element method · Discontinuous Galerkin method · Serendipity basis · $\mathcal{P}_p$ basis · Reduced cardinality basis · Exponential convergence

**Mathematics Subject Classification (2010)** 65N30 · 65N15 · 65N50

## 1 Introduction

Polynomial approximation on tensor product domains plays an important role in deriving the exponential rate of convergence with respect to the number of degrees

✉ Zhaonan Dong
zd14@le.ac.uk

1    Department of Mathematics, University of Leicester, Leicester, UK

of freedom for $hp$-version finite element methods (FEMs) [4, 16–20, 24, 25, 29] and $hp$-version discontinuous Galerkin finite element methods (DGFEMs) [21, 22, 26–28, 31]. In general, the proof of the exponential rate of convergence usually depends on the $hp$-approximation results for some suitable projection operators onto a local polynomial space consisting of polynomials with degree less or equal than $p$ in each variable (known as $\mathcal{Q}_p$ basis) over a tensor product element (quadrilateral/hexahedral elements), for dimension $d \geq 2$.

The key reason for using the $\mathcal{Q}_p$ basis over a tensor product element is because $hp$-optimal approximation results for the multi-dimensional projection operators can be derived by using the stability and approximation results of the one-dimensional projections via tensor product arguments. On the other hand, the $hp$-approximation results for $L^2$-orthogonal projections onto polynomial basis with total degree less or equal than $p$ ($\mathcal{P}_p$ basis) and $H^1$-projections onto serendipity basis ($\mathcal{S}_p$ basis) have not been fully explored. Typically, $hp$-error bounds for projections onto the $\mathcal{P}_p$ or $\mathcal{S}_p$ basis have been derived using the fact that there exists a $q \leq p$ such that the bases $\mathcal{P}_p$ or $\mathcal{S}_p$ contain $\mathcal{Q}_q$ as a subset, together with the help of the $hp$-optimal approximation results for the projections onto the basis $\mathcal{Q}_q$ (see Corollary 4.52 in [29]).

For instance, we consider the $L^2$-norm error bound of the two-dimensional $L^2$-orthogonal projection $\Pi_{\mathcal{Q}_p}$ onto the $\mathcal{Q}_p$ basis as an example (cf. [11, 22]). Let $\hat{\kappa} = (-1, 1)^2$ and $u \in H^l(\hat{\kappa})$, $l$ is an integer with $l \geq 0$. Then, the following estimate holds,

$$\|u - \Pi_{\mathcal{Q}_p} u\|^2_{L^2(\hat{\kappa})} \leq C(s)(p + 1)^{-2s} |u|^2_{H^s(\hat{\kappa})}, \tag{1}$$

where the constant $C(s)$ is independent of $p$ and $0 \leq s \leq \min\{p + 1, l\}$. It is straightforward to see that the above error bound is sharp in the sense that it is $p$-optimal in both Sobolev regularity index $l$ and polynomial approximation order $p$.

Next, we consider the $L^2$-norm error bound of $L^2$-orthogonal projection $\Pi_{\mathcal{P}_p}$ onto the $\mathcal{P}_p$ basis. Following the Lemma 6 in [2], we define $\Pi_{\mathcal{P}_p} = \Pi_{\mathcal{Q}_{\lfloor p/2 \rfloor}}$, with $\lfloor p/2 \rfloor$ denoting the largest integer which is less than or equal to $p/2$. Then, the following bound holds:

$$\|u - \Pi_{\mathcal{P}_p} u\|^2_{L^2(\hat{\kappa})} = \|u - \Pi_{\mathcal{Q}_{\lfloor p/2 \rfloor}} u\|^2_{L^2(\hat{\kappa})} \leq \tilde{C}(s)(\lfloor p/2 \rfloor + 1)^{-2s} |u|^2_{H^s(\hat{\kappa})}, \tag{2}$$

where the constant $\tilde{C}(s)$ is independent of $p$ and $0 \leq s \leq \min\{\lfloor p/2 \rfloor + 1, l\}$. We emphasize that for function $u \in H^l(\hat{\kappa})$, with $p$ sufficiently large, the above error bound is $p$-optimal because $s = l$. However, if function $u$ is sufficiently smooth or even analytic, then the above error bound is $p$-suboptimal by at least $p/2$ orders because $s = \lfloor p/2 \rfloor + 1$. The similar $p$-suboptimal error bound holds for $H^1$-projections onto $\mathcal{S}_p$ basis.

Using the $p$-suboptimal error bound for $L^2$-orthogonal projections onto the $\mathcal{P}_p$ basis and $H^1$-projections onto the $\mathcal{S}_p$ basis, it is possible to derive an exponential rate of convergence for $hp$-FEMs employing the $\mathcal{S}_p$ basis and $hp$-DGFEMs employing the $\mathcal{P}_p$ basis, but the resulting exponent is much smaller with respect to the number of degrees of freedom than the exponent of FEMs and DGFEMs employing the

$Q_p$ basis. This contradicts the numerical observation in work [8–10, 13], where it is observed that the error with respect to the number of degrees of freedom for DGFEMs with the $\mathcal{P}_p$ basis on tensor product elements has a steeper exponential convergence compared to DGFEMs with the $\mathcal{Q}_p$ basis, for sufficiently smooth problems. This situation has been numerically tested on many different examples. We also observed numerically that the ratio of the slope of the exponential error decay for DGFEMs with the $\mathcal{P}_p$ basis compared to that of the $\mathcal{Q}_p$ basis depends only on the space dimension. The same phenomenon is also observed when comparing conforming FEMs with the $\mathcal{S}_p$ basis and the $\mathcal{Q}_p$ basis.

The disagreement between the numerical observations and theoretical results implies that the error bound (2) is not a sharp bound for $\mathcal{P}_p$ and $\mathcal{S}_p$ bases. To address this, in this work, we derive an *hp*-optimal error bound for the $L^2$-orthogonal projection onto the $\mathcal{P}_p$ basis in the $L^2$-norm, and for the $H^1$-projection onto the $\mathcal{S}_p$ basis in the $L^2$-norm and $H^1$-seminorm.

The technique for proving the new error bounds is different from the existing techniques for *hp*-approximation with the $\mathcal{Q}_p$ basis, due to the lack of a tensor product structure in the $\mathcal{P}_p$ and $\mathcal{S}_p$ bases, thereby hindering the use of the usual tensor product arguments. The key tools used in this work are: a multi-dimensional orthogonal polynomial expansion and the careful selection of basis functions. To the best author's knowledge, the new error bounds for both projections never appeared in the literatures. The resulting bounds are *hp-optimal* with respect to both Sobolev regularity and polynomial approximation order. Moreover, it also shows that the $\mathcal{Q}_p$ basis contains in a sense "extra" basis functions that are unnecessary for optimal convergence. These basis functions do not increase the order in *p* of the error bound, but instead only reduce its "constant".

By using the new *hp*-optimal error bound for the $L^2$-orthogonal projection onto the $\mathcal{P}_p$ basis and the $H^1$-projection onto the $\mathcal{S}_p$ basis, we can prove that methods using $\mathcal{P}_p$ and $\mathcal{S}_p$ bases offer exponential convergence with a larger exponent with respect to the number of degrees of freedom than comparable methods using $\mathcal{Q}_p$ basis for piecewise analytic problem under *p*-refinement. Furthermore, the approximation results also show that there are a lot of basis functions in $\mathcal{Q}_p$ basis with no roles in improving the *hp*-optimal error bound, which can be generalized to other FEM with the local polynomial space employing reduced cardinality basis. Finally, we emphasize that we are using DGFEM employing $\mathcal{P}_p$ basis for quadrilateral and hexahedral elements also and this is the key novelty of the approach, since this is possible for DGFEM and essentially for serendipity spaces.

The remainder of this work is structured as follows. In Section 2, we introduce the required notation and the weighted Sobolev spaces together with some properties about the orthogonal polynomials. Then, the *p*-optimal error bound for the $L^2$-orthogonal projection onto the $\mathcal{P}_p$ basis in $L^2$-norm is proved in Section 3. In Section 4, we derive the *p*-optimal error bound for $H^1$-projection onto the $\mathcal{S}_p$ basis in both $L^2$-norm and $H^1$-seminorm. Section 5 is devoted to deriving the exponential rate of convergence for the $L^2$- and the $H^1$-projections employing different local polynomial bases. The sharpness of the approximation results is verified through a series of numerical examples in Section 6.

## 2 Preliminaries

### 2.1 Notation

We employ the multi-indices $i = (i_1, \ldots, i_d)$ and $\alpha = (\alpha_1, \ldots, \alpha_d)$, where each component is non-negative. We denote by $|\cdot|$ the $l_1$-norm of the multi-index $i$, with $|i| = \sum_{j=1}^{d} |i_k|$. Further, for multi-indices, the relation $i \geq \alpha$ means that $i_k \geq \alpha_k$ for all $k = 1, \ldots, d$.

Next, we define the following shorthand notation for the summations of indices. For multi-indices $i$ and $\alpha$ satisfying $i \geq \alpha$, we define $\sum_{i \geq \alpha}^{\infty} := \sum_{i_1=\alpha_1}^{\infty} \cdots \sum_{i_d=\alpha_d}^{\infty}$ and the summation for multi-indices $i$ satisfying $|i| \geq p$ is defined as $\sum_{|i|=p}^{\infty}$. Moreover, we also define a summation for a multi-index $i$ satisfying multiple conditions, e.g. multi-index $i$ satisfying the condition $i \geq \alpha$ and the condition $|i| \geq p$ is defined as $\sum_{|i|=p, i \geq \alpha}^{\infty}$.

We introduce a function $\Phi_d(m, n)$ which will be used frequently in this work, given by

$$\Phi_d(m, n) = \left( \frac{\Gamma(\frac{m-n}{d} + 1)}{\Gamma(\frac{m+n}{d} + 1)} \right)^d, \tag{3}$$

where $\Gamma$ is the Gamma function satisfying $\Gamma(n + 1) = n!$ for integer $n \geq 0$.

### 2.2 Weighted Sobolev spaces

For the reference element $\hat{\kappa} := (-1, 1)^d$, let $W^\alpha(x) = \prod_{k=1}^{d} W_k(x_k)^{\alpha_k}$, where the weight function $W_k(x_k) := (1 - x_k^2)^{1/2}$, for $k = 1, \ldots, d$, and $\alpha_k \geq 0$ are integers.

Next, we define the weighted Sobolev spaces $V^l(\hat{\kappa})$ as a closure of $C^\infty(\hat{\kappa})$ in the norm with the weights $W^\alpha$, defined by

$$\|u\|_{V^l(\hat{\kappa})}^2 = \sum_{|\alpha|=0}^{l} |u|_{V^l(\hat{\kappa})}^2, \quad \text{and} \quad |u|_{V^l(\hat{\kappa})}^2 = \sum_{|\alpha|=l} \|W^\alpha D^\alpha u\|_{L^2(\hat{\kappa})}^2. \tag{4}$$

It is easy to see that $|u|_{V^l(\hat{\kappa})} \leq |u|_{H^l(\hat{\kappa})}, \forall u \in H^l(\hat{\kappa})$, with some integer $l \geq 0$. We note that the above definition for weighted Sobolev spaces can be extended to the fractional order weighted Sobolev spaces and weighted Besov spaces by using the real interpolation techniques (cf. [7]).

For $u \in L^2(\hat{\kappa})$, we introduce the Legendre polynomial expansion over the reference element $\hat{\kappa}$, given by $u(x) = \sum_{|i|=0}^{\infty} a_i \prod_{k=1}^{d} L_{i_k}(x_k)$, where $x = (x_1, \ldots, x_d)$, and $L_{i_k}(x_k)$ denotes the Legendre polynomial with order $i_k$ over the variable $x_k$. The coefficients $a_i$ are defined by

$$a_i = \int_{\hat{\kappa}} u(x) \prod_{k=1}^{d} \frac{2i_k + 1}{2} L_{i_k}(x_k) \, dx. \tag{5}$$

The derivatives of the function $u$ can be expressed as

$$D^\alpha u(x) = \sum_{i \geq \alpha}^{\infty} a_i \prod_{k=1}^{d} L_{i_k}^{(\alpha_k)}(x_k). \tag{6}$$

The derivatives of the Legendre polynomials satisfy the orthogonality property

$$\int_{-1}^{1} (1 - \xi^2)^k L_i^{(k)}(\xi) L_j^{(k)}(\xi) \, d\xi = \frac{2\delta_{ij}}{2i + 1} \frac{\Gamma(i + k + 1)}{\Gamma(i - k + 1)}, \tag{7}$$

see [29, Lemma 3.10]. By employing (7), we have

$$\|W^\alpha D^\alpha u\|_{L^2(\hat{\kappa})}^2 = \sum_{i \geq \alpha}^{\infty} |a_i|^2 \prod_{k=1}^{d} \frac{2}{2i_k + 1} \frac{\Gamma(i_k + \alpha_k + 1)}{\Gamma(i_k - \alpha_k + 1)}. \tag{8}$$

Identity (8) establishes a link between the derivatives of the functions in the weighted $L^2$-norms and their Legendre polynomial expansions.

*Remark 1* The weighted Sobolev space in the above definition is a special case of the general Jacobi-weighted Sobolev spaces introduced in [5]. The key reason to introduce the Jacobi-weighted Sobolev spaces is to deal with the loss of orthogonality suffered by orthogonal polynomials in standard Sobolev spaces; the $L^2$-orthogonality is preserved in Jacobi-weighted Sobolev spaces. As we shall see in the forthcoming analysis, orthogonality plays a key role in deriving optimal error bounds in the polynomial order $p$.

## 3 The $L^2$-orthogonal projection operator onto the $\mathcal{P}_p$ basis

In this section, we derive an $hp$-optimal error bound for the $L^2$-orthogonal projection over the reference element $\hat{\kappa} := (-1, 1)^d$.

### 3.1 The $L^2$-orthogonal projection operator

For the reference element $\hat{\kappa}$, we define $\mathcal{P}_p(\hat{\kappa})$ and $\mathcal{Q}_p(\hat{\kappa})$ be the space of all polynomials with total degree less than or equal to $p$ and with separate degree less than or equal to $p$, respectively.

In order to distinguish the same projections onto spaces with different polynomial bases, we use subscripts to signify the basis type: we use $\Pi_{\mathcal{Q}_p} := \Pi_p^{(1)} \Pi_p^{(2)} \ldots \Pi_p^{(d)}$ to denote the $L^2$-projection onto $\mathcal{Q}_p$, which is constructed by using tensor product arguments together with the one-dimensional $L^2$-projection with respect to variable $x_k$, given by $\Pi_p^{(k)}$. On the other hand, the $L^2$-projection onto $\mathcal{P}_p$ is denoted by $\Pi_{\mathcal{P}_p}$.

First, we have the following $hp$-optimal approximation result for the $L^2$-orthogonal projection $\Pi_{\mathcal{Q}_p}$ (c.f. [22, Lemma 3.4]).

**Lemma 1** *Let $\hat{\kappa} = (-1, 1)^d$. Suppose that $u \in H^l(\hat{\kappa})$, for some interger $l \geq 0$. Let $\Pi_{\mathcal{Q}_p} u$ be the $L^2$-projection of $u$ onto $\mathcal{Q}_p(\hat{\kappa})$ with $p \geq 0$. Then, for any integer $s$, with $0 \leq s \leq \min\{p + 1, l\}$, and $W_k = W_k(x_k)$, we have:*

$$\|u - \Pi_{\mathcal{Q}_p} u\|_{L^2(\hat{\kappa})}^2 \leq \Phi_1(p + 1, s) \Big( \sum_{k=1}^{d} \|W_k^s D_k^s u\|_{\hat{\kappa}}^2 \Big) \leq \Phi_1(p + 1, s) |u|_{H^s(\hat{\kappa})}^2, \tag{9}$$

*where $\Phi_1(p+1, s)$ is defined in (3).*

*Proof* The result is proved by modifying the proof of Lemma 3.4 in [22]. Instead of using triangle inequality, we use the orthogonality and stability of the one-dimensional $L^2$-orthogonal projection, which leads to the error bound (9). □

We remark on the asymptotic behaviour of the gamma function. Making use of sharp double side inequalities for the gamma function (see Theorem 1.6. in [6]), for all positive real numbers $x \geq 1$, we have

$$\sqrt{2\pi} x^{x+\frac{1}{2}} e^{-x} \leq \Gamma(x+1) \leq e x^{x+\frac{1}{2}} e^{-x}, \tag{10}$$

and it follows

$$\Phi_d(p+1, s) \leq C(s)\left(\frac{d}{p+1}\right)^{2s}, \tag{11}$$

with $0 \leq s \leq \min\{p+1, l\}$ and $C(s)$ depending on the constant $s$ only. This implies that the error bound (9) is optimal in $p$ with respect to both the Sobolev regularity index $l$ and polynomial order $p$. In fact, by modifying the proof of Theorem 6.2 in [23], it is can be shown that the constant $C(s) = (\frac{e}{2})^{2s}$.

Next, we introduce a useful lemma which is the key tool in proving the optimal error bounds in $p$. The proof of the lemma is postponed until Section 3.2.

**Lemma 2** *Let $\xi = (\xi_1, \xi_2, \ldots, \xi_d)$ and $\rho = (\rho_1, \rho_2, \ldots, \rho_d)$ be two non-negative integer valued vectors with $\rho \geq \xi$, satisfying $|\rho| = M$, $|\xi| = m$ for $M, m \in \mathbb{N}$. Then, we have the (global) upper bound*

$$F(\xi, \rho) := \prod_{k=1}^{d} \frac{\Gamma(\rho_k - \xi_k + 1)}{\Gamma(\rho_k + \xi_k + 1)} \leq \Phi_d(M, m). \tag{12}$$

*Furthermore, the maximum value of $F(\xi, \rho)$ under the above constraints on $\xi$ and $\rho$ is attained at $\xi_k = m/d$, $\rho_k = M/d$, $k = 1, \ldots, d$.*

**Theorem 1** *Let $\hat{\kappa} = (-1, 1)^d$. Suppose that $u \in H^l(\hat{\kappa})$, for some integer $l \geq 0$. Let $\Pi_{\mathcal{P}_p} u$ be the $L^2$-projection of $u$ onto $\mathcal{P}_p(\hat{\kappa})$ with $p \geq 0$. Then, for any integer $s$, $0 \leq s \leq \min\{p+1, l\}$, we have:*

$$\|u - \Pi_{\mathcal{P}_p} u\|_{L^2(\hat{\kappa})}^2 \leq \Phi_d(p+1, s)|u|_{V^s(\hat{\kappa})}^2 \leq C(s)\left(\frac{d}{p+1}\right)^{2s}|u|_{H^s(\hat{\kappa})}^2. \tag{13}$$

*where $\Phi_d(p+1, s)$ is defined in (3).*

*Proof* Using the relation (7) , for any integer $s$, $0 \le s \le \min\{p+1, l\}$, we have

$$\|u - \Pi_{\mathcal{P}_p} u\|^2_{L^2(\hat{\kappa})} = \sum_{|i|=p+1}^{\infty} |a_i|^2 \prod_{k=1}^d \frac{2}{2i_k+1} \le \sum_{|\alpha|=s} \sum_{|i|=p+1, i \ge \alpha}^{\infty} |a_i|^2 \prod_{k=1}^d \frac{2}{2i_k+1}$$

$$= \sum_{|\alpha|=s} \sum_{|i|=p+1, i \ge \alpha}^{\infty} |a_i|^2 \Big( \prod_{k=1}^d \frac{2}{2i_k+1} \frac{\Gamma(i_k+\alpha_k+1)}{\Gamma(i_k-\alpha_k+1)} \Big)$$

$$\times \Big( \prod_{k=1}^d \frac{\Gamma(i_k-\alpha_k+1)}{\Gamma(i_k+\alpha_k+1)} \Big)$$

$$\le \Phi_d(p+1,s) \sum_{|\alpha|=s} \sum_{|i|=p+1, i \ge \alpha}^{\infty} |a_i|^2 \prod_{k=1}^d \frac{2}{2i_k+1} \frac{\Gamma(i_k+\alpha_k+1)}{\Gamma(i_k-\alpha_k+1)}$$

$$\le \Phi_d(p+1,s) \sum_{|\alpha|=s} \|W^\alpha D^\alpha u\|^2_{L^2(\hat{\kappa})}$$

$$= \Phi_d(p+1,s)|u|^2_{V^s(\hat{\kappa})} \le C(s)\Big(\frac{d}{p+1}\Big)^{2s} |u|^2_{H^s(\hat{\kappa})},$$

where in step 1, the index set is enlarged; indeed, some of the terms with multi-index $|i| \ge p+1$ have been used more than once; in step 3, we use Lemma 2, taking $\xi_k = \alpha_k \ge 0$, $\rho_k = i_k \ge 0$, $M = p+1$, $m = s$, together with the restriction $0 \le s \le \min\{p+1, l\}$; in step 4, we used (8) and in the last step, the bound holds from (11). □

*Remark 2* We point out that the above proof for the $L^2$-orthogonal projection $\Pi_{\mathcal{P}_p}$ on $d$-dimensional reference element is a natural extension of the proof for one-dimensional result, see [29] for details.

By comparing the $L^2$-norm bound (9) for the projection $\Pi_{\mathcal{Q}_p}$ and (13) for the projection $\Pi_{\mathcal{P}_p}$, it is easy to see that both bounds are $p$-optimal with respect to Sobolev regularity index $l$ and also for polynomial order $p$. Moreover, we can see that the bound in (13) will have a larger constant compared to the bound in (9), and this constant depends on the dimension $d$. This result will play a key role in deriving the exponential convergence for the $\mathcal{P}_p$ basis.

### 3.2 The proof of Lemma 2

The proof will be split into three steps.

**Step 1** The proof follows a constrained optimization procedure. We set,

$$L(\xi, \rho, \mu, \lambda) = F(\xi, \rho) + \mu(|\xi| - m) + \lambda(|\rho| - M), \tag{14}$$

and we calculate the stationary points. We consider the partial derivatives with respect to $\xi_k$ and $\rho_k$, $k = 1, \ldots, d$,

$$\frac{\partial L}{\partial \xi_k} = -\Big( \frac{\Gamma'(\rho_k - \xi_k + 1)}{\Gamma(\rho_k - \xi_k + 1)} + \frac{\Gamma'(\rho_k + \xi_k + 1)}{\Gamma(\rho_k + \xi_k + 1)} \Big) F(\xi, \rho) + \mu = 0,$$

and

$$\frac{\partial L}{\partial \rho_k} = \left( \frac{\Gamma'(\rho_k - \xi_k + 1)}{\Gamma(\rho_k - \xi_k + 1)} - \frac{\Gamma'(\rho_k + \xi_k + 1)}{\Gamma(\rho_k + \xi_k + 1)} \right) F(\xi, \rho) + \lambda = 0,$$

which satisfy the equations

$$\frac{\Gamma'(\rho_k - \xi_k + 1)}{\Gamma(\rho_k - \xi_k + 1)} = \frac{\mu - \lambda}{2F(\xi, \rho)} \quad \text{and} \quad \frac{\Gamma'(\rho_k + \xi_k + 1)}{\Gamma(\rho_k + \xi_k + 1)} = \frac{\mu + \lambda}{2F(\xi, \rho)}, \qquad (15)$$

with $k = 1, \ldots, d$, by using the fact that $F(\xi, \rho) > 0$. The right-hand sides of the two equations in (15) are independent of the index $k$. Moreover, the function $\phi(z) = \Gamma(z)' / \Gamma(z)$ is the so-called digamma function with the property that (see [1], (6.3.16))

$$\phi(z + 1) = -\gamma + \sum_{n=1}^{\infty} \frac{z}{n(n + z)} = -\gamma + \sum_{n=1}^{\infty} \left( \frac{1}{n} - \frac{1}{n + z} \right), \quad z \neq -1, -2, \ldots,$$

where $\gamma$ is the Euler-Mascheroni constant. For $z \geq 0$, the function $\phi(z + 1)$ is a continuous monotonically increasing function, which shows that (15) have only one solution. This solution is $\tilde{\xi}_k = m/d$ and $\tilde{\rho}_k = M/d$, $k = 1, \ldots, d$, and the $F(\xi, \rho)$ will have the extreme value at this stationary point, given by

$$F(\tilde{\xi}, \tilde{\rho}) = \Phi_d(M, m). \qquad (16)$$

**Step 2** In order to find the global maximum, we need to prove the following asymptotic relationship:

$$\Phi_n(M, m) \leq \Phi_d(M, m), \qquad n = 1, \ldots, d - 1. \qquad (17)$$

This is proven by considering three different cases. We first consider the special case $m = 0$. In this case, (17) holds trivially. Next, we consider the case $m = \delta M$, with $0 < \delta < 1$. By using the property (10) of gamma functions, we have the following bound:

$$\frac{\Phi_d(M, m)}{\Phi_n(M, m)} \geq \left( \frac{\sqrt{2\pi}}{e} \right)^{d+n} \left( \frac{d}{n} \right)^{2\delta M} \left( \frac{1 - \delta}{1 + \delta} \right)^{\frac{d-n}{2}}. \qquad (18)$$

By recalling that $0 < \delta < 1$ and $n = 1, \ldots, d - 1$, we have that $0 < \frac{1-\delta}{1+\delta} < 1$ and the function $(\frac{d}{n})^{2\delta M}$ is monotonically increasing with respect to $M$. For $M \geq \left( (d + n) \log \left( \frac{e}{\sqrt{2\pi}} \right) + \frac{d-n}{2} \log \left( \frac{1+\delta}{1-\delta} \right) \right) \left( 2\delta \log \left( \frac{d}{n} \right) \right)^{-1}$, the above quotient formula is greater than 1 and therefore (17) holds.

Finally, we consider the case $m = M$. Using the same techniques used to derive (18) together with the fact that $\Gamma(1) = 1$, we have

$$\frac{\Phi_d(M, m)}{\Phi_n(M, m)} = \frac{(\Gamma(\frac{2M}{n} + 1))^n}{(\Gamma(\frac{2M}{d} + 1))^d} \geq \frac{(\sqrt{2\pi})^n}{e^d} \left( \frac{d}{2M} \right)^{\frac{d-n}{2}} \left( \frac{d}{n} \right)^{2M + \frac{n}{2}}. \qquad (19)$$

By using the fact that exponentially increasing functions grow faster than polynomials, we know that for sufficiently large $M$ the right-hand side of (19) is greater than 1 and therefore (17) holds.

**Step 3** Finally, we need to show that the extreme value (16) is the global maximum value of $F(\xi, \rho)$ under the constraints $|\xi| = m$ and $|\rho| = M$.

First, we can see that the function $F(\xi, \rho)$ is symmetric and continuous with respect to $\xi$ and $\rho$. The constraints $|\xi| = m$ and $|\rho| = M$ restrict the domain of $\xi$ and $\rho$ to be a $(d - 1)$-dimensional simplex, which is convex and compact. So the maximum value of the function $F(\xi, \rho)$ over the domain will be obtained only at the boundary of the domain or the stationary point of $F(\xi, \rho)$. We have calculated the function value at the stationary point in (16) already, so now we just need to check the function values on the boundary of the domain.

This may be proved by induction. We start with the case $d = 2$: the domain of $\xi$ and $\rho$ satisfying the constrains are two straight lines, $\rho_1 + \rho_2 = M$ and $\xi_1 + \xi_2 = m$. Here, the stationary point is the mid-point of each of the two lines $\tilde{\xi} = (m/2, m/2)$, $\tilde{\rho} = (M/2, M/2)$, and the boundary of the domain consists of the points $\xi^b = (0, m)$, $\rho^b = (0, M)$ or $\xi^b = (m, 0)$, $\rho^b = (M, 0)$, due to the constraints $\rho \geq \xi$. Using the symmetry of the function and of the domain, we know that at the two boundary points of the domain, $F(\xi, \rho)$ will attain the same value, with $F(\xi^b, \rho^b) = \Phi_1(M, m)$. By using the asymptotic relation (17), we find

$$F(\xi^b, \rho^b) = \Phi_1(M, m) \leq \Phi_2(M, m) = F(\tilde{\xi}, \tilde{\rho}).$$

The above relation shows that the extreme value (16) is the global maximum value under the constraints for $d = 2$.

Next, we consider the case $d = 3$, where the domain of each of $\xi$ and $\rho$ will be a triangle. In this case, the stationary point of $F(\xi, \rho)$ is when $\xi$ and $\rho$ are located at the barycentre of their respective triangle. The boundary of each domain consists of three straight lines. We need to calculate the maximum value of $F(\xi, \rho)$ on the boundary of the domain. By using the symmetry of $F(\xi, \rho)$, and that fact that $|\xi| = m$ and $|\rho| = M$, we only need to consider one part of domain boundary where $\xi_3 = 0$ and $\rho_3 = 0$. Then, the maximum of $F(\xi, \rho)$ on the domain boundary can be viewed as exactly the same problem with the same constraints as in the case $d = 2$. Consequently, the maximum value of $F(\xi, \rho)$ along the boundary of the domain is $F(\xi^b, \rho^b) = \Phi_2(M, m)$. Again, by using the same techniques as for $d = 2$, we deduce that

$$F(\xi^b, \rho^b) = \Phi_2(M, m) \leq \Phi_3(M, m) = F(\tilde{\xi}, \tilde{\rho}).$$

The above relation shows that the extreme value (16) is the global maximum value under the constraints for $d = 3$. For the general $d$-dimensional case, the proof can be carried out in a similar way. The key observation is that the maximum value of $F(\xi, \rho)$ on the boundary of $d$-dimensional domain will be at the stationary point of $F(\xi, \rho)$ on the $(d - 1)$-dimensional domain. By using the relation

$$\Phi_{d-1}(M, m) \leq \Phi_d(M, m),$$

the proof is complete.

## 4 The $H^1$-projection operator onto the $\mathcal{S}_p$ basis

In this section, we shall consider the $H^1$-projection over the reference element $\hat{\kappa} := (-1, 1)^d$ with $d = 2, 3$. Since the three-dimensional results depend on the two-dimensional results, we start with the two-dimensional case.

### 4.1 The $H^1$-projection operator on the reference square

First, we introduce the two-dimensional serendipity finite element space (cf. [29])

$$\mathcal{S}_p(\hat{\kappa}) := \mathcal{P}_p(\hat{\kappa}) + \mathrm{span}\{x_1^p x_2, x_1 x_2^p\}, \quad p \geq 1. \tag{20}$$

We can see in Fig. 1 that the serendipity space $\mathcal{S}_p$ contains two more basis functions than the $\mathcal{P}_p$ basis for $p \geq 2$. Another way to define the serendipity basis is to consider the decomposition of the $C^0$ finite element space with $\mathcal{Q}_p$ basis over $\hat{\kappa}$. For polynomial order $p$, the $\mathcal{S}_p$ basis has the same number of nodal basis functions and edge basis functions as the $\mathcal{Q}_p$ basis, but the $\mathcal{S}_p$ basis only has internal moment basis functions (those with zero value along the element boundary $\partial\hat{\kappa}$) whose total degree is less than or equal $p$ (cf. [29, 30]). We note that serendipity FEMs can be defined in a dimension-independent fashion (see [3]).

Similarly to the case of the $L^2$-projection, we use $\pi_{\mathcal{Q}_p} := \pi_p^{(1)} \pi_p^{(2)}$ to denote the $H^1$-projection onto the $\mathcal{Q}_p$ basis, which can be constructed via a tensor product of one-dimensional $H^1$-projection with respect to variable $x_k$, given by $\pi_p^{(k)}$. Similarly, the $H^1$-projection onto the $\mathcal{S}_p$ basis is denoted by $\pi_{\mathcal{S}_p}$, which is defined in (25).
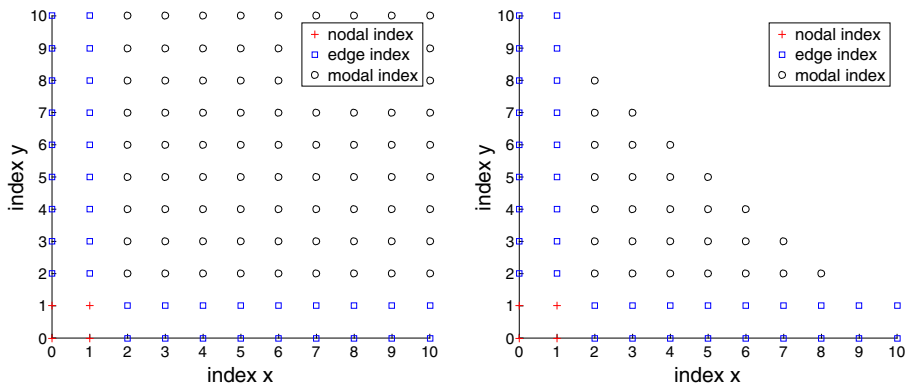


**Fig. 1** $\mathcal{Q}_p$ (left) and $\mathcal{S}_p$ (right) with polynomial order 10

Now, we construct the two-dimensional $H^1$-projection explicitly by using the one-dimensional $H^1$-projection and tensor product arguments (see [21, 29]). For $u \in H^l(\hat{\kappa})$, $l \geq 2$, the projection $\pi_{\mathcal{Q}_p} u \in \mathcal{Q}_p(\hat{\kappa})$, $p \geq 1$, is defined by

$$\pi_{\mathcal{Q}_p} u(x_1, x_2) := \int_{-1}^{x_1} \int_{-1}^{x_2} \Pi_{\mathcal{Q}_{p-1}} \partial_1 \partial_2 u(x_1, x_2) \, dx_1 \, dx_2 + \int_{-1}^{x_1} \Pi_{p-1}^{(1)} \partial_1 u(x_1, -1) \, dx_1$$

$$+ \int_{-1}^{x_2} \Pi_{p-1}^{(2)} \partial_2 u(-1, x_2) \, dx_2 + u(-1, -1)$$

$$= \sum_{i_1=0}^{p-1} \sum_{i_2=0}^{p-1} a_{i_1 i_2} \psi_{i_1}(x_1) \psi_{i_2}(x_2) + \sum_{i_1=0}^{p-1} b_{i_1} \psi_{i_1}(x_1) + \sum_{i_2=0}^{p-1} c_{i_2} \psi_{i_2}(x_2)$$

$$+ u(-1, -1); \tag{21}$$

the projections $\Pi_{\mathcal{Q}_{p-1}}$ and $\Pi_{p-1}^{(k)}$ are the two-dimensional and one-dimensional $L^2$-orthogonal projections, respectively, the coefficients $a_{i_1 i_2}$, $b_{i_1}$ and $c_{i_2}$ are given by:

$$a_{i_1 i_2} = \frac{2i_1 + 1}{2} \frac{2i_2 + 1}{2} \int_{\hat{\kappa}} \partial_1 \partial_2 u(x_1, x_2) L_{i_1}(x_1) L_{i_2}(x_2) \, dx,$$

$$b_{i_1} = \frac{2i_1 + 1}{2} \int_{-1}^{1} \partial_1 u(x_1, -1) L_{i_1}(x_1) \, dx_1,$$

$$c_{i_2} = \frac{2i_2 + 1}{2} \int_{-1}^{1} \partial_2 u(-1, x_2) L_{i_2}(x_2) \, dx_2, \tag{22}$$

and the polynomial function $\psi_j(z) = \int_{-1}^{z} L_j(z) \, dz$ with degree $j + 1$, and satisfies $\psi_j(\pm 1) = 0$ for $j \geq 1$. Moreover, for $j \geq 1$, $\psi_j(z) = -\frac{1}{j(j+1)}(1 - z^2)L_j'(z)$ has the following properties (cf. [29]),

$$\int_{-1}^{1} \psi_j(z) \psi_k(z) \frac{1}{1 - z^2} \, dz = \frac{2\delta_{jk}}{j(j+1)(2i+1)}. \tag{23}$$

Next, we rearrange the relation (21) by separating the internal moment basis functions:

$$\pi_{\mathcal{Q}_p} u(x_1, x_2) := \sum_{i_1=1}^{p-1} \sum_{i_2=1}^{p-1} a_{i_1 i_2} \psi_{i_1}(x_1) \psi_{i_2}(x_2) + \sum_{i_1=0}^{p-1} a_{i_1 0} \psi_{i_1}(x_1) \psi_0(x_2) + u(-1, -1)$$

$$+ \sum_{i_2=1}^{p-1} a_{0 i_2} \psi_0(x_1) \psi_{i_2}(x_2) + \sum_{i_1=0}^{p-1} b_{i_1} \psi_{i_1}(x_1) + \sum_{i_2=0}^{p-1} c_{i_2} \psi_{i_2}(x_2), \tag{24}$$

so that the first double summation in (24) only contains the internal moment basis functions. From the definition of $\mathcal{S}_p$, $\pi_{\mathcal{S}_p}$ can be constructed by removing the internal moment basis functions with polynomial order greater than $p$ in $\pi_{\mathcal{Q}_p}$. More specifically, $\pi_{\mathcal{S}_p} u \in \mathcal{S}_p(\hat{\kappa})$, $p \geq 4$, is defined by

$$
\pi_{\mathcal{S}_p} u(x_1, x_2) := \sum_{\substack{|i|=2 \\ i_k \geq 1, k=1,2}}^{p-2} a_{i_1 i_2} \psi_{i_1}(x_1) \psi_{i_2}(x_2) + \sum_{i_1=0}^{p-1} a_{i_1 0} \psi_{i_1}(x_1) \psi_0(x_2) + u(-1,-1)
$$

$$
+ \sum_{i_2=1}^{p-1} a_{0 i_2} \psi_0(x_1) \psi_{i_2}(x_2) + \sum_{i_1=0}^{p-1} b_{i_1} \psi_{i_1}(x_1) + \sum_{i_2=0}^{p-1} c_{i_2} \psi_{i_2}(x_2). \quad (25)
$$

For $1 \leq p \leq 3$, the first term in (25) will vanish, because there are no internal moment basis functions for the serendipity basis in that case. In this work, we focus on the high-order polynomial cases, so we only consider the $H^1$-projection $\pi_{\mathcal{S}_p}$ for $p \geq 4$.

Next, we recall the following approximation lemma for $\pi_{\mathcal{Q}_p}$ from [21].

**Lemma 3** *Let $\hat{\kappa} = (-1, 1)^2$. Suppose that $u \in H^{l+1}(\hat{\kappa})$, for some $l \geq 1$. Let $\pi_{\mathcal{Q}_p} u$ be the $H^1$-projection of $u$ onto $\mathcal{Q}_p(\hat{\kappa})$ with $p \geq 1$. Then, we have*

$$
\pi_{\mathcal{Q}_p} u = u \quad \text{at the vertices of } \hat{\kappa}, \quad (26)
$$

*and the following error estimates hold:*

$$
\|u - \pi_{\mathcal{Q}_p} u\|_{L^2(\hat{\kappa})}^2 \leq \frac{2}{p(p+1)} \Phi_1(p, s) \left( \|\partial_1^{s+1} u\|_{L^2(\hat{\kappa})}^2 + 2\|\partial_2^{s+1} u\|_{L^2(\hat{\kappa})}^2 \right)
$$

$$
+ \frac{4}{p^2(p+1)^2} \Phi_1(p, s-1) \|\partial_1 \partial_2^s u\|_{L^2(\hat{\kappa})}^2
$$

$$
\leq C(s) \left( \frac{1}{p} \right)^{2s+2} |u|_{H^{s+1}(\hat{\kappa})}^2, \quad (27)
$$

*and*

$$
\|\nabla(u - \pi_{\mathcal{Q}_p} u)\|_{L^2(\hat{\kappa})}^2 \leq 2\Phi_1(p, s) \left( \|\partial_1^{s+1} u\|_{L^2(\hat{\kappa})}^2 + \|\partial_2^{s+1} u\|_{L^2(\hat{\kappa})}^2 \right)
$$

$$
+ \frac{8}{p(p+1)} \Phi_1(p, s-1) \left( \|\partial_1^s \partial_2 u\|_{L^2(\hat{\kappa})}^2 + \|\partial_1 \partial_2^s u\|_{L^2(\hat{\kappa})}^2 \right) \leq C(s) \left( \frac{1}{p} \right)^{2s} |u|_{H^{s+1}(\hat{\kappa})}^2, \quad (28)
$$

*for any integer $s$, $1 \leq s \leq \min\{p, l\}$.*

Then, we derive the *hp*-error bound for the $H^1$-projection $\pi_{\mathcal{S}_p}$ for $p \geq 4$.

**Theorem 2** *Let $\hat{\kappa} = (-1, 1)^2$. Suppose that $u \in H^{l+1}(\hat{\kappa})$, for some $l \geq 1$. Let $\pi_{\mathcal{S}_p} u$ be the $H^1$ projection of $u$ onto $\mathcal{S}_p(\hat{\kappa})$ with $p \geq 4$. Then, we have*

$$
\pi_{\mathcal{S}_p} u = u \quad \text{at the vertices of } \hat{\kappa}, \quad (29)
$$

*and for any integer s,* $1 \le s \le \min\{p, l\}$, *p sufficiently large, the following error estimates hold:*

$$
\begin{aligned}
\|u - \pi_{\mathcal{S}_p} u\|^2_{L^2(\hat{\kappa})} &\le \frac{4}{p(p+1)} \Phi_1(p, s) \left( \|\partial_1^{s+1} u\|^2_{L^2(\hat{\kappa})} + 2\|\partial_2^{s+1} u\|^2_{L^2(\hat{\kappa})} \right) \\
&\quad + \frac{8}{p^2(p+1)^2} \Phi_1(p, s-1) \|\partial_1 \partial_2^s u\|^2_{L^2(\hat{\kappa})} \\
&\quad + 72 \Phi_2(p+1, s+1) |\partial_1 \partial_2 u|^2_{V^{s-1}(\hat{\kappa})} \\
&\le C(s) \left( \frac{2}{p+1} \right)^{2s+2} |u|^2_{H^{s+1}(\hat{\kappa})},
\end{aligned}
\tag{30}
$$

*and*

$$
\begin{aligned}
\|\nabla(u - \pi_{\mathcal{S}_p} u)\|^2_{L^2(\hat{\kappa})} &\le 4 \Phi_1(p, s) \left( \|\partial_1^{s+1} u\|^2_{L^2(\hat{\kappa})} + \|\partial_2^{s+1} u\|^2_{L^2(\hat{\kappa})} \right) \\
&\quad + \frac{16}{p(p+1)} \Phi_1(p, s-1) \left( \|\partial_1^s \partial_2 u\|^2_{L^2(\hat{\kappa})} + \|\partial_1 \partial_2^s u\|^2_{L^2(\hat{\kappa})} \right) \\
&\quad + 24 \Phi_2(p, s) |\partial_1 \partial_2 u|^2_{V^{s-1}(\hat{\kappa})} \le C(s) \left( \frac{2}{p} \right)^{2s} |u|^2_{H^{s+1}(\hat{\kappa})}.
\end{aligned}
\tag{31}
$$

*Proof* The key observation is the fact that the serendipity basis $\mathcal{S}_p$ differs from the $\mathcal{Q}_p$ basis only at the internal moment basis functions which vanish along the boundary of $\hat{\kappa}$. Indeed, using (24) and (25), we have

$$
\left( \pi_{\mathcal{Q}_p} u - \pi_{\mathcal{S}_p} u \right)(x_1, x_2) = \sum_{\substack{|i|=p-1 \\ p-1 \ge i_k \ge 1, k=1,2}}^{2(p-1)} a_{i_1 i_2} \psi_{i_1}(x_1) \psi_{i_2}(x_2).
\tag{32}
$$

Using the fact that $\psi_j(\pm 1) = 0$ for $j \ge 1$, we deduce that $(\pi_{\mathcal{Q}_p} u - \pi_{\mathcal{S}_p} u)|_{\partial \hat{\kappa}} = 0$. Thus, (29) is proved.

Next, we derive (30). The first step is the use of the triangle inequality,

$$
\|u - \pi_{\mathcal{S}_p} u\|^2_{L^2(\hat{\kappa})} \le 2\|u - \pi_{\mathcal{Q}_p} u\|^2_{L^2(\hat{\kappa})} + 2\|\pi_{\mathcal{Q}_p} u - \pi_{\mathcal{S}_p} u\|^2_{L^2(\hat{\kappa})}.
\tag{33}
$$

Thus, we only need to consider the error from the second term in the above bound. By using the orthogonality relation (23) of $\psi_j(x)$ for $j \ge 1$ and $1 \le s \le \min\{p, l\}$, we have

$$
\begin{aligned}
\|\pi_{\mathcal{Q}_p} u - \pi_{\mathcal{S}_p} u\|^2_{L^2(\hat{\kappa})} &\le \|(\pi_{\mathcal{Q}_p} u - \pi_{\mathcal{S}_p} u) W^{-1}\|^2_{L^2(\hat{\kappa})} \\
&= \sum_{\substack{|i|=p-1 \\ p-1 \ge i_k \ge 1, k=1,2}}^{2(p-1)} |a_{i_1 i_2}|^2 \prod_{k=1}^{2} \frac{2}{2i_k + 1} \frac{1}{i_k(i_k+1)} \\
&\le \sum_{|\alpha|=s-1} \sum_{\substack{|i|=p-1, i \ge \alpha \\ i_k \ge 1, k=1,2}}^{\infty} |a_{i_1 i_2}|^2 \prod_{k=1}^{2} \frac{2}{2i_k + 1} \frac{1}{i_k(i_k+1)},
\end{aligned}
\tag{34}
$$

where in step 2, we enlarged the summation index sets by adding the high-order internal moment basis functions with coefficients $a_{i_1 i_2}$, $i_k \geq 1$ for $k = 1, 2$ and $|i| \geq p - 1$. Thus, we have

$$
\begin{aligned}
\|\pi_{\mathcal{Q}_p} u - \pi_{\mathcal{S}_p} u\|^2_{L^2(\hat{\kappa})} &\leq \sum_{|\alpha|=s-1} \sum_{\substack{|i|=p-1, i \geq \alpha \\ i_k \geq 1, k=1,2}}^{\infty} |a_{i_1 i_2}|^2 \Big( \prod_{k=1}^{2} \frac{2}{2i_k + 1} \frac{\Gamma(i_k + \alpha_k + 1)}{\Gamma(i_k - \alpha_k + 1)} \Big) \\
&\quad \times \Big( \prod_{k=1}^{2} \frac{1}{i_k(i_k + 1)} \frac{\Gamma(i_k - \alpha_k + 1)}{\Gamma(i_k + \alpha_k + 1)} \Big) \\
&\leq \sum_{|\alpha|=s-1} \sum_{\substack{|i|=p-1, i \geq \alpha \\ i_k \geq 1, k=1,2}}^{\infty} |a_{i_1 i_2}|^2 \Big( \prod_{k=1}^{2} \frac{2}{2i_k + 1} \frac{\Gamma(i_k + \alpha_k + 1)}{\Gamma(i_k - \alpha_k + 1)} \Big) \\
&\quad \times \Big( \prod_{k=1}^{2} \frac{\Gamma(i_k - \alpha_k + 1)}{\Gamma(i_k + \alpha_k + 3)} \Big) \times 36,
\end{aligned}
\tag{35}
$$

where we used $\frac{1}{i_k(i_k+1)} \leq \frac{6}{(i_k + \alpha_k + 1)(i_k + \alpha_k + 2)}$, since $i_k \geq \alpha_k$ and $i_k \geq 1$. Now, we have

$$
\begin{aligned}
\|\pi_{\mathcal{Q}_p} u - \pi_{\mathcal{S}_p} u\|^2_{L^2(\hat{\kappa})} &\leq \sum_{|\alpha|=s-1} \sum_{|i|=p-1, i \geq \alpha}^{\infty} |a_{i_1 i_2}|^2 \Big( \prod_{k=1}^{2} \frac{2}{2i_k + 1} \frac{\Gamma(i_k + \alpha_k + 1)}{\Gamma(i_k - \alpha_k + 1)} \Big) \\
&\quad \times \Big( \prod_{k=1}^{2} \frac{\Gamma(i_k - \alpha_k + 1)}{\Gamma(i_k + \alpha_k + 3)} \Big) \times 36 \\
&\leq 36 \Phi_2(p+1, s+1) \sum_{|\alpha|=s-1} \|W^\alpha D^\alpha (\partial_1 \partial_2 u)\|^2_{L^2(\hat{\kappa})} \\
&= 36 \Phi_2(p+1, s+1) |\partial_1 \partial_2 u|^2_{V^{s-1}(\hat{\kappa})} \\
&\leq C(s) \Big( \frac{2}{p+1} \Big)^{2s+2} |u|^2_{H^{s+1}(\hat{\kappa})};
\end{aligned}
\tag{36}
$$

in step 1, we enlarge the index set by adding functions with coefficients $a_{i_1 i_2}$ whose index satisfying the relation $|i| \geq p-1$, $\prod_{k=1}^{2} i_k = 0$, while in step 2, we use Lemma 2, with $\xi_1 = \alpha_1 + 1 \geq 1$, $\xi_2 = \alpha_2 + 1 \geq 1$, $\rho_1 = i_1 + 1 \geq 1$, $\rho_2 = i_2 + 1 \geq 1$, $M = p+1$, and $m = s+1$, together with the restriction $1 \leq s \leq \min\{p, l\}$; in step 3, we use (8) and (22) to build up the link between the derivatives of $u$ and coefficients $a_{i_1 i_2}$ and in the last step, we use (11).

Using the same techniques, we can derive the error estimate for the $H^1$-seminorm. We have

$$
\begin{aligned}
\|\partial_1 (\pi_{\mathcal{Q}_p} u - \pi_{\mathcal{S}_p} u)\|^2_{L^2(\hat{\kappa})} &\leq \|\partial_1 (\pi_{\mathcal{Q}_p} u - \pi_{\mathcal{S}_p} u) W_2^{-1}\|^2_{L^2(\hat{\kappa})} \\
&\leq \sum_{|\alpha|=s-1} \sum_{\substack{|i|=p-1, i \geq \alpha \\ i_k \geq 1, k=1,2}}^{\infty} |a_{i_1 i_2}|^2 \frac{1}{i_2(i_2 + 1)} \prod_{k=1}^{2} \frac{2}{2i_k + 1}
\end{aligned}
\tag{37}
$$

In the last step, we enlarge the summation index sets by adding the high-order internal moment basis functions with coefficients $a_{i_1 i_2}$, $i_k \geq 1$ for $k = 1, 2$ and $|i| \geq p - 1$. Thus, we have

$$\|\partial_1(\pi_{\mathcal{Q}_p} u - \pi_{\mathcal{S}_p} u)\|^2_{L^2(\hat{\kappa})} \leq \sum_{\substack{|\alpha|=s-1}} \sum_{\substack{|i|=p-1, i \geq \alpha \\ i_k \geq 1, k=1,2}}^{\infty} |a_{i_1 i_2}|^2 \Big( \prod_{k=1}^{2} \frac{2}{2i_k + 1} \frac{\Gamma(i_k + \alpha_k + 1)}{\Gamma(i_k - \alpha_k + 1)} \Big)$$

$$\times \Big( \frac{\Gamma(i_1 - \alpha_1 + 1)}{\Gamma(i_1 + \alpha_1 + 1)} \frac{\Gamma(i_2 - \alpha_2 + 1)}{\Gamma(i_2 + \alpha_2 + 3)} \Big) \times 6$$

$$\leq \sum_{\substack{|\alpha|=s-1}} \sum_{\substack{|i|=p-1, i \geq \alpha}}^{\infty} |a_{i_1 i_2}|^2 \Big( \prod_{k=1}^{2} \frac{2}{2i_k + 1} \frac{\Gamma(i_k + \alpha_k + 1)}{\Gamma(i_k - \alpha_k + 1)} \Big)$$

$$\times \Big( \frac{\Gamma(i_1 - \alpha_1 + 1)}{\Gamma(i_1 + \alpha_1 + 1)} \frac{\Gamma(i_2 - \alpha_2 + 1)}{\Gamma(i_2 + \alpha_2 + 3)} \Big) \times 6$$

$$\leq 6\Phi_2(p, s) \sum_{\substack{|\alpha|=s-1}} \|W^\alpha D^\alpha (\partial_1 \partial_2 u)\|^2_{L^2(\hat{\kappa})}$$

$$= 6\Phi_2(p, s)|\partial_1 \partial_2 u|^2_{V^{s-1}(\hat{\kappa})} \leq C(s) \Big( \frac{2}{p} \Big)^{2s} |u|^2_{H^{s+1}(\hat{\kappa})}, \quad (38)$$

where in step 2, we enlarge the index set again; in step 3, we use Lemma 2, taking $\xi_1 = \alpha_1 \geq 0$, $\xi_2 = \alpha_2 + 1 \geq 1$, $\rho_1 = i_1 \geq 0$, $\rho_2 = i_2 + 1 \geq 1$, $M = p$, and $m = s$, together with the restriction $1 \leq s \leq \min\{p, l\}$.

Therefore, we have the bound

$$\|\nabla(\pi_{\mathcal{Q}_p} u - \pi_{\mathcal{S}_p} u)\|^2_{L^2(\hat{\kappa})} \leq 12\Phi_2(p, s)|\partial_1 \partial_2 u|^2_{V^{s-1}(\hat{\kappa})} \leq C(s) \Big( \frac{2}{p} \Big)^{2s} |u|^2_{H^{s+1}(\hat{\kappa})} \quad (39)$$

Finally, using (36), (39) and Lemma 3, the bounds (30) and (31) follow.                                                      □

## 4.2 The $H^1$-projection operator on the reference cube

In this section, we shall consider the $H^1$-projection operator over the reference cube $\hat{\kappa} := (-1, 1)^3$. First, we introduce the 3D serendipity finite element space.

A simple way to define the serendipity basis is to consider a decomposition of the $C^0$ finite element space with $\mathcal{Q}_p$ basis over $\hat{\kappa}$. For polynomial order $p$, the $\mathcal{S}_p$ basis has the same number of nodal basis functions and edge basis functions as the $\mathcal{Q}_p$ basis, but the $\mathcal{S}_p$ basis only has face basis functions (those with zero value on 12 edges and eight vertices) and internal moment basis functions (those with zero value along the element boundary $\partial \hat{\kappa}$) whose total degree is less than or equal $p$. The number of basis functions of $\mathcal{S}_p$ basis is calculated in the following way

$$Dof(\mathcal{S}_p(\hat{\kappa})) := 8 + 12 \times (p - 1) + 6 \times \frac{(p - 2)(p - 3)}{2}$$

$$+ \frac{(p - 3)(p - 4)(p - 5)}{6}, \quad (40)$$

here, we note that for $p = 1$, the serendipity basis only contains eight nodal basis functions and $\mathcal{S}_1(\hat{\kappa}) := \mathcal{Q}_1(\hat{\kappa})$. For $p \geq 2$, the serendipity basis contains $(p-1)$ edge

basis functions for each of the 12 edges. For $p \geq 4$, the serendipity basis contains $(p-2)(p-3)/2$ face basis functions for each of the six faces. For $p \geq 6$, the serendipity basis contains $(p-3)(p-4)(p-5)/6$ internal moment basis functions.

Similarly to the 2D case, we use $\pi_{\mathcal{Q}_p} := \pi_p^{(1)} \pi_p^{(2)} \pi_p^{(3)}$ to denote the $H^1$-projection onto the $\mathcal{Q}_p$ basis. The $H^1$-projection onto the $\mathcal{S}_p$ basis is denoted by $\pi_{\mathcal{S}_p}$. Additionally, we introduce some new notation for the forthcoming analysis. The projection $\pi_{\mathcal{S}_p}^{(1,2)}$ shall denote the $H^1$-projection onto the serendipity spaces $\mathcal{S}_p$ with variables $(x_1, x_2)$ only, and the projections $\pi_{\mathcal{S}_p}^{(1,3)}$ and $\pi_{\mathcal{S}_p}^{(2,3)}$ are defined in an analogous manner.

First, we explicitly construct the three-dimensional projection $\pi_{\mathcal{Q}_p} = \pi_p^{(1)} \pi_p^{(2)} \pi_p^{(3)}$. For $u \in H^l(\hat{\kappa})$, $l \geq 3$, the projection $\pi_{\mathcal{Q}_p} u \in \mathcal{Q}_p(\hat{\kappa})$, $p \geq 1$, is defined by

$$
\pi_{\mathcal{Q}_p} u(x_1, x_2, x_3) := \int_{-1}^{x_1} \int_{-1}^{x_2} \int_{-1}^{x_3} \Pi_{\mathcal{Q}_{p-1}} \partial_1 \partial_2 \partial_3 u(x_1, x_2, x_3)\, dx_1\, dx_2\, dx_3
$$

$$
+ \int_{-1}^{x_1} \int_{-1}^{x_2} \Pi_{p-1}^{(1)} \Pi_{p-1}^{(2)} \partial_1 \partial_2 u(x_1, x_2, -1)\, dx_1\, dx_2
$$

$$
+ \int_{-1}^{x_1} \int_{-1}^{x_3} \Pi_{p-1}^{(1)} \Pi_{p-1}^{(3)} \partial_1 \partial_3 u(x_1, -1, x_3)\, dx_1\, dx_3
$$

$$
+ \int_{-1}^{x_2} \int_{-1}^{x_3} \Pi_{p-1}^{(2)} \Pi_{p-1}^{(3)} \partial_2 \partial_3 u(-1, x_2, x_3)\, dx_2\, dx_3
$$

$$
+ \int_{-1}^{x_1} \Pi_{p-1}^{(1)} \partial_1 u(x_1, -1, -1)\, dx_1
$$

$$
+ \int_{-1}^{x_2} \Pi_{p-1}^{(2)} \partial_2 u(-1, x_2, -1)\, dx_2
$$

$$
+ \int_{-1}^{x_3} \Pi_{p-1}^{(3)} \partial_3 u(-1, -1, x_3)\, dx_3 + u(-1, -1, -1).
$$

Then, the following Legendre polynomial expansion holds:

$$
\pi_{\mathcal{Q}_p} u(x_1, x_2, x_3) := \sum_{i_1=0}^{p-1} \sum_{i_2=0}^{p-1} \sum_{i_3=0}^{p-1} a_{i_1 i_2 i_3} \psi_{i_1}(x_1) \psi_{i_2}(x_2) \psi_{i_3}(x_3) + u(-1, -1, -1)
$$

$$
+ \sum_{i_1=0}^{p-1} \sum_{i_2=0}^{p-1} b_{i_1 i_2} \psi_{i_1}(x_1) \psi_{i_2}(x_2) + \sum_{i_1=0}^{p-1} \sum_{i_3=0}^{p-1} c_{i_1 i_3} \psi_{i_1}(x_1) \psi_{i_3}(x_3)
$$

$$
+ \sum_{i_2=0}^{p-1} \sum_{i_3=0}^{p-1} d_{i_2 i_3} \psi_{i_2}(x_2) \psi_{i_3}(x_3) + \sum_{i_1=0}^{p-1} e_{i_1} \psi_{i_1}(x_1) + \sum_{i_2=0}^{p-1} f_{i_2} \psi_{i_2}(x_2)
$$

$$
+ \sum_{i_3=0}^{p-1} g_{i_3} \psi_{i_3}(x_3), \tag{41}
$$

with coefficients $a_{i_1 i_2 i_3}$, $b_{i_1 i_2}$, $c_{i_1 i_3}$, $d_{i_2 i_3}$, given by

$$a_{i_1 i_2 i_3} = \frac{2i_1 + 1}{2} \frac{2i_2 + 1}{2} \frac{2i_3 + 1}{2} \int_{\hat{\kappa}} \partial_1 \partial_2 \partial_3 u(x_1, x_2, x_3) L_{i_1}(x_1) L_{i_2}(x_2) L_{i_3}(x_3) \, dx,$$

$$b_{i_1 i_2} = \frac{2i_1 + 1}{2} \frac{2i_2 + 1}{2} \int_{-1}^{1} \int_{-1}^{1} \partial_1 \partial_2 u(x_1, x_2, -1) L_{i_1}(x_1) L_{i_2}(x_2) \, dx_1 \, dx_2,$$

$$c_{i_1 i_3} = \frac{2i_1 + 1}{2} \frac{2i_3 + 1}{2} \int_{-1}^{1} \int_{-1}^{1} \partial_1 \partial_3 u(x_1, -1, x_3) L_{i_1}(x_1) L_{i_3}(x_3) \, dx_1 \, dx_3,$$

$$d_{i_2 i_3} = \frac{2i_2 + 1}{2} \frac{2i_3 + 1}{2} \int_{-1}^{1} \int_{-1}^{1} \partial_2 \partial_3 u(-1, x_2, x_3) L_{i_2}(x_2) L_{i_3}(x_3) \, dx_2 \, dx_3, \quad (42)$$

together with $e_{i_1}$, $f_{i_2}$ and $g_{i_3}$

$$e_{i_1} = \frac{2i_1 + 1}{2} \int_{-1}^{1} \partial_1 u(x_1, -1, -1) L_{i_1}(x_1) \, dx_1,$$

$$f_{i_2} = \frac{2i_2 + 1}{2} \int_{-1}^{1} \partial_2 u(-1, x_2, -1) L_{i_2}(x_2) \, dx_2,$$

$$g_{i_3} = \frac{2i_3 + 1}{2} \int_{-1}^{1} \partial_3 u(-1, -1, x_3) L_{i_3}(x_3) \, dx_3. \quad (43)$$

Now, we separate the face basis functions and internal moment basis functions from (41).

$$\begin{aligned}
\pi_{\mathcal{Q}_p} u(x_1, x_2, x_3) := & \sum_{i_1=1}^{p-1} \sum_{i_2=1}^{p-1} \sum_{i_3=1}^{p-1} a_{i_1 i_2 i_3} \psi_{i_1}(x_1) \psi_{i_2}(x_2) \psi_{i_3}(x_3) \\
& + \sum_{i_1=1}^{p-1} \sum_{i_2=1}^{p-1} \left( a_{i_1 i_2 0} \psi_{i_1}(x_1) \psi_{i_2}(x_2) \psi_0(x_3) + b_{i_1 i_2} \psi_{i_1}(x_1) \psi_{i_2}(x_2) \right) \\
& + \sum_{i_1=1}^{p-1} \sum_{i_3=1}^{p-1} \left( a_{i_1 0 i_3} \psi_{i_1}(x_1) \psi_0(x_2) \psi_{i_3}(x_3) + c_{i_1 i_3} \psi_{i_1}(x_1) \psi_{i_3}(x_3) \right) \\
& + \sum_{i_2=1}^{p-1} \sum_{i_3=1}^{p-1} \left( a_{0 i_2 i_3} \psi_0(x_1) \psi_{i_2}(x_2) \psi_{i_3}(x_3) + d_{i_2 i_3} \psi_{i_2}(x_2) \psi_{i_3}(x_3) \right) \\
& + \text{edge basis} + \text{nodal basis.} \quad (44)
\end{aligned}$$

Here, the first triple summation terms contains all the internal moment basis functions only. Three double summation terms contain all the face basis functions. The edge basis functions and nodal basis functions will not be written explicitly because they play no role in the analysis.

From the definition of $\mathcal{S}_p$, $\pi_{\mathcal{S}_p} u$ can be constructed by removing the face basis functions and internal moment basis functions with polynomial order greater than $p$ in $\pi_{\mathcal{Q}_p} u$. More specifically, $\pi_{\mathcal{S}_p} u \in \mathcal{S}_p(\hat{\kappa})$, $p \geq 6$, is defined by

$$
\begin{aligned}
\pi_{\mathcal{S}_p} u(x_1, x_2, x_3) := & \sum_{\substack{|i|=3 \\ i_k \geq 1, k=1,2,3}}^{p-3} a_{i_1 i_2 i_3} \psi_{i_1}(x_1) \psi_{i_2}(x_2) \psi_{i_3}(x_3) \\
& + \sum_{\substack{i_1+i_2=2 \\ i_1 \geq 1, i_2 \geq 1}}^{p-2} \left( a_{i_1 i_2 0} \psi_{i_1}(x_1) \psi_{i_2}(x_2) \psi_0(x_3) + b_{i_1 i_2} \psi_{i_1}(x_1) \psi_{i_2}(x_2) \right) \\
& + \sum_{\substack{i_1+i_3=2 \\ i_1 \geq 1, i_3 \geq 1}}^{p-2} \left( a_{i_1 0 i_3} \psi_{i_1}(x_1) \psi_0(x_2) \psi_{i_3}(x_3) + c_{i_1 i_3} \psi_{i_1}(x_1) \psi_{i_3}(x_3) \right) \\
& + \sum_{\substack{i_2+i_3 \geq 2 \\ i_2 \geq 1, i_3 \geq 1}}^{p-2} \left( a_{0 i_2 i_3} \psi_0(x_1) \psi_{i_2}(x_2) \psi_{i_3}(x_3) + d_{i_2 i_3} \psi_{i_2}(x_2) \psi_{i_3}(x_3) \right) \\
& + \text{edge basis} + \text{nodal basis}
\end{aligned}
\tag{45}
$$

For $1 \leq p \leq 3$, both face basis functions and internal moment basis functions in (45) will vanish. For $4 \leq p \leq 5$, internal moment basis functions in (45) will vanish. Similar to the 2D case, we only consider the $H^1$-projection $\pi_{\mathcal{S}_p}$ for $p \geq 6$.

Next, by using the stability and approximation results for one-dimensional $H^1$-projection in [21], we can derive the following approximation results for $\pi_{\mathcal{Q}_p}$.

**Lemma 4** *Let $\hat{\kappa} = (-1, 1)^3$. Suppose that $u \in H^{l+1}(\hat{\kappa})$, for some $l \geq 2$. Let $\pi_{\mathcal{Q}_p} u$ be the $H^1$-projection of $u$ onto $\mathcal{Q}_p(\hat{\kappa})$ with $p \geq 1$. Then, we have*

$$
\pi_{\mathcal{Q}_p} u = u \quad \text{at the vertices of } \hat{\kappa},
\tag{46}
$$

*and the following error estimates hold:*

$$
\begin{aligned}
\|u - \pi_{\mathcal{Q}_p} u\|_{L^2(\hat{\kappa})}^2 \leq & \frac{8}{p(p+1)} \Phi_1(p, s) \\
& \times \left( \|\partial_1^{s+1} u\|_{L^2(\hat{\kappa})}^2 + \|\partial_2^{s+1} u\|_{L^2(\hat{\kappa})}^2 + \|\partial_3^{s+1} u\|_{L^2(\hat{\kappa})}^2 \right) \\
& + \frac{8}{p^2(p+1)^2} \Phi_1(p, s-1) \left( \|\partial_1 \partial_2^s u\|_{L^2(\hat{\kappa})}^2 + \|\partial_1 \partial_3^s u\|_{L^2(\hat{\kappa})}^2 + \|\partial_2 \partial_3^s u\|_{L^2(\hat{\kappa})}^2 \right) \\
& + \frac{8}{p^3(p+1)^3} \Phi_1(p, s-2) \|\partial_1 \partial_2 \partial_3^{s-1} u\|_{L^2(\hat{\kappa})}^2 \leq C(s) \left( \frac{1}{p} \right)^{2s+2} |u|_{H^{s+1}(\hat{\kappa})}^2,
\end{aligned}
\tag{47}
$$

*and*

$$\|\nabla(u - \pi_{\mathcal{Q}_p} u)\|^2_{L^2(\hat{\kappa})} \leq 2\Phi_1(p, s)\Big(\|\partial_1^{s+1} u\|^2_{L^2(\hat{\kappa})} + \|\partial_2^{s+1} u\|^2_{L^2(\hat{\kappa})} + \|\partial_3^{s+1} u\|^2_{L^2(\hat{\kappa})}\Big)$$

$$+ \frac{8}{p(p+1)}\Phi_1(p, s-1)\Big(\|\partial_1 \partial_2^s u\|^2_{L^2(\hat{\kappa})} + \|\partial_2 \partial_3^s u\|^2_{L^2(\hat{\kappa})} + \|\partial_3 \partial_1^s u\|^2_{L^2(\hat{\kappa})}$$

$$+ \|\partial_1 \partial_3^s u\|^2_{L^2(\hat{\kappa})} + \|\partial_2 \partial_1^s u\|^2_{L^2(\hat{\kappa})} + \|\partial_3 \partial_2^s u\|^2_{L^2(\hat{\kappa})}\Big)$$

$$+ \frac{8}{p^2(p+1)^2}\Phi_1(p, s-2)\Big(\|\partial_1 \partial_2 \partial_3^{s-1} u\|^2_{L^2(\hat{\kappa})} + \|\partial_1 \partial_2 \partial_3^{s-1} u\|^2_{L^2(\hat{\kappa})}$$

$$+ \|\partial_1 \partial_2 \partial_3^{s-1} u\|^2_{L^2(\hat{\kappa})}\Big) \leq C(s)\Big(\frac{1}{p}\Big)^{2s} |u|^2_{H^{s+1}(\hat{\kappa})}, \tag{48}$$

*for any integer $s$, $2 \leq s \leq \min\{p, l\}$.*

Then, we derive the *hp*-error bound for the $H^1$-projection $\pi_{\mathcal{S}_p}$ for $p \geq 6$.

**Theorem 3** *Let $\hat{\kappa} = (-1, 1)^3$. Suppose that $u \in H^{l+1}(\hat{\kappa})$, for some $l \geq 2$. Let $\pi_{\mathcal{S}_p} u$ be the $H^1$ projection of $u$ onto $\mathcal{S}_p(\hat{\kappa})$ with $p \geq 6$. Then, we have*

$$\pi_{\mathcal{S}_p} u = u \quad \text{at the vertices of } \hat{\kappa}, \tag{49}$$

*and for any integer $s$, $2 \leq s \leq \min\{p, l\}$, $p$ sufficiently large, the following error estimates hold:*

$$\|u - \pi_{\mathcal{S}_p} u\|^2_{L^2(\hat{\kappa})} \leq 2\|u - \pi_{\mathcal{Q}_p} u\|^2_{L^2(\hat{\kappa})} + 2\|\pi_{\mathcal{Q}_p} u - \pi_{\mathcal{S}_p} u\|^2_{L^2(\hat{\kappa})}.$$

$$\leq C_1 \Phi_3(p+1, s+1)|u|^2_{H^{s+1}(\hat{\kappa})} \leq C(s)\Big(\frac{3}{p+1}\Big)^{2s+2}|u|^2_{H^{s+1}(\hat{\kappa})}, \tag{50}$$

*and*

$$\|\nabla(u - \pi_{\mathcal{S}_p} u)\|^2_{L^2(\hat{\kappa})} \leq 2\|\nabla(u - \pi_{\mathcal{Q}_p} u)\|^2_{L^2(\hat{\kappa})} + 2\|\nabla(\pi_{\mathcal{Q}_p} - \pi_{\mathcal{S}_p} u)\|^2_{L^2(\hat{\kappa})}$$

$$\leq C_2 \Phi_3(p, s)|u|^2_{H^{s+1}(\hat{\kappa})} \leq C(s)\Big(\frac{3}{p}\Big)^{2s}|u|^2_{H^{s+1}(\hat{\kappa})}. \tag{51}$$

*Here, $C_1$ and $C_2$ are positive constants independent of $p$, $l$ and $s$.*

*Proof* See the proof of Theorem 4.4 in [14]. □

*Remark 3* We again make a comparison between the bounds in the $L^2$-norm and $H^1$-seminorm, given in Lemma 3 for $d = 2$ and Lemma 4 for $d = 3$ respectively for $\pi_{\mathcal{Q}_p}$, and Theorem 2 for $d = 2$ and Theorem 3 for $d = 3$ respectively for $\pi_{\mathcal{S}_p}$. Similarly to the comparisons for the $L^2$-projection onto $\mathcal{P}_p$ and $\mathcal{Q}_p$, both bounds are $p$-optimal in both Sobolev regularity and polynomial order. We can also see that the bounds for $\pi_{\mathcal{S}_p}$ have a larger constant than those for $\pi_{\mathcal{Q}_p}$, and this constant depends on dimension $d$. Moreover, we point out that the optimal approximation results for the $H^1$-projection with $\mathcal{S}_p$ basis in Theorems 2 and 3 directly imply the *hp*-optimal error bound for the $L^2$-norm on the trace of $\hat{\kappa}$ for $\pi_{\mathcal{S}_p}$.

*Remark 4* We note that in the Theorems 2 and 3, the minimum Sobolev regularity requirement for defining $H^1$-projection is $u \in H^d(\hat{\kappa})$ for the reference element. In fact, this regularity requirement can be relaxed by using the the tensor product Sobolev spaces (cf. [15, 29]). In this work, we do not consider the minimum regularity assumptions because we only consider the standard Sobolev spaces.

### 4.3 The $H^1$-projection operator onto the $\mathcal{P}_p$ basis

Finally, we present the error bound for $\pi_{\mathcal{P}_p}$ which we shall define now. The key observation is that the $\mathcal{P}_p$ basis with polynomial order $p$ contains the $\mathcal{S}_{p+1-d}$ basis for $p \geq d$, see [3]. Then, we can simply define $\pi_{\mathcal{P}_p} = \pi_{\mathcal{S}_{p+1-d}}$ for $d = 2, 3$.

**Corollary 1** *Let $\hat{\kappa} = (-1, 1)^d$, $d = 2, 3$. Suppose that $u \in H^{l+1}(\hat{\kappa})$, for some $l \geq d - 1$. Let $\pi_{\mathcal{P}_p} u := \pi_{\mathcal{S}_{p+1-d}} u$ be the $H^1$ projection of $u$ onto $\mathcal{P}_p(\hat{\kappa})$ with $p \geq 3d - 1$. Then, we have:*

$$\pi_{\mathcal{P}_p} u = u \quad \text{at the vertices of } \hat{\kappa}, \tag{52}$$

*and the following error estimates hold:*

$$\|u - \pi_{\mathcal{P}_p} u\|_{L^2(\hat{\kappa})}^2 = \|u - \pi_{\mathcal{S}_{p+1-d}} u\|_{L^2(\hat{\kappa})}^2 \leq C(s) \left( \frac{d}{p+1-d} \right)^{2s+2} |u|_{H^{s+1}(\hat{\kappa})}^2. \tag{53}$$

*and*

$$\|\nabla(u - \pi_{\mathcal{P}_p} u)\|_{L^2(\hat{\kappa})}^2 = \|\nabla(u - \pi_{\mathcal{S}_{p+1-d}} u)\|_{L^2(\hat{\kappa})}^2 \leq C(s) \left( \frac{d}{p-d} \right)^{2s} |u|_{H^{s+1}(\hat{\kappa})}^2. \tag{54}$$

*for any integer $s$, $d - 1 \leq s \leq \min\{p + 1 - d, l\}$, $p$ sufficiently large.*

*Remark 5* We emphasize that the above error bound for the $\pi_{\mathcal{P}_p}$ projection is $p$-suboptimal by one order for $d = 2$ and two orders for $d = 3$ for sufficiently smooth functions, but it is $p$-optimal for functions with finite Sobolev regularity in the case $l \leq p + 1 - d$. However, sub-optimality by one or two orders in $p$ is better than using the $\pi_{\lfloor p/d \rfloor}^Q$ projection, as suggested by [29] (see Corollary 4.52 on p. 190), which is sub-optimal in $p$ by at least $p/2$ orders for sufficiently smooth functions for $d = 2$. Moreover, the one- or two-order sub-optimality in $p$ for analytic functions does not influence the exponent of the exponential rate of convergence, as we shall see below.

## 5 Exponential convergence for analytic solutions

We shall be concerned with the proof of exponential convergence for serendipity FEMs and DGFEMs with the $\mathcal{P}_p$ basis over tensor product elements. For simplicity, we only consider the case when the given problem is piecewise analytic over the whole computational domain. Exponential convergence is then achieved by fixing the computational mesh, and increasing the polynomial order $p$. Only parallelepiped meshes are considered, which are the affine family obtained from the reference

element $\hat{\kappa} = (-1, 1)^d$. The analysis of FEMs and DGFEMs with a general *hp*-refinement strategy is beyond the scope of this analysis (see [25–28] for the analysis for both methods employing the $\mathcal{Q}_p$ basis).

The proof of exponential convergence for FEMs and DGFEMs depends on proving exponential convergence of $L^2$- and $H^1$-projections for piecewise analytic functions under *p*-refinement. The $H^1$-projection $\pi_{\mathcal{S}_p}$ onto $\mathcal{S}_p$ can be directly applied to *p*-FEMs for second-order elliptic problems with the same optimal rate as the $H^1$ projection $\pi_{\mathcal{Q}_p}$ (see [29] for details). For deriving error bounds of DGFEMs using the $L^2$- and $H^1$- projections onto $\mathcal{Q}_p$, we refer to [15, 21, 22]. Following similar techniques, we can prove the corresponding *hp*-bounds for DGFEMs employing the $\mathcal{P}_p$ basis, albeit with sub-optimal rate in *p*. The sub-optimality in *p* is due to the fact that the *p*-optimal bound for $L^2$-projection onto $\mathcal{P}_p$ basis over the trace of the tensor product elements is still open. Additionally, the $H^1$-projection onto the $\mathcal{P}_p$ basis is suboptimal in *p* by $d-1$ orders for sufficiently smooth functions. However, we point out that the sub-optimality in *p* by $d-1$ order, with $d = 2, 3$, does not influence the exponent of the exponential rate of convergence.

Next, we focus on deriving the exponential convergence for the $L^2$-projections in the $L^2$-norm and $H^1$-projections in the $L^2$-norm and $H^1$-seminorm on analytic problems under *p*-refinement on shape-regular *d*-parallelepiped meshes. The extension to anisotropic meshes will be consider in the future.

Let $\kappa$ be a parallelepiped element. For a function $u$ having an analytic extension into an open neighbourhood of $\bar{\kappa}$, we have:

$$\exists R_\kappa > 0, \quad C(u) > 0, \quad \forall s_\kappa : |u|_{H^{s_\kappa}(\kappa)} \le C(u)(R_\kappa)^{s_\kappa} \Gamma(s_\kappa + 1)|\kappa|^{1/2}, \quad (55)$$

where $|\kappa|$ denotes the measure of element $\kappa$, cf. [12, Theorem 1.9.3].

**Lemma 5** *Let* $u : \kappa \to \mathbb{R}$ *have an analytic extension to an open neighbourhood of* $\bar{\kappa}$. *Also let* $p_\kappa \ge 0$ *and* $0 \le s_\kappa \le p_\kappa + 1$ *be two positive numbers such that* $s_\kappa = \epsilon(p_\kappa + 1), 0 \le \epsilon \le 1$ *and* $d = 2, 3$. *Then, the following bounds hold:*

$$\|u - \Pi_{\mathcal{Q}_{p_\kappa}} u\|^2_{L^2(\kappa)} \le C(h_\kappa)^{2s_\kappa} \Phi_1(p_\kappa + 1, s_\kappa)|u|^2_{H^{s_\kappa}(\hat{\kappa})}$$
$$\le C(u)(p_\kappa + 1)e^{-2b_{1,\kappa}(p_\kappa + 1)}|\kappa|,$$

*and*

$$\|u - \Pi_{\mathcal{P}_{p_\kappa}} u\|^2_{L^2(\kappa)} \le C(h_\kappa)^{2s_\kappa} \Phi_d(p_\kappa + 1, s_\kappa)|u|^2_{H^{s_\kappa}(\hat{\kappa})}$$
$$\le C(u)(p_\kappa + 1)e^{-2b_{2,\kappa}(p_\kappa + 1)}|\kappa|.$$

*Here, C and C(u) are positive constants depending elemental shape regularity, and C(u) also depends on u.* $F_1(R_\kappa, \epsilon) = \frac{(1-\epsilon)^{1-\epsilon}}{(1+\epsilon)^{1+\epsilon}}(\epsilon R_\kappa)^{2\epsilon}$, $\epsilon_{\min} = 1/\sqrt{1 + R_\kappa^2}$, $b_{1,\kappa} := \frac{1}{2}|\log F_1(R_\kappa, \epsilon_{\min})| + \epsilon_{\min}|\log h_\kappa|$ *and* $b_{2,\kappa} := b_{1,\kappa} - \epsilon_{\min} \log d$.

*Proof* Using standard scaling arguments for $\kappa$ together with Lemma 1 and Theorem 1, we have the approximation results for the $L^2$-projection over $\kappa$. For brevity, we

set $q_\kappa = p_\kappa + 1$. By employing the relation (11) and the fact $|u|_{V^l(\kappa)} \le |u|_{H^l(\kappa)}$, we have the bounds:

$$
\begin{aligned}
\Phi_1(p_\kappa+1, s_\kappa)|u|^2_{H^{s_\kappa}(\hat{\kappa})} &\le C(u)(R_\kappa)^{2s_\kappa} \Gamma(s_\kappa+1)^2 \frac{\Gamma(q_\kappa - s_\kappa + 1)}{\Gamma(q_\kappa + s_\kappa + 1)}|\kappa| \\
&\le C(u)(R_\kappa)^{2\epsilon q_\kappa} \frac{(\epsilon q_\kappa)^{2\epsilon q_\kappa + 1}}{e^{2\epsilon q_\kappa}} \frac{((1-\epsilon)q_\kappa)^{(1-\epsilon)q_\kappa} e^{-(1-\epsilon)q_\kappa}}{((1+\epsilon)q_\kappa)^{(1+\epsilon)q_\kappa} e^{-(1+\epsilon)q_\kappa}}|\kappa| \\
&\le C(u)q_\kappa (F_1(R_\kappa, \epsilon))^{q_\kappa}|\kappa|,
\end{aligned}
$$

where

$$
F_1(R_\kappa, \epsilon) = \frac{(1-\epsilon)^{1-\epsilon}}{(1+\epsilon)^{1+\epsilon}}(\epsilon R_\kappa)^{2\epsilon}.
$$

Recalling (55), we have $R_\kappa > 0$,

$$
\min_{0<\epsilon<1} F_1(R_\kappa, \epsilon) = F_1(R_\kappa, \epsilon_{\min}) = \left(\frac{R_\kappa}{\sqrt{1+R_\kappa^2}+1}\right)^2 < 1, \quad \epsilon_{\min} = \frac{1}{\sqrt{1+R_\kappa^2}}. \tag{56}
$$

Thus, we have

$$
\frac{\Gamma(p_\kappa - s_\kappa + 2)}{\Gamma(p_\kappa + s_\kappa + 2)}|u|^2_{H^{s_\kappa}(\hat{\kappa})} \le C(u)q_\kappa e^{-|\log F_1(R_\kappa, \epsilon_{\min})|q_\kappa}|\kappa|. \tag{57}
$$

Therefore, we have the exponential convergence for the $L^2$-projection $\Pi_{\mathcal{Q}_{p_\kappa}}$, via

$$
\|u - \Pi_{\mathcal{Q}_{p_\kappa}}u\|^2_{L^2(\kappa)} \le C(u)(p_\kappa+1)e^{-2b_{1,\kappa}(p_\kappa+1)}|\kappa|, \tag{58}
$$

with $b_{1,\kappa} := \frac{1}{2}|\log F_1(R_\kappa, \epsilon_{\min})| + \epsilon_{\min}|\log h_\kappa|$. Similarly, for the $L^2$-projection $\Pi_{\mathcal{P}_{p_\kappa}}$, Stirling's formula implies

$$
\begin{aligned}
\Phi_d(p_\kappa+1, s_\kappa)|u|^2_{H^{s_\kappa}(\hat{\kappa})} &\le C(u)(R_\kappa)^{2s_\kappa} \Gamma(s_\kappa+1)^2 \left(\frac{\Gamma(\frac{q_\kappa - s_\kappa}{d}+1)}{\Gamma(\frac{q_\kappa + s_\kappa}{d}+1)}\right)^d |\kappa| \\
&\le C(u)(R_\kappa)^{2\epsilon q_\kappa} \frac{(\epsilon q_\kappa)^{2\epsilon q_\kappa + 1}}{e^{2\epsilon q_\kappa}} \\
&\quad \times \frac{((1-\epsilon)q_\kappa)^{(1-\epsilon)q_\kappa}(ed)^{-(1-\epsilon)q_\kappa}}{((1+\epsilon)q_\kappa)^{(1+\epsilon)q_\kappa}(ed)^{-(1+\epsilon)q_\kappa}}|\kappa| \\
&\le C(u)q_\kappa (F_2(R_\kappa, \epsilon))^{q_\kappa}|\kappa|,
\end{aligned}
$$

where,

$$
F_2(R_\kappa, \epsilon) = \frac{(1-\epsilon)^{1-\epsilon}}{(1+\epsilon)^{1+\epsilon}}(\epsilon R_\kappa d)^{2\epsilon},
$$

with the minimum,

$$
\min_{0<\epsilon<1} F_2(R_\kappa, \epsilon) = \left(\frac{R_\kappa d}{\sqrt{1+(R_\kappa d)^2}+1}\right)^2 < 1.
$$

In order to compare with the slope of projection $\Pi_{\mathcal{Q}_{p_\kappa}}$, here, we will use the same $\epsilon_{\min}$. We have

$$
\min_{0<\epsilon<1} F_2(R_\kappa, \epsilon) \le F_2(R_\kappa, \epsilon_{\min}) = F_1(R_\kappa, \epsilon_{\min})d^{2\epsilon_{\min}}.
$$

Thus, we have

$$\|u - \Pi_{\mathcal{P}_{p_\kappa}} u\|^2_{L^2(\kappa)} \leq C(u)(p+1)e^{-2b_{2,\kappa}(p_\kappa+1)}|\kappa|, \tag{59}$$

with slope $b_{2,\kappa} := \frac{1}{2}|\log F_1(R_\kappa, \epsilon_{\min})| + \epsilon_{\min}(|\log h_\kappa| - \log d)$. □

Next, we begin to derive the exponential convergence for $H^1$-projections.

**Lemma 6** *Let $u : \kappa \to \mathbb{R}$ have an analytic extension to an open neighbourhood of $\bar{\kappa}$. Also let $p_\kappa \geq 2d$ and $(d-1) \leq s_\kappa \leq p_\kappa$ be two positive numbers such that $s_\kappa = \epsilon p_\kappa$, $0 < \epsilon \leq 1$ and $d = 2, 3$. Then, the following bounds hold:*

$$\|u - \pi_{\mathcal{Q}_{p_\kappa}} u\|^2_{L^2(\kappa)} \leq C(h_\kappa)^{2s_\kappa+2}\Phi_1(p_\kappa+1, s_\kappa+1)|u|^2_{H^{s_\kappa+1}(\hat{\kappa})} \tag{60}$$
$$\leq C(u)p_\kappa e^{-2b_{1,\kappa} p_\kappa}|\kappa|,$$

$$\|u - \pi_{\mathcal{S}_{p_\kappa}} u\|^2_{L^2(\kappa)} \leq C(h_\kappa)^{2s_\kappa+2}\Phi_d(p_\kappa+1, s_\kappa+1)|u|^2_{H^{s_\kappa+1}(\hat{\kappa})} \tag{61}$$
$$\leq C(u)p_\kappa e^{-2b_{2,\kappa} p_\kappa}|\kappa|,$$

*and*

$$\|\nabla(u - \pi_{\mathcal{Q}_{p_\kappa}} u)\|^2_{L^2(\kappa)} \leq C(h_\kappa)^{2s_\kappa}\Phi_1(p_\kappa, s_\kappa)|u|^2_{H^{s_\kappa+1}(\hat{\kappa})} \leq C(u)p_\kappa^3 e^{-2b_{1,\kappa} p_\kappa}|\kappa|,$$

$$\|\nabla(u - \pi_{\mathcal{S}_{p_\kappa}} u)\|^2_{L^2(\kappa)} \leq C(h_\kappa)^{2s_\kappa}\Phi_d(p_\kappa, s_\kappa)|u|^2_{H^{s_\kappa+1}(\hat{\kappa})} \leq C(u)p_\kappa^3 e^{-2b_{2,\kappa} p_\kappa}|\kappa|.$$

*Here, $C$ and $C(u)$ are positive constants depending elemental shape regularity, and $C(u)$ also depends on $u$. $F_1(R_\kappa, \epsilon) = \frac{(1-\epsilon)^{1-\epsilon}}{(1+\epsilon)^{1+\epsilon}}(\epsilon R_k)^{2\epsilon}$, $\epsilon_{\min} = 1/\sqrt{1 + R_\kappa^2}$, $b_{1,\kappa} := \frac{1}{2}|\log F_1(R_\kappa, \epsilon_{\min})| + \epsilon_{\min}|\log h_\kappa|$ and $b_{2,\kappa} := b_\kappa^1 - \epsilon_{\min} \log d$.*

*Proof* The proof follows by the same techniques used in Lemma 5. □

In the above Lemmas 5 and 6, we can see that the $L^2$-norm error for both $L^2$-projections $\Pi_{\mathcal{Q}_{p_\kappa}}$ and $\Pi_{\mathcal{P}_{p_\kappa}}$, and the $L^2$-norm and $H^1$-seminorm errors for the $H^1$-projections $\pi_{\mathcal{S}_{p_\kappa}}$ and $\pi_{\mathcal{Q}_{p_\kappa}}$ decay exponentially for analytic functions under $p$-refinement. If we measure the error against $p$, the exponent $b_{1,\kappa}$ for the $\mathcal{Q}_p$ basis is slightly greater than the exponent $b_{2,\kappa}$ for the $\mathcal{P}_p$ basis and $\mathcal{S}_p$ basis by a small factor of $(\log d)/\sqrt{1 + R_\kappa^2}$. By using Lemmas 5 and 6, we can also derive the following theorem.

**Theorem 4** *Let $u$ be an analytic function as defined in (55), and exponent $b_{1,\kappa}$ and $b_{2,\kappa}$ defined in Lemma 5. Then, there exists $C > 0$ such that following bounds hold:*

$$\|u - \Pi_{\mathcal{Q}_{p_\kappa}} u\|^2_{L^2(\kappa)} \leq Ce^{-2b_{1,\kappa}\sqrt[d]{Dof}}, \tag{62}$$

$$\|u - \Pi_{\mathcal{P}_{p_\kappa}} u\|^2_{L^2(\kappa)} \leq Ce^{-2(b_{2,\kappa}\sqrt[d]{d!})\sqrt[d]{Dof}}, \tag{63}$$

*and*

$$\|u - \pi_{\mathcal{Q}_{p_\kappa}} u\|_{L^2(\kappa)}^2 \le C e^{-2b_{1,\kappa} \sqrt[d]{Dof}}, \tag{64}$$

$$\|u - \pi_{\mathcal{S}_{p_\kappa}} u\|_{L^2(\kappa)}^2 \le C e^{-2(b_{2,\kappa} \sqrt[d]{d!}) \sqrt[d]{Dof}}, \tag{65}$$

*and*

$$\|\nabla(u - \pi_{\mathcal{Q}_{p_\kappa}} u)\|_{L^2(\kappa)}^2 \le C e^{-2b_{1,\kappa} \sqrt[d]{Dof}}, \tag{66}$$

$$\|\nabla(u - \pi_{\mathcal{S}_{p_\kappa}} u)\|_{L^2(\kappa)}^2 \le C e^{-2(b_{2,\kappa} \sqrt[d]{d!}) \sqrt[d]{Dof}}. \tag{67}$$

*Proof* By recalling the relationship between degrees of freedom and polynomial order $p$ for both the $\mathcal{Q}_p$ and $\mathcal{P}_p$ bases, we have

$$Dof(\mathcal{Q}_p) = (p+1)^d, \tag{68}$$

and

$$Dof(\mathcal{P}_p) = \binom{p+d}{d} = \frac{(p+1)^d}{d!} + \mathcal{O}((p+1)^{d-1}). \tag{69}$$

Then, (62) and (63) follow from Lemma 5.

By using relations (20) and (40), we have the asymptotic relation

$$Dof(\mathcal{S}_p) \approx \frac{p^d}{d!} + \mathcal{O}(p^{d-1}). \tag{70}$$

The relations (64), (65), (66) and (67) follow from the Lemma 6.                    □

For $d = 2, 3$, if the following condition

$$\frac{1}{2} |\log F_1(R_\kappa, \epsilon_{\min})| + \epsilon_{\min} |\log h_\kappa| \gg \epsilon_{\min} \log d, \tag{71}$$

holds, then we have $b_{2,\kappa} \approx b_{1,\kappa}$. It is easy to see that for small $R_\kappa$ or small mesh size $h$, the condition (71) will be satisfied. Moreover, we point out that an analytic function having sufficiently small $R_\kappa$ is equivalent to the function having an analytic continuation into a sufficiently large open neighbourhood of $\bar{\kappa}$, see [12] for details.

Now, if we consider the error in terms of $\sqrt[d]{Dof}$ for the above bounds, the exponent for the exponential convergence rate of the $\mathcal{P}_p$ basis and the $\mathcal{S}_p$ basis are larger than the exponent for the $\mathcal{Q}_p$ basis by a fixed factor of $\sqrt[d]{d!}$.

We have observed a steeper slope in error against $\sqrt[d]{Dof}$ for FEMs with $\mathcal{S}_p$ basis and DGFEMs with $\mathcal{P}_p$ basis. For $d = 2$, this suggests a typical ratio between convergence slopes of DGFEMs with $\mathcal{P}_p$ and $\mathcal{Q}_p$ bases, FEMs with $\mathcal{S}_p$ and $\mathcal{Q}_p$ bases, to be $\sqrt{2!} \approx 1.414$. For $d = 3$, this ratio is $\sqrt[3]{3!} \approx 1.817$. The numerical examples in Section 6 show that the ratio is slightly worse than the ideal ratio, but it is not far from the ideal ratio.

## 6 Numerical examples

We present some numerical examples to confirm the theoretical analysis in the previous sections. All the numerical examples are computed by Matlab on the High Performance Computing facility ALICE of the University of Leicester. For simplicity of presentation, we use DGFEM(P) and DGFEM(Q) to denote the DGFEMs with local polynomial basis consisting of either $\mathcal{P}_p$ or $\mathcal{Q}_p$ polynomials and use FEM(S) and FEM(Q) to denote the FEMs with local polynomial basis consisting of either $\mathcal{S}_p$ or $\mathcal{Q}_p$ polynomials.

The comparisons are mainly made between the slope of FEM(S) and FEM(Q) over square meshes for $d = 2$ and hexahedral meshes for $d = 3$ under *p*-refinement. The slopes of the convergence lines are calculated by taking the average of the last two slopes of the line segments of each convergence line. We will also present an example comparing DGFEM(P) and DGFEM(Q). For more numerical examples for DGFEMs, see [13, 14].

### 6.1 Example 1

In the first example, we investigate the computational efficiency of DGFEM(P) and DGFEM(Q) schemes. To this end, we consider a partial differential equation with non-negative characteristic form of mixed type. Let $\Omega = (-1, 1)^2$, and consider the PDE problem:

$$\begin{cases} -x^2 u_{yy} + u_x + u = 0, & \text{for} -1 \le x \le 1, y > 0, \\ u_x + u = 0, & \text{for} -1 \le x \le 1, y \le 0, \end{cases} \quad (72)$$
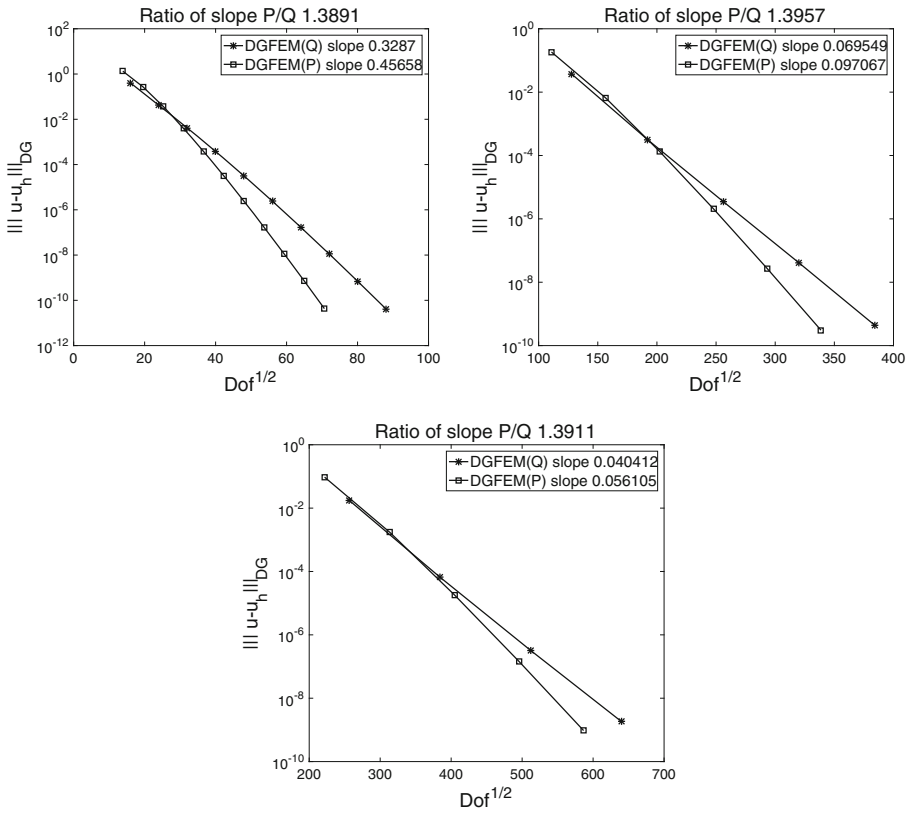
with exact solution:

$$u(x, y) = \begin{cases} \sin(\frac{1}{2}\pi(1 + y)) \exp(-(x + \frac{\pi^2 x^3}{12})), & \text{for} -1 \le x \le 1, y > 0, \\ \sin(\frac{1}{2}\pi(1 + y)) \exp(-x), & \text{for} -1 \le x \le 1, y \le 0. \end{cases} \quad (73)$$

This problem is hyperbolic in the region $y \le 0$ and parabolic for $y > 0$. In order to ensure continuity of the normal flux across $y = 0$, where the partial differential equation changes type, the exact solution has a discontinuity across the line $y = 0$, cf. [9, 15].

By following [9], we use the symmetric interior penalty DGFEMs employing a special class of quadrilateral meshes for which the discontinuity in the exact solution lies on element interfaces. In this setting, we modify the discontinuity-penalization parameter $\sigma$, so that $\sigma$ vanishes on edges which form part of the interface $y = 0$; this ensures that the (physical) discontinuity present in the exact solution is not penalized within by the numerical scheme.

In this case, the exact solution is piecewise analytic on the two parts of the domain. In Fig. 2, we observe that the DG–norm $|||u - u_h|||_{\text{DG}}$ decays exponentially for both DGFEM(P) and DGFEM(Q) under *p*-refinement on 64, 4096 and 16384 uniform square elements. The definition of DG–norm $|||\cdot|||_{\text{DG}}$ can be found in [9]. Moreover, the slope of the convergence line for the DGFEM(P) is greater than the line of

**Fig. 2** Example 1: Convergence of the DGFEMs under $p$-refinement on uniform square elements ($\||u - u_h\||_{DG}$). $8 \times 8$ mesh (left); $64 \times 64$ mesh (right); $128 \times 128$ mesh (bottom)
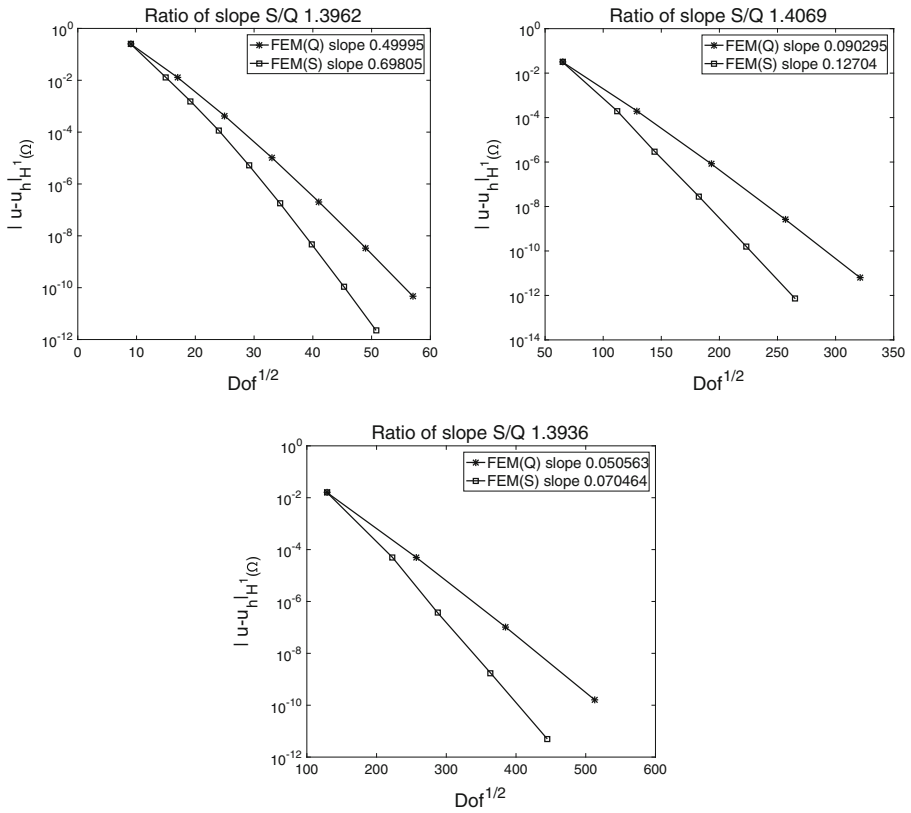
DGFEM(Q) in error against $\sqrt{Dof}$. The ratio between the two slopes is about 1.39 on coarse meshes and fine meshes. The numerical observation confirms the theoretical results in Theorem 4.

## 6.2 Example 2

In the second example, we investigate the computational efficiency of FEM(S) and FEM(Q) on standard tensor-product elements (quadrilaterals in 2D and hexahedra in 3D).

Firstly, we consider the following two-dimensional Poisson problem: let $\Omega = (0, 1)^2$ and select $f = 2\pi^2 \sin(\pi x) \sin(\pi y)$, so that the exact solution is given by $u = \sin(\pi x) \sin(\pi y)$.

In this case, the exact solution is piecewise analytic on the domain. In Fig. 3, we observe that the $H^1$-seminorm $|u - u_h|_{H^1(\Omega)}$ decays exponentially for both FEM(S)
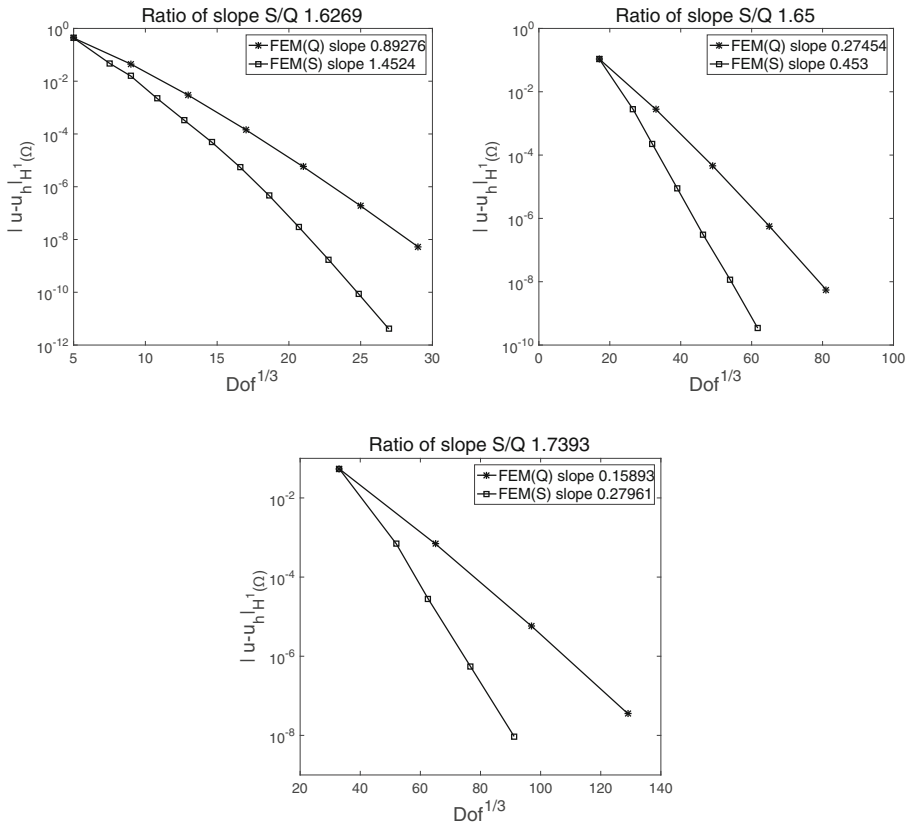
**Fig. 3** Example 2: Convergence of the FEMs under *p*-refinement on uniform square elements ($|u - u_h|_{H^1(\Omega)}$). $8 \times 8$ mesh (left); $64 \times 64$ mesh (right); $128 \times 128$ mesh (bottom)

and FEM(Q) under *p*-refinement on 64, 4096 and 16384 uniform square elements. Again, we observe that the slope of the convergence line for the FEM(S) is greater than the line of FEM(Q) in error against $\sqrt{Dof}$. The ratio between the two slopes is about 1.39 on coarse meshes and fine meshes.

We now consider the three-dimensional variant of the above problem. Let $\Omega = (0, 1)^3$ and select $f = 3\pi^2 \sin(\pi x) \sin(\pi y) \sin(\pi z)$, so that the exact solution is given by $u = \sin(\pi x) \sin(\pi y) \sin(\pi z)$.

In Fig. 4, we observe that the $H^1$-seminorm $|u - u_h|_{H^1(\Omega)}$ decays exponentially for both FEM(S) and FEM(Q) under *p*-refinement on 64, 4096 and 32768 uniform hexahedral elements. Moreover, we observe that the slope of the convergence line for the FEM(S) is greater than the line of FEM(Q) in error against $\sqrt[3]{Dof}$. The ratio between the two slopes is about 1.62 on coarse meshes and 1.73 on fine meshes. The numerical observation confirms the theoretical results in Theorem 4.

**Fig. 4** Example 2: Convergence of the FEMs under *p*-refinement on uniform hexahedral elements ($|u - u_h|_{H^1(\Omega)}$). $4 \times 4 \times 4$ mesh (left); $16 \times 16 \times 16$ mesh (right); $32 \times 32 \times 32$ mesh (bottom)

# References

1. Abramowitz, M., Stegun, I.A.: Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables, Volume 55 of National Bureau of Standards Applied Mathematics Series. For Sale by the Superintendent of Documents, U.S. Government Printing Office, Washington, D.C. (1964)
2. Ainsworth, M., Pinchedez, K.: *hp*-approximation theory for BDFM and RT finite elements on quadrilaterals. SIAM J. Numer. Anal. **40**(6), 2047–2068 (2003). 2002

3. Arnold, D.N., Awanou, G.: The serendipity family of finite elements. Found. Comput. Math. **11**(3), 337–344 (2011)
4. Babuška, I., Guo, B.Q.: The *h-p* version of the finite element method for domains with curved boundaries. SIAM J. Numer. Anal. **25**(4), 837–861 (1988)
5. Babuška, I., Guo, B.Q.: Direct and inverse approximation theorems for the p-version of the finite element method in the framework of weighted besov spaces. part I: Approximability of functions in the weighted besov spaces. SIAM J. Numer. Anal. **39**(5), 1512–1538 (2002)
6. Batir, N.: Inequalities for the gamma function. Arch. Math. **91**(6), 554–563 (2008)
7. Bergh, J., Löfström, J.: Interpolation Spaces. An Introduction. Springer-Verlag, Berlin-New York Grundlehren der Mathematischen Wissenschaften, No. 223 (1976)
8. Cangiani, A., Dong, Z., Georgoulis, E.H.: *hp*-version space-time discontinuous Galerkin methods for parabolic problems on prismatic meshes. SIAM J. Sci. Comput. **39**(4), A1251–A1279 (2017)
9. Cangiani, A., Dong, Z., Georgoulis, E.H., Houston, P.: *hp*-version discontinuous Galerkin methods for advection-diffusion-reaction problems on polytopic meshes. M2AN Math. Model. Numer. Anal. **50**(3), 699–725 (2016)
10. Cangiani, A., Georgoulis, E.H., Houston, P.: *hp*-version discontinuous Galerkin methods on polygonal and polyhedral meshes. Math. Models Methods Appl. Sci. **24**(10), 2009–2041 (2014)
11. Canuto, C., Quarteroni, A.: Approximation results for orthogonal polynomials in Sobolev spaces. Math. Comput. **38**(157), 67–86 (1982)
12. Davis, P.J.: Interpolation and Approximation. Courier Corporation (1975)
13. Dong, Z.: Discontinuous Galerkin Methods on Polytopic Meshes. PhD thesis, University of Leicester (2016)
14. Dong, Z.: On the exponent of exponential convergence of *hp*-finite element spaces. arXiv:1704.08046 (2017)
15. Georgoulis, E.H.: Discontinuous Galerkin Methods on Shape-regular and Anisotropic Meshes. D.Phil. Thesis, University of Oxford (2003)
16. Gui, W., Babuška, I.: The *h*, *p* and *h-p* versions of the finite element method in 1 dimension. I–III. Numer. Math. **49**(6), 577–683 (1986)
17. Guo, B.Q.: The *h-p* version of the finite element method for elliptic equations of order 2*m*. Numer. Math. **53**(1-2), 199–224 (1988)
18. Guo, B.Q.: The *h-p* version of the finite element method for solving boundary value problems in polyhedral domains. In: Boundary Value Problems and Integral Equations in Nonsmooth Domains (Luminy, 1993), vol. 167 of Lecture Notes in Pure and Appl. Math., pp. 101–120. Dekker, New York (1995)
19. Guo, B.Q., Babuška, I.: The hp version of the finite element method. Part I: the basic approximation results. Comput. Mech. **1**(1), 21–41 (1986)
20. Guo, B.Q., Babuška, I.: The hp version of the finite element method. Part II: general results and applications. Comput. Mech. **1**(1), 203–220 (1986)
21. Houston, P., Schwab, C., Süli, E.: Stabilized *hp*-finite element methods for first-order hyperbolic problems. SIAM J. Numer. Anal. **37**(5), 1618–1643 (2000). (electronic)
22. Houston, P., Schwab, C., Süli, E.: Discontinuous *hp*-finite element methods for advection-diffusion-reaction problems. SIAM J. Numer. Anal. **39**(6), 2133–2163 (2002). (electronic)
23. Kretzschmar, F., Moiola, A., Perugia, I., Schnepp, S.M.: A priori, error analysis of space-time Trefftz discontinuous Galerkin methods for wave problems. IMA J. Numer Anal. **36**(4), 1599–1635 (2016)
24. Melenk, J.M., Schwab, C.: An *hp*–finite element method for convection-diffusion problems in one dimension. IMA J. Numer. Anal. **19**(3), 425–453 (1999)
25. Schötzau, D., Schwab, C.: Exponential convergence for *hp*-version and spectral finite element methods for elliptic problems in polyhedra. Math. Models Methods Appl. Sci. **25**(9), 1617–1661 (2015)
26. Schötzau, D., Schwab, C., Wihler, T.P.: *hp*-dGFEM for second-order elliptic problems in polyhedra I Stability on geometric meshes. SIAM J. Numer. Anal. **51**(3), 1610–1633 (2013)
27. Schötzau, D., Schwab, C., Wihler, T.P.: *hp*-dGFEM for second order elliptic problems in polyhedra II Exponential convergence. SIAM J. Numer. Anal. **51**(4), 2005–2035 (2013)
28. Schötzau, D., Schwab, C., Wihler, T.P.: *hp*-dGFEM, for second-order mixed elliptic problems in polyhedra. Math. Comput. **85**(299), 1051–1083 (2016)
29. Schwab, C.: *p*– and *hp*–Finite element methods: Theory and applications in solid and fluid mechanics Oxford University Press: Numerical mathematics and scientific computation (1998)
30. Szabó, B., Babuška, I.: Finite Element Analysis. A Wiley-Interscience Publication. Wiley, New York (1991)
31. Wihler, T.P., Frauenfelder, P., Schwab, C.: Exponential convergence of the *hp*-DGFEM, for diffusion problems. Comput. Math. Appl. **46**, 183–205 (2003)