

Orthogonal polynomial bases in the mixed virtual element method

Original

Orthogonal polynomial bases in the mixed virtual element method / Berrone, S.; Scialo, S.; Teora, G.. - In: NUMERICAL METHODS FOR PARTIAL DIFFERENTIAL EQUATIONS. - ISSN 0749-159X. - (2024), pp. 1-26. [10.1002/num.23144]

Availability:

This version is available at: 11583/2992315 since: 2024-09-09T08:10:15Z

Publisher:

Wiley

Published

DOI:10.1002/num.23144

Terms of use:

This article is made available under terms and conditions as specified in the corresponding bibliographic description in the repository

Publisher copyright

(Article begins on next page)

RESEARCH ARTICLE

WILEY

Orthogonal polynomial bases in the mixed virtual element method

Stefano Berrone | Stefano Scialò | Gioana Teora 

Dipartimento di Scienze Matematiche “G. L. Lagrange”, Politecnico di Torino, Turin, Italy

Correspondence

Gioana Teora, Dipartimento di Scienze Matematiche “G. L. Lagrange”, Politecnico di Torino, Turin, Italy.

Email: gioana.teora@polito.it

Funding information

Ministero dell’Istruzione, dell’Università e della Ricerca; Ministero dell’Università e della Ricerca; Gruppo Nazionale per il Calcolo Scientifico.

Abstract

The use of orthonormal polynomial bases has been found to be efficient in preventing ill-conditioning of the system matrix in the primal formulation of Virtual Element Methods (VEM) for high values of polynomial degree and in presence of badly-shaped polygons. However, we show that using the natural extension of a orthogonal polynomial basis built for the primal formulation is not sufficient to cure ill-conditioning in the mixed case. Thus, in the present work, we introduce an orthogonal vector-polynomial basis which is built ad hoc for being used in the mixed formulation of VEM and which leads to very high-quality solution in each tested case. Furthermore, a numerical experiment related to simulations in Discrete Fracture Networks (DFN), which are often characterised by very badly-shaped elements, is proposed to validate our procedures.

KEYWORDS

ill-conditioning, mixed VEM, orthogonal polynomial basis

1 | INTRODUCTION

The Mixed Virtual Element Methods were introduced originally in [11] for the Poisson problem in the two-dimensional case and then were extended to more general elliptic equations in [6]. In the Mixed Virtual Element Space, two discrete spaces are introduced for approximating the pressure variable and the velocity field, respectively. The first space is a scalar-polynomial space, while a vector-polynomial basis is required to build the local projection matrices and for defining the internal degrees of freedom needed to obtain an approximation of the velocity field. It was observed that using the classical choice

This is an open access article under the terms of the [Creative Commons Attribution-NonCommercial-NoDerivs](https://creativecommons.org/licenses/by-nc-nd/4.0/) License, which permits use and distribution in any medium, provided the original work is properly cited, the use is non-commercial and no modifications or adaptations are made. © 2024 The Author(s). *Numerical Methods for Partial Differential Equations* published by Wiley Periodicals LLC.

of scaled monomials in the definition of internal degrees of freedom in the primal Virtual Element construction [3, 4, 7, 11], the system matrix could become ill-conditioned in presence of *badly-shaped* polygons [9] (e.g., collapsing edges and bulks) and for high values of local polynomial degree [15].

The present work aims at defining new polynomial bases for the mixed VEM construction yielding well-conditioned local projection matrices also in presence of badly-shaped elements. More precisely, we propose two different approaches for building a vector-polynomial basis, which we briefly called “Partial” and “Ortho”, respectively. The first one is the natural extension to the mixed case of the approach presented in [15] for the primal VEM, which allows us to build a vector-polynomial basis that is only partially L^2 -orthonormalized. We show that the use of such basis is not sufficient to cure the ill-conditioning of the system matrix related to the mixed formulation of VEM in all circumstances [1], even if, in the primal setting, it reveals to be the best alternative. Thus, we introduce the Ortho approach which aims to orthogonalize the gradients of a proper scalar-polynomial basis in order to obtain a full orthonormal vector-polynomial basis. We show that this approach leads to the best local and global performances, throughout different numerical experiments characterised by challenging geometries.

The outline of the present paper is the following. We define the model problem in Section 2 and its mixed VEM approximation in Section 3. In Section 4, we describe how to build the new polynomial bases, while in Section 5 we show an efficient implementation of the method, totally matrix-based. Finally, in Section 6, we perform some numerical experiments that show the advantages of using the new polynomial bases.

Let us introduce some notations used throughout the paper. Given $k \in \mathbb{N}$, we use $(\cdot, \cdot)_{k,\sigma}$ and $\|\cdot\|_{k,\sigma}$ to indicate the inner product and the norm in the Sobolev space $H^k(\sigma)$ on some open subset $\sigma \subset \mathbb{R}^2$, respectively. Furthermore, if $\mathbf{v} = [v_1, v_2]^T$ and $\mathbf{u} = [u_1, u_2]^T$ are vectors in $L^2(\sigma) \times L^2(\sigma)$, we define

$$(\mathbf{v}, \mathbf{u})_{0,\sigma} = \int_{\sigma} (v_1 u_1 + v_2 u_2), \quad \|\mathbf{v}\|_{0,\sigma} = \sqrt{(\mathbf{v}, \mathbf{v})_{0,\sigma}}. \quad (1)$$

Let $\Omega \subset \mathbb{R}^2$ be a bounded convex polygonal domain with boundary Γ and let \mathbf{n}_{Γ} be the outward unit normal vector to Γ , then we define the functional spaces

$$H(\operatorname{div}; \Omega) = \{\mathbf{v} \in L^2(\Omega) \times L^2(\Omega) : \nabla \cdot \mathbf{v} \in L^2(\Omega)\}, \quad (2)$$

$$H_{0,\Gamma_N}(\operatorname{div}; \Omega) = \{\mathbf{v} \in H(\operatorname{div}, \Omega) : \mathbf{v} \cdot \mathbf{n}_{\Gamma} = 0 \text{ on } \Gamma_N \subseteq \Gamma\}, \quad (3)$$

$$H(\operatorname{rot}; \Omega) = \{\mathbf{v} \in L^2(\Omega) \times L^2(\Omega) : \operatorname{rot} \mathbf{v} \in L^2(\Omega)\}. \quad (4)$$

Furthermore, let $H^{-\frac{1}{2}}(\Gamma)$ be the dual space of the Sobolev space $H^{\frac{1}{2}}(\Gamma)$, the symbol $\langle \cdot, \cdot \rangle_{\pm \frac{1}{2}, \Gamma}$ denotes the duality pairing between $H^{-\frac{1}{2}}(\Gamma)$ and $H^{\frac{1}{2}}(\Gamma)$.

2 | THE CONTINUOUS PROBLEM AND THE MIXED VARIATIONAL FORMULATION

Let $\boldsymbol{\kappa}$ be a symmetric uniformly positive definite tensor over Ω , γ a sufficiently smooth function $\Omega \rightarrow \mathbb{R}$ and \mathbf{b} a smooth vector valued function $\Omega \rightarrow \mathbb{R}^2$. We consider the following problem:

$$\begin{cases} \nabla \cdot (-\boldsymbol{\kappa} \nabla p + \mathbf{b}p) + \gamma p = f & \text{in } \Omega, \\ p = g_D & \text{on } \Gamma_D, \\ (-\boldsymbol{\kappa} \nabla p + \mathbf{b}p) \cdot \mathbf{n}_{\Gamma_N} = g_N & \text{on } \Gamma_N, \end{cases} \quad (5)$$

where Γ_D and Γ_N are the Dirichlet and the Neumann boundary, respectively, such that $\Gamma_D \cup \Gamma_N = \Gamma$ and $|\Gamma_D \cap \Gamma_N| = 0$.

In order to introduce the mixed variational formulation, we define

$$K = \kappa^{-1}, \quad \beta = Kb, \tag{6}$$

and we re-write problem (5) as

$$\begin{cases} Ku = -\nabla p + \beta p & \text{in } \Omega, \\ \nabla \cdot u + \gamma p = f & \text{in } \Omega, \\ p = g_D & \text{on } \Gamma_D, \\ u \cdot n_{\Gamma_N} = g_N & \text{on } \Gamma_N. \end{cases} \tag{7}$$

Thus, the mixed variational formulation of (5) reads:

Find $u = u_0 + u_N$, with $u_0 \in V = H_{0,\Gamma_N}(\text{div}; \Omega)$, and $p \in Q = L^2(\Omega)$ such that

$$\begin{cases} (Ku, v)_{0,\Omega} - (p, \nabla \cdot v)_{0,\Omega} - (\beta p, v)_{0,\Omega} = -\langle g_D, v \cdot n_{\Gamma_D} \rangle_{\frac{1}{2}, \Gamma_D} & \forall v \in V \\ (\nabla \cdot u, q)_{0,\Omega} + (\gamma p, q)_{0,\Omega} = (f, q)_{0,\Omega} & \forall q \in Q \end{cases} \tag{8}$$

where $u_N \in H(\text{div}; \Omega)$ is a chosen function that satisfies $u_N \cdot n_{\Gamma_N} = g_N$ on Γ_N .

3 | THE MIXED VIRTUAL ELEMENT METHOD

Let \mathcal{T}_h be a decomposition of Ω into star-shaped polygons E . We will denote by x_E , h_E and $\mathcal{E}_{h,E}$ the centroid, the diameter and the set of edges of E , respectively. We further set $N^{E,e} = \#\mathcal{E}_{h,E}$, and, as usual, we fix $h = \max_{E \in \mathcal{T}_h} h_E$.

Moreover, $\mathbb{P}_k(E)$ is the set of all polynomials defined on E of degree less or equal to $k \geq 0$ and $n_k = \dim \mathbb{P}_k(E) = \frac{(k+1)(k+2)}{2}$. For the ease of the notation, we fix $\mathbb{P}_{-1} = \{0\}$ and $n_{-1} = 0$ and we introduce the natural function $\ell : \mathbb{N}^2 \rightarrow \mathbb{N}$ which associates

$$(0, 0) \mapsto 1, \quad (1, 0) \rightarrow 2, \quad (0, 1) \mapsto 3, \quad (2, 0) \mapsto 4, \quad (1, 1) \rightarrow 5, \dots \tag{9}$$

A classical choice of the basis for $\mathbb{P}_k(E)$ that can be found in virtual element literature (see [3, 4, 7]) is the set of the scaled monomials, which can be defined as

$$\mathcal{M}_k(E) = \left\{ m_\alpha^k = \left(\frac{x - x_E}{h_E} \right)^\alpha, \forall \alpha = \ell(\alpha) \in \mathbb{N}^2 \text{ s.t. } \alpha = 1, \dots, n_k \right\}, \tag{10}$$

where the function ℓ is defined in (9).

As in [6], we introduce the (vector) polynomial space

$$\mathcal{G}_k^\nabla(E) = \nabla \mathbb{P}_{k+1}(E) = \{ \mathbf{g}_\alpha^{\nabla,k} \}_{\alpha=1}^{n_k^\nabla} \tag{11}$$

and its complement $\mathcal{G}_k^\perp(E) = \{ \mathbf{g}_\alpha^{\perp,k} \}_{\alpha=1}^{n_k^\perp}$ in $[\mathbb{P}_k(E)]^2$, which satisfies

$$[\mathbb{P}_k(E)]^2 = \mathcal{G}_k^\nabla(E) \oplus \mathcal{G}_k^\perp(E), \tag{12}$$

where \oplus is the direct sum operator, and

$$\dim [\mathbb{P}_k(E)]^2 = 2n_k, \tag{13}$$

$$n_k^\nabla = \dim \mathcal{G}_k^\nabla(E) = n_k + (k + 1), \tag{14}$$

$$n_k^\perp = \dim \mathcal{G}_k^\perp(E) = n_k - (k + 1). \tag{15}$$

Now, following [6], for any integer $k \geq 0$, we define the local mixed virtual element space for the velocity variable \mathbf{u} as

$$V_{h,k}(E) = \{ \mathbf{v}_h \in H(\text{div}; E) \cap H(\text{rot}; E) \text{ s.t. } \mathbf{v}_h \cdot \mathbf{n}_e \in \mathbb{P}_k(e) \ \forall e \in \mathcal{E}_{h,E}, \\ \nabla \cdot \mathbf{v}_h \in \mathbb{P}_k(E), \text{ rot } \mathbf{v}_h \in \mathbb{P}_{k-1}(E) \}. \tag{16}$$

It is easy to see that $[\mathbb{P}_k(E)]^2 \subset V_{h,k}(E)$.

The following set of local degrees of freedom is unisolvent for $V_{h,k}(E)$ (see [5, 7]): given $\mathbf{v}_h \in V_{h,k}(E)$,

- **Edge dofs:** chosen $k + 1$ Gauss quadrature points $\mathbf{x}_i^{e,Q}$ internal on each edge $e \in \mathcal{E}_{h,E}$:

$$\text{dof}_i^e(\mathbf{v}_h) = (\mathbf{v}_h \cdot \mathbf{n}_e)(\mathbf{x}_i^{e,Q}) \quad \forall i = 1, \dots, k + 1. \tag{17}$$

We note that this choice automatically ensures the continuity of the flux $\mathbf{v}_h \cdot \mathbf{n}$ across two adjacent elements.

- **Internal ∇ dofs:**

$$\text{dof}_\alpha^\nabla(\mathbf{v}_h) = \frac{1}{|E|} \int_E \mathbf{v}_h \cdot \mathbf{g}_\alpha^{\nabla,k-1} \quad \forall \alpha = 1, \dots, n_{k-1}^\nabla. \tag{18}$$

- **Internal \perp dofs:**

$$\text{dof}_\alpha^\perp(\mathbf{v}_h) = \frac{1}{|E|} \int_E \mathbf{v}_h \cdot \mathbf{g}_\alpha^{\perp,k} \quad \forall \alpha = 1, \dots, n_k^\perp. \tag{19}$$

Let it be $N_E^{\text{dof}} = \dim V_{h,k}(E) = N^{E,e}(k + 1) + n_{k-1}^\nabla + n_k^\perp$, we denote henceforth the local Lagrangian mixed VE basis corresponding to the defined degrees of freedom:

$$\{ \boldsymbol{\varphi}_i \}_{i=1}^{N_E^{\text{dof}}} = \left\{ \{ \boldsymbol{\varphi}_i^e \}_{i=1}^{k+1} \right\}_{e \in \mathcal{E}_{h,E}} \cup \{ \boldsymbol{\varphi}_i^\nabla \}_{i=1}^{n_{k-1}^\nabla} \cup \{ \boldsymbol{\varphi}_i^\perp \}_{i=1}^{n_k^\perp}, \tag{20}$$

where the dofs numbering first counts the edge dofs, then the internal ∇ dofs and lastly the internal \perp dofs.

As in [6], we define the local mixed virtual element space $Q_{h,k}(E)$ for the pressure variable p as the space of polynomials $\mathbb{P}_k(E)$, that is, $Q_{h,k}(E) = \mathbb{P}_k(E)$. In the next section, we will provide further details regarding the selection of the local basis functions for the local pressure space.

Finally, we define the global mixed virtual element spaces for both velocity and pressure variables as

$$V_{h,k} = \{ \mathbf{v}_h \in H_{0,\Gamma_N}(\text{div}; \Omega) \text{ s.t. } \mathbf{v}_h|_E \in V_{h,k}(E) \ \forall E \in \mathcal{T}_h \}, \tag{21}$$

$$Q_{h,k} = \{ q_h \in L^2(\Omega) \text{ s.t. } q_h|_E \in Q_{h,k}(E) \ \forall E \in \mathcal{T}_h \}. \tag{22}$$

3.1 | The discrete mixed variational formulation

The $L^2(E)$ -projection operator $\Pi_k^0 : V_{h,k} \rightarrow [\mathbb{P}_k(\mathcal{T}_h)]^2$ is locally defined as

$$(\Pi_k^0 \mathbf{v}_h, \mathbf{p}_k)_{0,E} = (\mathbf{v}_h, \mathbf{p}_k)_{0,E} \ \forall \mathbf{p}_k \in [\mathbb{P}_k(E)]^2 \text{ and } \forall E \in \mathcal{T}_h. \tag{23}$$

and, as shown in [5], the projection $\Pi_k^0 \mathbf{v}_h$ of a virtual function $\mathbf{v}_h \in V_{h,k}$ can be explicitly computed from the knowledge of its degrees of freedom (17)–(19).

Now, the local discrete counterpart of the continuous bilinear form

$$a(\mathbf{u}, \mathbf{v}) = (\mathbf{K}\mathbf{u}, \mathbf{v})_{0,\Omega}, \quad \forall \mathbf{u}, \mathbf{v} \in V, \tag{24}$$

reads

$$a_h(\mathbf{u}_h, \mathbf{v}_h) = \sum_{E \in \mathcal{T}_h} a_h^E(\mathbf{u}_h, \mathbf{v}_h) \tag{25}$$

$$= \sum_{E \in \mathcal{T}_h} \left[(\mathbf{K} \Pi_k^0 \mathbf{u}_h, \Pi_k^0 \mathbf{v}_h)_{0,E} + S^E(\mathbf{u}_h, \mathbf{v}_h) \right], \tag{26}$$

where the stabilization term $S^E(\cdot, \cdot)$ is any symmetric and positive definite bilinear form that satisfies, $\forall \mathbf{v}_h \in \mathbf{V}_{h,k}$

$$\alpha_* a_{|E}(\mathbf{v}_h, \mathbf{v}_h) \leq S^E(\mathbf{v}_h, \mathbf{v}_h) \leq \alpha^* a_{|E}(\mathbf{v}_h, \mathbf{v}_h)$$

for some constants $\alpha_*, \alpha^* > 0$ that are depending on \mathbf{K} but independent of h . As in [6, 7], we will choose

$$S^E(\mathbf{u}_h, \mathbf{v}_h) = \bar{\mathbf{K}}|E| \sum_{r=1}^{N_{\text{dof}}^E} \text{dof}_r((I - \Pi_k^0) \mathbf{u}_h) \text{dof}_r((I - \Pi_k^0) \mathbf{v}_h),$$

where $\bar{\mathbf{K}}$ is the largest singular value of \mathbf{K} on E .

Finally, the mixed VEM approximation of (8) is given by:

Find $\mathbf{u}_h = \mathbf{u}_{0,h} + \mathbf{u}_{N,h}$, with $\mathbf{u}_{0,h} \in \mathbf{V}_{h,k}$, and $p_h \in Q_{h,k}$ such that $\forall \mathbf{v}_h \in \mathbf{V}_{h,k}$ and $\forall q_h \in Q_{h,k}$:

$$\begin{cases} a_h(\mathbf{u}_h, \mathbf{v}_h) - (p_h, \nabla \cdot \mathbf{v}_h)_{0,\Omega} - (\beta p_h, \Pi_k^0 \mathbf{v}_h)_{0,\Omega} = -\langle g_D, \mathbf{v}_h \cdot \mathbf{n}_{\Gamma_D} \rangle_{\pm \frac{1}{2}, \Gamma_D} \\ (\nabla \cdot \mathbf{u}_h, q_h)_{0,\Omega} + (\gamma p_h, q_h)_{0,\Omega} = (f, q_h)_{0,\Omega} \end{cases} \tag{27}$$

where $\mathbf{u}_{N,h} \in \{\mathbf{v} \in H(\text{div}; \Omega) : \mathbf{v}|_E \in \mathbf{V}_{h,k}(E) \forall E \in \mathcal{T}_h\}$ is a proper function that satisfies $\text{dof}_i^e(\mathbf{u}_{N,h}) = \text{dof}_i^e(\mathbf{u}_N)$, $\forall i = 1, \dots, k + 1$ and for each edge $e \in \mathcal{E}_h$, belonging to the Neumann boundary Γ_N .

The problem (27) has unique solution $(\mathbf{u}_h, p_h) \in \mathbf{V}_{h,k} \times Q_{h,k}$ and, for h sufficiently small, the following a priori error estimates hold true

$$\|p - p_h\|_0 = O(h^{k+1}), \quad \|\mathbf{u} - \mathbf{u}_h\|_0 = O(h^{k+1}). \tag{28}$$

Furthermore, the following superconvergence result holds true.

Theorem 3.1 (Superconvergence result). *Let p_h the solution to (27) and let $p_I \in Q_{h,k}$ be the interpolant of p . Then, for h sufficiently small,*

$$\|p_I - p_h\|_0 = O(h^{k+2}). \tag{29}$$

4 | POLYNOMIAL BASIS

In this section, we show different procedures for building some polynomial bases for both $\mathbb{P}_k(E)$ and $[\mathbb{P}_k(E)]^2$. In the following, the left superscript will denote the underlying polynomial basis used for the space $\mathbb{P}_k(E)$. In particular, we will use the symbol \mathfrak{m} to indicate the scaled monomial basis (10) and the symbol \mathfrak{q} to refer to the orthogonal basis $\{\mathfrak{q}_\alpha^k\}_{\alpha=1}^{N_k}$ for $\mathbb{P}_k(E)$, whose construction will be detailed in the following. Finally, we will use the symbol \mathfrak{p} to denote a generic polynomial basis $\{\mathfrak{p}_\alpha^k\}_{\alpha=1}^{N_k}$ for $\mathbb{P}_k(E)$.

As mentioned in the previous section, the standard choice for the polynomial basis is the set of scaled monomials, defined in (10) (see [3, 4, 7]). Orthogonal polynomial bases have already been

introduced in the virtual element literature, since it has been proven that modifying the definition of internal moments by choosing an $L^2(E)$ -orthonormal basis for $\mathbb{P}_k(E)$ can largely improve the reliability of the virtual element method for higher order approximations and in the presence of badly-shaped polygons ([9, 15]).

An efficient procedure for building an $L^2(E)$ -orthonormal polynomial basis for $\mathbb{P}_k(E)$ is presented in [2]. This procedure consists in the application of the modified Gram-Schmidt (MGS) orthogonalization process to the monomial Vandermonde matrix $\mathbf{m}\mathbf{V}^{E,k} \in \mathbb{R}^{N^{E,Q} \times n_k}$ associated with a quadrature formula $\left\{ \left(\mathbf{x}_i^{E,Q}, \omega_i^{E,Q} \right) \right\}_{i=1}^{N^{E,Q}}$ on a given element E with $N^{E,Q}$ nodes. As suggested in [2, 13], the process must be applied twice in order to make the orthonormalization error $\|\mathbf{I} - \mathfrak{q}\mathbf{H}^{E,k}\|$ independent of the condition number of $\mathbf{m}\mathbf{V}^{E,k}$ matrix, where $\mathfrak{q}\mathbf{H}^{E,k} \in \mathbb{R}^{n_k \times n_k}$ is the mass matrix of the resulting $L^2(E)$ -orthonormal polynomial basis $\{\mathfrak{q}\alpha^k\}_{\alpha=1}^{n_k}$ for $\mathbb{P}_k(E)$.

In particular, the whole procedure defines a matrix $\mathbf{L}^{E,k} \in \mathbb{R}^{n_k \times n_k}$ on each element E such that

$$\mathfrak{q}\mathbf{V}^{E,k} = \mathbf{m}\mathbf{V}^{E,k} (\mathbf{L}^{E,k})^T, \quad (30)$$

where $\mathfrak{q}\mathbf{V}^{E,k}$ is the Vandermonde matrix associated with the new polynomial basis. More precisely, we first apply the MGS process to the $\mathbf{m}\mathbf{V}^{E,k}$ matrix, that is, we define an upper triangular matrix $\mathbf{R}_1^{E,k} \in \mathbb{R}^{n_k \times n_k}$ and an orthonormal matrix $\mathbf{Q}_1^{E,k} \in \mathbb{R}^{N^{E,Q} \times n_k}$ such that

$$\mathbf{m}\mathbf{V}^{E,k} = \mathbf{Q}_1^{E,k} \mathbf{R}_1^{E,k}.$$

Then we apply MGS to the $\mathbf{Q}_1^{E,k}$ matrix properly rescaled by quadrature weights in order to obtain an $L^2(E)$ -orthonormal basis:

$$(\mathbf{Z}^{E,Q})^{1/2} \mathbf{Q}_1^{E,k} = \mathbf{Q}_2^{E,k} \mathbf{R}_2^{E,k},$$

where $\mathbf{Z}^{E,Q} \in \mathbb{R}^{N^{E,Q} \times N^{E,Q}}$ is the diagonal matrix whose diagonal entries are the quadrature weights. Being $\mathbf{L}^{E,k} = (\mathbf{R}_2^{E,k} \mathbf{R}_1^{E,k})^{-T}$, we note that

$$\begin{aligned} \mathbf{Q}_2^{E,k} &= (\mathbf{Z}^{E,Q})^{1/2} \mathbf{Q}_1^{E,k} (\mathbf{R}_2^{E,k})^{-1} \\ &= (\mathbf{Z}^{E,Q})^{1/2} \mathbf{m}\mathbf{V}^{E,k} (\mathbf{L}^{E,k})^T \\ &= (\mathbf{Z}^{E,Q})^{1/2} \mathfrak{q}\mathbf{V}^{E,k} \end{aligned}$$

which means that $\mathbf{Q}_2^{E,k}$ is the Vandermonde matrix $\mathfrak{q}\mathbf{V}^{E,k}$ rescaled by the square root of the quadrature weights, thus,

$$\mathfrak{q}\mathbf{H}^{E,k} = (\mathfrak{q}\mathbf{V}^{E,k})^T \mathbf{Z}^{E,Q} \mathfrak{q}\mathbf{V}^{E,k} = (\mathbf{Q}_2^{E,k})^T \mathbf{Q}_2^{E,k} = \mathbf{I}. \quad (31)$$

Remark 1. Note that $\mathfrak{q}\mathbf{V}^{E,k} = \mathbf{m}\mathbf{V}^{E,k} (\mathbf{L}^{E,k})^T$ means that the columns of $(\mathbf{L}^{E,k})^T$ are the coefficients that provide each orthonormal polynomial as a linear combination of the scaled monomials.

4.1 | Polynomial bases for $[\mathbb{P}_k(E)]^2$

For simplicity, in the following, we will drop the superscript E when no ambiguity occurs. Now, we introduce some auxiliary matrices that we will use in the following. Let ${}^{\mathfrak{p}}\mathbf{D}_x^{k+1}, {}^{\mathfrak{p}}\mathbf{D}_y^{k+1} \in \mathbb{R}^{n_{k+1} \times n_k}$ be the matrices which collect the coefficients of the partial derivatives of $\{\mathfrak{p}\alpha^{k+1}\}_{\alpha=1}^{n_{k+1}}$, that is,

$$\frac{\partial \mathfrak{p}\alpha^{k+1}}{\partial x} = \sum_{\beta=1}^{n_k} ({}^{\mathfrak{p}}\mathbf{D}_x^{k+1})_{\alpha\beta} \mathfrak{p}\beta^k, \quad \frac{\partial \mathfrak{p}\alpha^{k+1}}{\partial y} = \sum_{\beta=1}^{n_k} ({}^{\mathfrak{p}}\mathbf{D}_y^{k+1})_{\alpha\beta} \mathfrak{p}\beta^k, \quad \forall \alpha = 1, \dots, n_{k+1}. \quad (32)$$

Note that, if the MGS basis is used, matrices ${}^{\mathbb{Q}}\mathbf{D}_*^{k+1}$ can be derived from the (easily computable) monomial ones ${}^{\mathbb{M}}\mathbf{D}_*^{k+1}$ as:

$${}^{\mathbb{Q}}\mathbf{D}_*^{k+1} = \mathbf{L}^{k+1} {}^{\mathbb{M}}\mathbf{D}_*^{k+1} (\mathbf{L}^k)^{-1} \quad \forall * \in \{x, y\},$$

since, $\forall \alpha = 1, \dots, n_{k+1}$,

$$\begin{aligned} \frac{\partial {}^{\mathbb{Q}}\alpha^{k+1}}{\partial *}&= \sum_{\beta=1}^{n_{k+1}} \mathbf{L}_{\alpha\beta}^{k+1} \frac{\partial {}^{\mathbb{M}}\beta^{k+1}}{\partial *}&= \sum_{\beta=1}^{n_{k+1}} \sum_{\gamma=1}^{n_k} \mathbf{L}_{\alpha\beta}^{k+1} ({}^{\mathbb{M}}\mathbf{D}_*^{k+1})_{\beta\gamma} {}^{\mathbb{M}}\gamma^k \\ &= \sum_{\beta=1}^{n_{k+1}} \sum_{\gamma=1}^{n_k} \sum_{s=1}^{n_k} \mathbf{L}_{\alpha\beta}^{k+1} ({}^{\mathbb{M}}\mathbf{D}_*^{k+1})_{\beta\gamma} (\mathbf{L}^k)_{\gamma s}^{-1} {}^{\mathbb{Q}}\alpha^k. \end{aligned}$$

Remark 2. As highlighted in [2], the modified Gram-Schmidt algorithm allows us to obtain a hierarchical sequence of bases $\{\{q_\alpha^k\}_{\alpha=1}^{n_k}\}_{k \geq 0}$, that is,

$$\{q_\alpha^k\}_{\alpha=1}^{n_k} \subset \{q_\alpha^{k+1}\}_{\alpha=1}^{n_{k+1}}. \tag{33}$$

As a consequence, we only need to compute \mathbf{L}^{k+1} and then we can set

$$\mathbf{L}^k = \mathbf{L}^{k+1}(1 : n_k, 1 : n_k), \tag{34}$$

being $\mathbf{L}^{k+1}(1 : n_k, 1 : n_k)$ the matrix obtained from the first n_k rows and columns of \mathbf{L}^{k+1} . Indeed, we define the Vandermonde matrices of both k and $k + 1$ orders with respect to the same quadrature formula.

Starting from a generic polynomial basis for $\mathbb{P}_k(E)$, as shown in [7], an easily computable basis $\{p_I^k\}_{I=1}^{2n_k}$ for $[\mathbb{P}_k(E)]^2$ can be built as

$$p_I^k = \begin{cases} \begin{bmatrix} \mathbb{P}_I^k \\ 0 \end{bmatrix} & I = 1, \dots, n_k \\ \begin{bmatrix} 0 \\ \mathbb{P}_{I-n_k}^k \end{bmatrix} & I = n_k + 1, \dots, 2n_k \end{cases} \tag{35}$$

The Vandermonde matrix ${}^p\mathbf{V}^k \in \mathbb{R}^{2N^Q \times 2n_k}$ associated with the $\{p_I^k\}_{I=1}^{2n_k}$ polynomial basis functions can be written as

$${}^p\mathbf{V}^k = \begin{bmatrix} {}^p\mathbf{V}^k & \mathbf{O} \in \mathbb{R}^{N^Q \times n_k} \\ \mathbf{O} \in \mathbb{R}^{N^Q \times n_k} & {}^p\mathbf{V}^k \end{bmatrix}, \tag{36}$$

where

- the top left ${}^p\mathbf{V}^k$ matrix contains the evaluations of the x -components of p_I^k for all $I = 1, \dots, n_k$;
- the top right \mathbf{O} matrix contains the evaluations of the y -components of p_I^k for all $I = 1, \dots, n_k$, that are all zeros;
- the bottom left \mathbf{O} matrix contains the evaluations of the x -components of p_I^k for all $I = n_k + 1, \dots, 2n_k$, that are all zeros;
- the bottom right ${}^p\mathbf{V}^k$ matrix contains the evaluations of the y -components of p_I^k for all $I = n_k + 1, \dots, 2n_k$.

Note that if $\{p_\alpha^k\}_{\alpha=1}^{n_k}$ is an $L^2(E)$ -orthonormal polynomial basis for $\mathbb{P}_k(E)$, then $\{p_I^k\}_{I=1}^{2n_k}$ is an $L^2(E)$ -orthonormal polynomial basis for $[\mathbb{P}_k(E)]^2$.

Now, functions belonging to $\mathcal{G}_k^\nabla(E)$, as defined in (11), can be written as

$$\mathbf{g}_\alpha^{\nabla,k} = \nabla \mathbb{p}_{\alpha+1}^{k+1} = \sum_{l=1}^{2n_k} \mathbf{T}_{\alpha l}^{\nabla,k} \mathbf{p}_l^k \quad \forall \alpha = 1, \dots, n_k^\nabla, \tag{37}$$

where $\mathbf{T}^{\nabla,k} \in \mathbb{R}^{n_k^\nabla \times 2n_k}$ is the coefficient matrix of gradients of the polynomial functions $\{\mathbb{p}_\alpha^{k+1}\}_{\alpha=2}^{n_{k+1}}$ with respect to the polynomial basis $\{\mathbf{p}_l^k\}_{l=1}^{2n_k}$ of $[\mathbb{P}_k(E)]^2$ and the corresponding Vandermonde matrix is

$$\mathbf{g}^{\nabla,k} = \mathbf{p} \mathbf{V}^k (\mathbf{T}^{\nabla,k})^T. \tag{38}$$

Based on definitions (32), the $\mathbf{T}^{\nabla,k}$ matrix reads

$$\mathbf{T}^{\nabla,k} = \left[\mathbb{D}_x^{k+1}(2 : n_{k+1}, :) \quad \mathbb{D}_y^{k+1}(2 : n_{k+1}, :) \right],$$

with $\mathbb{D}_*^{k+1}(2 : n_{k+1}, :)$ the sub-matrix of \mathbb{D}_*^{k+1} obtained extracting rows from 2 to n_{k+1} and all columns.

Now, we can complete a basis $\mathcal{G}_k(E)$ for $[\mathbb{P}_k(E)]^2$ by adding the set of functions $\mathcal{G}_k^\perp(E) = \{\mathbf{g}_\alpha^{\perp,k}\}_{\alpha=1}^{n_k^\perp}$ defined in such a way (12) is satisfied. As suggested in [7], $\mathbf{g}_\alpha^{\perp,k}$ function can be defined as

$$\mathbf{g}_\alpha^{\perp,k} = \sum_{l=1}^{2n_k} \mathbf{T}_{\alpha l}^{\perp,k} \mathbf{p}_l^k \quad \forall \alpha = 1, \dots, n_k^\perp, \tag{39}$$

where $\mathbf{T}^{\perp,k} \in \mathbb{R}^{n_k^\perp \times 2n_k}$ is the matrix whose rows define an *euclidean* orthonormal basis for the nullspace of $\mathbf{T}^{\nabla,k}$ matrix. Thus, by considering the Singular Value Decomposition of $\mathbf{T}^{\nabla,k} = \mathbf{U} \mathbf{\Sigma}(\mathbf{V})^T$, we can define $\mathbf{T}^{\perp,k}$ as

$$\mathbf{T}^{\perp,k} = [\mathbf{V}(:, n_k^\nabla + 1 : 2n_k)]^T, \tag{40}$$

where $\mathbf{V}(:, n_k^\nabla + 1 : 2n_k)$ is the submatrix of \mathbf{V} made up of all its rows and of the columns running from the $(n_k^\nabla + 1)$ -th to the $2n_k$ -th. As a consequence,

$$\mathbf{T}^{\nabla,k} (\mathbf{T}^{\perp,k})^T = \mathbf{O}. \tag{41}$$

The Vandermonde matrix $\mathbf{g}^{\mathbf{V}^k} \in \mathbb{R}^{2N^Q \times 2n_k}$ associated with the basis functions $\mathcal{G}_k(E) = \{\mathbf{g}_l^k\}_{l=1}^{2n_k} = \{\mathbf{g}_\alpha^{\nabla,k}\}_{\alpha=1}^{n_k^\nabla} \cup \{\mathbf{g}_\beta^{\perp,k}\}_{\beta=1}^{n_k^\perp}$ reads

$$\mathbf{g}^{\mathbf{V}^k} = \begin{bmatrix} \mathbf{g}^{\mathbf{V}^{\nabla,k}} & \mathbf{g}^{\mathbf{V}^{\perp,k}} \end{bmatrix} = \mathbf{p} \mathbf{V}^k \begin{bmatrix} (\mathbf{T}^{\nabla,k})^T & (\mathbf{T}^{\perp,k})^T \end{bmatrix}. \tag{42}$$

We define $\mathbf{G}^k \in \mathbb{R}^{2n_k \times 2n_k}$ as the mass matrix related to the $\mathcal{G}_k(E)$ basis, whose entries are given by

$$\mathbf{G}_{IJ}^k = \int_E \mathbf{g}_I^k \cdot \mathbf{g}_J^k, \quad \forall I, J = 1, \dots, 2n_k. \tag{43}$$

In matrix form, \mathbf{G}^k reads

$$\begin{aligned} \mathbf{G}^k &= (\mathbf{g}^{\mathbf{V}^k})^T \begin{bmatrix} \mathbf{Z}^Q & \mathbf{O} \\ \mathbf{O} & \mathbf{Z}^Q \end{bmatrix} \mathbf{g}^{\mathbf{V}^k} \\ &= \begin{bmatrix} \mathbf{T}^{\nabla,k} \mathbf{p} \mathbf{H}^k (\mathbf{T}^{\nabla,k})^T & \mathbf{T}^{\nabla,k} \mathbf{p} \mathbf{H}^k (\mathbf{T}^{\perp,k})^T \\ \mathbf{T}^{\perp,k} \mathbf{p} \mathbf{H}^k (\mathbf{T}^{\nabla,k})^T & \mathbf{T}^{\perp,k} \mathbf{p} \mathbf{H}^k (\mathbf{T}^{\perp,k})^T \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{G}^{\nabla,\nabla} & \mathbf{G}^{\nabla,\perp} \\ \mathbf{G}^{\perp,\nabla} & \mathbf{G}^{\perp,\perp} \end{bmatrix}, \end{aligned} \tag{44}$$

where

$${}^p\mathbf{H}^k = ({}^p\mathbf{V}^k)^T \begin{bmatrix} \mathbf{Z}^{\mathcal{Q}} & \mathbf{O} \\ \mathbf{O} & \mathbf{Z}^{\mathcal{Q}} \end{bmatrix} {}^p\mathbf{V}^k$$

is the mass matrix related to $\{\mathbf{p}_l\}_{l=1}^{2n_k}$ basis.

Now, we make the following observations.

- 1 If we choose the standard set of scaled monomials $\mathcal{M}_k(E)$ as the basis for $\mathbb{P}_k(E)$ and we set

$${}^1\mathbf{T}^{\nabla,k} = \left[\mathbb{m}\mathcal{D}_x^{k+1}(2 : n_{k+1}, :) \quad \mathbb{m}\mathcal{D}_y^{k+1}(2 : n_{k+1}, :) \right] = {}^1\mathbf{U} {}^1\mathbf{\Sigma} ({}^1\mathbf{V})^T, \quad {}^1\mathbf{T}^{\perp,k} = [{}^1\mathbf{V}(:, n_k^{\nabla} + 1 : 2n_k)]^T,$$

it is known that the corresponding ${}^1\mathbf{G}^k$ matrix, that is,

$${}^1\mathbf{G}^k = \begin{bmatrix} {}^1\mathbf{T}^{\nabla,k} \mathbf{m}\mathbf{H}^k ({}^1\mathbf{T}^{\nabla,k})^T & {}^1\mathbf{T}^{\nabla,k} \mathbf{m}\mathbf{H}^k ({}^1\mathbf{T}^{\perp,k})^T \\ {}^1\mathbf{T}^{\perp,k} \mathbf{m}\mathbf{H}^k ({}^1\mathbf{T}^{\nabla,k})^T & {}^1\mathbf{T}^{\perp,k} \mathbf{m}\mathbf{H}^k ({}^1\mathbf{T}^{\perp,k})^T \end{bmatrix} \quad (45)$$

will be ill-conditioned for high polynomial degrees. Anyway, we want to highlight that, in this so-called *monomial approach*, the rows of ${}^1\mathbf{T}^{\perp,k}$ are orthonormal to each other and are orthogonal to the rows of ${}^1\mathbf{T}^{\nabla,k}$ with respect to the euclidean scalar product, by construction.

- 2 On the other hand, if we choose the MGS basis $\{\mathfrak{q}_\alpha^k\}_{\alpha=1}^{n_k}$, that is, the $L^2(E)$ -orthonormal polynomial basis for $\mathbb{P}_k(E)$ introduced in [2], the corresponding mass matrix ${}^q\mathbf{H}^k$ will be the identity matrix \mathbf{I} . Thus, by setting

$${}^2\mathbf{T}^{\nabla,k} = \left[\mathfrak{q}\mathcal{D}_x^{k+1}(2 : n_{k+1}, :) \quad \mathfrak{q}\mathcal{D}_y^{k+1}(2 : n_{k+1}, :) \right] = {}^2\mathbf{U} {}^2\mathbf{\Sigma} ({}^2\mathbf{V})^T, \quad {}^2\mathbf{T}^{\perp,k} = [{}^2\mathbf{V}(:, n_k^{\nabla} + 1 : 2n_k)]^T, \quad (46)$$

the related ${}^2\mathbf{G}^k$ matrix, in infinite precision, is given by

$${}^2\mathbf{G}^k = \begin{bmatrix} {}^2\mathbf{T}^{\nabla,k} \mathbf{I} ({}^2\mathbf{T}^{\nabla,k})^T & {}^2\mathbf{T}^{\nabla,k} \mathbf{I} ({}^2\mathbf{T}^{\perp,k})^T \\ {}^2\mathbf{T}^{\perp,k} \mathbf{I} ({}^2\mathbf{T}^{\nabla,k})^T & {}^2\mathbf{T}^{\perp,k} \mathbf{I} ({}^2\mathbf{T}^{\perp,k})^T \end{bmatrix} = \begin{bmatrix} {}^2\mathbf{T}^{\nabla,k} ({}^2\mathbf{T}^{\nabla,k})^T & \mathbf{O} \\ \mathbf{O} & \mathbf{I} \end{bmatrix}, \quad (47)$$

where the last equation is a consequence of property (41). In conclusion, if an $L^2(E)$ -orthonormal polynomial basis for $\mathbb{P}_k(E)$ is used, ${}^2\mathcal{G}_k^k(E)$ will be partially $L^2(E)$ -orthonormalized, since the ${}^2\mathbf{g}_\alpha^{\perp,k}$ functions are a set of $L^2(E)$ -orthonormal functions and are $L^2(E)$ -orthogonal to ${}^2\mathbf{g}_\alpha^{\nabla,k}$ functions, but the ${}^2\mathbf{g}_\alpha^{\nabla,k}$ are not naturally $L^2(E)$ -orthogonal to each other. We denote this as *partial-orthonormal approach*.

4.1.1 | A full $L^2(E)$ -orthonormal approach

In this section, we will show a procedure that allows to full $L^2(E)$ -orthonormalize the $\mathcal{G}_k(E)$ basis. For this purpose, we first note from (47) that in order to make $\mathcal{G}_k^{\nabla}(E)$ an $L^2(E)$ -orthonormal set of functions, it is sufficient to orthonormalize the rows of ${}^2\mathbf{T}^{\nabla,k}$ matrix defined in (46) with respect to the euclidean scalar product. In order to do this, we apply only once the modified Gram-Schmidt algorithm to $({}^2\mathbf{T}^{\nabla,k})^T$. More precisely, we factorize $({}^2\mathbf{T}^{\nabla,k})^T$ as

$$({}^2\mathbf{T}^{\nabla,k})^T = \mathbf{Q}^{\nabla,k} \mathbf{R}^{\nabla,k},$$

and, then we set

$${}^3\mathbf{T}^{\nabla,k} = (\mathbf{Q}^{\nabla,k})^T = \mathbf{L}^{\nabla,k} {}^2\mathbf{T}^{\nabla,k}, \quad {}^3\mathbf{T}^{\perp,k} = {}^3\mathbf{V}(:, n_k^{\nabla} + 1 : 2n_k)^T, \quad (48)$$

where ${}^3\mathbf{T}^{\nabla,k} = {}^3\mathbf{U} {}^3\mathbf{\Sigma} {}^3\mathbf{V}^T$ and $\mathbf{L}^{\nabla,k} = (\mathbf{R}^{\nabla,k})^{-T} \in \mathbb{R}^{n_k^{\nabla} \times n_k^{\nabla}}$.

By proceeding in this way, the ${}^3\mathbf{G}^k$ matrix will become the identity matrix. We will call this procedure as *full-orthonormal approach*.

Remark 3. It is worth mentioning that in order to orthonormalize the rows of ${}^2\mathbf{T}^{\nabla,k}$ we could resort to its computed Singular Value Decomposition (46) and set

$$\begin{aligned} {}^3\mathbf{T}^{\nabla,k} &= {}^2\mathbf{V}(:, 1 : n_k^{\nabla})^T = ({}^2\mathbf{\Sigma}(:, 1 : n_k^{\nabla}))^{-1} {}^2\mathbf{U}^T {}^2\mathbf{T}^{\nabla,k}, \\ {}^3\mathbf{T}^{\perp,k} &= {}^2\mathbf{V}(:, n_k^{\nabla} + 1 : 2n_k)^T = {}^2\mathbf{T}^{\perp,k}. \end{aligned}$$

In this way, we obtain that ${}^3\mathbf{G}^k = \mathbf{I}$, but the SVD process is not hierarchical and, as a consequence, the set of functions $\{{}^3\mathbf{g}_{\alpha}^{\nabla,k-1}\}_{\alpha=1}^{n_k^{\nabla}}$ used for defining the internal ∇ dofs is not $L^2(E)$ -orthonormalized for free. As highlighted in [9, 15], this type of change can improve only the condition number of elemental matrices, but does not ensure to improve the global performance of the method.

Instead, since the MGS is a hierarchical procedure (see Remark 2), we obtain that the functions

$$\{{}^3\mathbf{g}_{\alpha}^{\nabla,k-1}\}_{\alpha=1}^{n_k^{\nabla}} \subset \{{}^3\mathbf{g}_{\alpha}^{\nabla,k}\}_{\alpha=1}^{n_k^{\nabla}} \quad (49)$$

used in the definition of internal ∇ dofs are a set of $L^2(E)$ -orthonormal functions.

5 | SOME IMPLEMENTATION DETAILS

In this section, we will show how to compute the local matrices needed for assembling the local system matrix related to the discrete problem (27) with the aforementioned approaches, following a procedure similar to the one shown in [7].

5.1 | The divergence term $(\nabla \cdot \mathbf{u}_h, q_h)_{0,E}$

By choosing $\mathbf{u}_h = \boldsymbol{\varphi}_i$ and $q_h = \mathbb{p}_{\alpha}^k$, we define $\mathbf{W} \in \mathbb{R}^{n_k \times N_E^{\text{dof}}}$ as the matrix whose entries read, $\forall i = 1, \dots, N_E^{\text{dof}}$, $\alpha = 1, \dots, n_k$,

$$\mathbf{W}_{ai} = \int_E \nabla \cdot \boldsymbol{\varphi}_i \mathbb{p}_{\alpha}^k \quad (50)$$

$$\begin{aligned} &= - \int_E \boldsymbol{\varphi}_i \cdot \nabla \mathbb{p}_{\alpha}^k + \int_{\partial E} \boldsymbol{\varphi}_i \cdot \mathbf{n}_{\partial E} \mathbb{p}_{\alpha}^k \\ &:= (\mathbf{W}_1)_{ai} + (\mathbf{W}_2)_{ai}. \end{aligned} \quad (51)$$

Since in the monomial and partial-orthonormal approaches we have

$$\begin{aligned} \mathbb{p} &= \mathfrak{m}, & \nabla \mathfrak{m}_{\alpha}^k &= \mathbf{1} \mathbf{g}_{\alpha-1}^{\nabla,k-1} & \forall \alpha &= 2, \dots, n_k, \\ \mathbb{p} &= \mathfrak{q}, & \nabla \mathfrak{q}_{\alpha}^k &= \mathbf{2} \mathbf{g}_{\alpha-1}^{\nabla,k-1} & \forall \alpha &= 2, \dots, n_k, \end{aligned} \quad (52)$$

by exploiting the internal ∇ degrees of freedom (18), we have

$$\mathbf{W}_1 = \begin{bmatrix} \mathbf{O} \in \mathbb{R}^{1 \times N_E^{\text{dof}}} & & \\ \mathbf{O} \in \mathbb{R}^{n_k^{\nabla} \times N_E^{\text{dof}}(k+1)} & -|E| \mathbf{I} \in \mathbb{R}^{n_k^{\nabla} \times n_k^{\nabla}} & \mathbf{O} \in \mathbb{R}^{n_k^{\nabla} \times n_k^{\nabla}} \end{bmatrix}.$$

Concerning the full-orthonormal approach, we observe that

$$\begin{aligned} \mathbb{p} &= \mathbb{q}, \quad \nabla \mathbb{q}_{\alpha+1}^k = \sum_{l=1}^{2n_{k-1}} 2 \mathbf{T}_{\alpha l}^{\nabla, k} \mathbf{q}_l^{k-1} = \sum_{l=1}^{2n_{k-1}} \sum_{\beta=1}^{n_{k-1}^{\nabla}} (\mathbf{L}^{\nabla, k-1})_{\alpha\beta}^{-1} 3 \mathbf{T}_{\beta l}^{\nabla, k} \mathbf{q}_l^{k-1} \\ &= \sum_{\beta=1}^{n_{k-1}^{\nabla}} (\mathbf{L}^{\nabla, k-1})_{\alpha\beta}^{-1} 3 \mathbf{g}_{\beta}^{\nabla, k-1}, \quad \forall \alpha = 1, \dots, n_{k-1}^{\nabla}. \end{aligned} \tag{53}$$

Thus, thanks to the hierarchical property shown in Remark 3, the \mathbf{W}_1 matrix becomes

$$\mathbf{W}_1 = \begin{bmatrix} \mathbf{O} \in \mathbb{R}^{1 \times \Lambda_E^{\text{dof}}} \\ \mathbf{O} \in \mathbb{R}^{n_{k-1}^{\nabla} \times N^{E,e}(k+1)} & -|E| (\mathbf{L}^{\nabla, k})^{-1} (1 : n_{k-1}^{\nabla}, 1 : n_{k-1}^{\nabla}) & \mathbf{O} \in \mathbb{R}^{n_{k-1}^{\nabla} \times n_k^{\perp}} \end{bmatrix}.$$

Concerning the second term in the right-hand side of Equation (51), we observe that the integrand function is a known polynomial of degree $2k$ on each edge of E , which can be computed exactly (up to machine precision) by exploiting the edge dofs (17). Then, the second term \mathbf{W}_2 reads

$$\mathbf{W}_2 = \left[(\mathbb{p} \mathbf{V}^{\partial, k})^T \mathbf{Z}^{\partial, Q} \quad \mathbf{O} \in \mathbb{R}^{n_k \times (n_{k-1}^{\nabla} + n_k^{\perp})} \right],$$

where $\mathbb{p} \mathbf{V}^{\partial, k} \in \mathbb{R}^{N^{E,e}(k+1) \times n_k}$ is the Vandermonde matrix related to the polynomials $\{\mathbb{p}_{\alpha}^k\}_{\alpha=1}^{n_k}$ and the boundary quadrature formula $\left\{ \left\{ \left(\mathbf{x}_i^{e, Q}, \omega_i^{e, Q} \right) \right\}_{i=1}^{k+1} \right\}_{e \in \mathcal{E}_{h,E}}$ defined on ∂E , whose quadrature nodes coincide with the nodes used in the edge dofs (17) and the matrix $\mathbf{Z}^{\partial, Q} \in \mathbb{R}^{N^{E,e}(k+1) \times N^{E,e}(k+1)}$ is the diagonal matrix whose non-zero entries coincide with the related boundary quadrature weights properly arranged.

5.2 | Computation of $L^2(E)$ -projection of basis functions of $V_{h,k}(E)$

Now, we compute the projection $\Pi_k^0 \varphi_i$ of the Lagrangian basis functions of $V_{h,k}(E)$ in terms of the $\mathcal{G}_k(E)$ basis, that is, we determine the matrix $\mathbf{\Pi}_k^0 \in \mathbb{R}^{2n_k \times N_E^{\text{dof}}}$ such that

$$\Pi_k^0 \varphi_i = \sum_{l=1}^{2n_k} (\mathbf{\Pi}_k^0)_{li} \mathbf{g}_l^k. \tag{54}$$

By replacing (54) in the definition (23) of Π_0^k , we obtain, $\forall i = 1, \dots, N_E^{\text{dof}}, \forall J = 1, \dots, 2n_k$,

$$\sum_{l=1}^{2n_k} (\mathbf{\Pi}_k^0)_{li} (\mathbf{g}_l^k, \mathbf{g}_J^k)_{0,E} = (\varphi_i, \mathbf{g}_J^k)_{0,E}.$$

or, analogously, in matrix form

$$\mathbf{G}^k \mathbf{\Pi}_k^0 = \mathbf{B}, \tag{55}$$

where $\mathbf{B} \in \mathbb{R}^{2n_k \times N_E^{\text{dof}}}$ is the matrix whose entries are defined as

$$\mathbf{B}_{Ji} = (\varphi_i, \mathbf{g}_J^k)_{0,E}, \quad \forall i = 1, \dots, N_E^{\text{dof}}, \forall J = 1, \dots, 2n_k.$$

Now, as in [7], we split \mathbf{B} as

$$\mathbf{B} = \begin{bmatrix} \mathbf{B}^{\nabla} \\ \mathbf{B}^{\perp} \end{bmatrix}.$$

The term $\mathbf{B}^\perp \in \mathbb{R}^{n_k^\perp \times N_E^{\text{dof}}}$, whose entries are

$$\mathbf{B}_{\alpha i}^\perp = (\boldsymbol{\varphi}_i, \mathbf{g}_\alpha^{\perp, k})_{0, E}, \quad \forall \alpha = 1, \dots, n_k^\perp, \quad \forall i = 1, \dots, N_E^{\text{dof}},$$

can be readily computed as

$$\mathbf{B}^\perp = \left[\mathbf{O} \in \mathbb{R}^{n_k^\perp \times (N_E^{\text{e}}(k+1) + n_{k-1}^\perp)} \quad |E| \mathbf{I} \in \mathbb{R}^{n_k^\perp \times n_k^\perp} \right],$$

by exploiting the internal \perp degrees of freedom (19).

Regarding the first term $\mathbf{B}^\nabla \in \mathbb{R}^{n_k^\nabla \times N_E^{\text{dof}}}$, we note that, $\forall \alpha = 1, \dots, n_k^\nabla, \forall i = 1, \dots, N_E^{\text{dof}}$

- in the monomial and in the partial-orthonormal approaches, recalling Equation (52),

$$\begin{aligned} \mathbf{B}_{\alpha i}^\nabla &= \int_E \boldsymbol{\varphi}_i \cdot \mathbf{g}_\alpha^{\nabla, k} = \int_E \boldsymbol{\varphi}_i \cdot \nabla \mathbb{p}_{\alpha+1}^{k+1} \\ &= - \int_E \nabla \cdot \boldsymbol{\varphi}_i \mathbb{p}_{\alpha+1}^{k+1} + \int_{\partial E} \boldsymbol{\varphi}_i \cdot \mathbf{n}_{\partial E} \mathbb{p}_{\alpha+1}^{k+1} \\ &:= (\mathbf{B}_1^\nabla)_{\alpha i} + (\mathbf{B}_2^\nabla)_{\alpha i}, \end{aligned} \quad (56)$$

- in the full-orthonormal approach, evoking Equation (53), we have

$$\mathbf{B}_{\alpha i}^\nabla = \int_E \boldsymbol{\varphi}_i \cdot \mathbf{g}_\alpha^{\nabla, k} = \sum_\beta \mathbf{L}_{\alpha\beta}^{\nabla, k} \int_E \boldsymbol{\varphi}_i \cdot \nabla \mathbb{p}_{\beta+1}^{k+1} = (\mathbf{L}^{\nabla, k} \mathbf{B}_1^\nabla)_{\alpha i} + (\mathbf{L}^{\nabla, k} \mathbf{B}_2^\nabla)_{\alpha i}.$$

Since the integrand of the second term of (56) is a known polynomial of degree $2k + 1$ on each edge of E , it can be integrated exactly by exploiting the edge dofs (17). Thus,

$$\mathbf{B}_2^\nabla = \left[(\mathbb{p} \mathbf{V}^{\partial, k+1}(:, 2 : n_{k+1}))^T \mathbf{Z}^{\partial, Q} \quad \mathbf{O} \in \mathbb{R}^{n_k^\nabla \times (n_{k-1}^\nabla + n_k^\perp)} \right].$$

In order to compute the term \mathbf{B}_1^∇ , we first note that $\nabla \cdot \boldsymbol{\varphi}_i$ is a known polynomial of degree k on E . Then, it can be written as

$$\nabla \cdot \boldsymbol{\varphi}_i = \sum_{\alpha=1}^{n_k} \Lambda_{\alpha i} \mathbb{p}_\alpha^k, \quad \forall i = 1, \dots, N_E^{\text{dof}}.$$

whose coefficient matrix $\Lambda \in \mathbb{R}^{n_k \times N_E^{\text{dof}}}$ can be retrieved from

$$\sum_{\alpha=1}^{n_k} \Lambda_{\alpha i} \int_E \mathbb{p}_\alpha^k \mathbb{p}_\beta^k = \int_E \nabla \cdot \boldsymbol{\varphi}_i \mathbb{p}_\beta^k, \quad \forall \beta = 1, \dots, n_k, \quad \forall i = 1, \dots, N_E^{\text{dof}}. \quad (57)$$

Equation (57) can be written in matrix form as

$$\mathbb{p} \mathbf{H}^k \Lambda = \mathbf{W}, \quad (58)$$

since the term at right hand side of (57) coincides with the definition of the \mathbf{W} matrix given in Section 5.1. We highlight that, in the partial and in the full orthonormal approaches, in infinite precision, it holds

$$\Lambda = \mathbf{W}.$$

Finally, matrix \mathbf{B}_1^∇ is given by

$$(\mathbf{B}_1^\nabla)_{\alpha i} = - \int_E \nabla \cdot \boldsymbol{\varphi}_i \mathbb{p}_{\alpha+1}^{k+1} = - \sum_{\beta=1}^{n_k} \Lambda_{\beta i} \int_E \mathbb{p}_\beta^k \mathbb{p}_{\alpha+1}^{k+1},$$

or, equivalently,

$$\mathbf{B}_1^\nabla = -\mathbb{P}\mathbf{H}^{k+1}(2 : n_{k+1}, 1 : n_k)\mathbf{\Lambda} = -\mathbb{P}\mathbf{H}^{k+1}(2 : n_{k+1}, 1 : n_k)(\mathbb{P}\mathbf{H}^k)^{-1}\mathbf{W}.$$

5.3 | The diffusion term $a_h^E(\mathbf{u}_h, \mathbf{v}_h)$

The local discrete diffusion term, defined in Equation (26), is the sum of a consistency and a stability term. Thus, as in [7], we can separately define the consistency matrix

$$(\mathbf{K}_C^a)_{ij} = (\mathbf{K}\Pi_0^k \boldsymbol{\varphi}_i, \Pi_0^k \boldsymbol{\varphi}_j)_{0,E}$$

and the stability matrix

$$(\mathbf{K}_S^a)_{ij} = \bar{K}|E| \sum_{r=1}^{N_E^{\text{dof}}} \text{dof}_r((I - \Pi_0^k) \boldsymbol{\varphi}_i) \text{dof}_r((I - \Pi_0^k) \boldsymbol{\varphi}_j).$$

By writing the tensor $\mathbf{K} = \begin{bmatrix} K_{xx} & K_{xy} \\ K_{xy} & K_{yy} \end{bmatrix}$, we define the matrix $\mathbf{Z}^{Q,K_*} \in \mathbb{R}^{N^Q \times N^Q}$ as the diagonal matrix whose non zeros entries are

$$\mathbf{Z}_{ii}^{Q,K_*} = K_* (\mathbf{x}_i^Q) \omega_i^Q, \quad \forall i = 1, \dots, N^Q \quad \forall * \in \{xx, xy, yy\}.$$

Then, the consistency matrix can be computed as

$$\mathbf{K}_C^a = (\Pi_k^0)^T \mathbf{G}^{k,K} \Pi_k^0,$$

where the matrix $\mathbf{G}^{k,K} \in \mathbb{R}^{2n_k \times 2n_k}$ reads

$$\mathbf{G}^{k,K} = (\mathbf{g}\mathbf{V}^k)^T \begin{bmatrix} \mathbf{Z}^{Q,K_{xx}} & \mathbf{Z}^{Q,K_{xy}} \\ \mathbf{Z}^{Q,K_{xy}} & \mathbf{Z}^{Q,K_{yy}} \end{bmatrix} \mathbf{g}\mathbf{V}^k,$$

that is, $\mathbf{G}^{k,K}$ is the mass matrix \mathbf{G}^k related to the basis polynomial basis $\mathcal{G}_k(E)$ of degree k weighted by the tensor \mathbf{K} .

Now, we define the matrix $\mathbf{D} \in \mathbb{R}^{N_E^{\text{dof}} \times 2n_k}$, whose entries are

$$\mathbf{D}_{il} = \text{dof}_i(\mathbf{g}_l^k), \quad \forall i = 1, \dots, N_E^{\text{dof}}, \forall l = 1, \dots, 2n_k.$$

Matrix \mathbf{D} can be split as

$$\mathbf{D} = \begin{bmatrix} \mathbf{D}^\partial \\ \mathbf{D}^\nabla \\ \mathbf{D}^\perp \end{bmatrix},$$

being matrices $\mathbf{D}^\nabla \in \mathbb{R}^{n_{k-1}^\nabla \times 2n_k}$ and $\mathbf{D}^\perp \in \mathbb{R}^{n_k^\perp \times 2n_k}$ given by

$$\mathbf{D}^\nabla = \frac{1}{|E|} \mathbf{G}^k(1 : n_{k-1}^\nabla, :), \quad \mathbf{D}^\perp = \frac{1}{|E|} \mathbf{G}^k(n_k^\nabla + 1 : 2n_k, :).$$

Matrix $\mathbf{D}^\partial \in \mathbb{R}^{N^{E,e}(k+1) \times 2n_k}$, instead, $\forall i = 1, \dots, k+1, \forall e \in \mathcal{E}_{h,E}, \forall l = 1, \dots, 2n_k$, reads

$$\mathbf{D}[\partial]_{il} = \text{dof}_i^e(\mathbf{g}_l^k) = (\mathbf{g}_l^k \cdot \mathbf{n}_e)(\mathbf{x}_i^{e,Q}),$$

or, equivalently,

$$\begin{aligned} \mathbf{D}^\partial &= \begin{bmatrix} \mathbf{N}_x & \mathbf{N}_y \end{bmatrix} \mathbf{g} \mathbf{V}^{\partial,k} \\ &= \begin{bmatrix} \mathbf{N}_x & \mathbf{N}_y \end{bmatrix} \begin{bmatrix} \mathbb{P} \mathbf{V}^{\partial,k} & \mathbf{O} \in \mathbb{R}^{N^{E,e}(k+1) \times n_k} \\ \mathbf{O} \in \mathbb{R}^{N^{E,e}(k+1) \times n_k} & \mathbb{P} \mathbf{V}^{\partial,k} \end{bmatrix} \begin{bmatrix} (\mathbf{T}^{\nabla,k})^T & (\mathbf{T}^{\perp,k})^T \end{bmatrix}, \end{aligned}$$

where $\mathbf{N}_* \in \mathbb{R}^{N^{E,e}(k+1) \times N^{E,e}(k+1)}$, for all $* \in \{x, y\}$, is the diagonal matrix whose diagonal entries are the $*$ -component of the normal vectors to the edges of E , properly arranged.

Finally, the stability matrix can be computed as

$$\mathbf{K}_S^a = \bar{\mathbf{K}} |E| (\mathbf{I} - \mathbf{D} \mathbf{\Pi}_k^0)^T (\mathbf{I} - \mathbf{D} \mathbf{\Pi}_k^0).$$

5.4 | The advection term $-(\beta p_h, \mathbf{\Pi}_0^k \mathbf{v}_h)_{0,E}$

The matrix $\mathbf{A}^\beta \in \mathbb{R}^{N_E^{\text{dof}} \times n_k}$ corresponding to the local advection term

$$-(\beta p_h, \mathbf{v}_h)_{0,E},$$

$\forall i = 1, \dots, N_E^{\text{dof}}, \forall \alpha = 1, \dots, n_k$, is defined as

$$\mathbf{A}_{i\alpha}^\beta = - \int_E \beta \cdot \mathbf{\Pi}_0^k \boldsymbol{\varphi}_i \mathbb{P}_\alpha^k = - \sum_{l=1}^{2n_k} (\mathbf{\Pi}_k^0)_{li} \int_E \beta \cdot \mathbf{g}_l^k \mathbb{P}_\alpha^k.$$

Given $\boldsymbol{\beta} = [\beta_x, \beta_y]^T$, in matrix form \mathbf{A}^β reads

$$\mathbf{A}^\beta = (\mathbf{\Pi}_k^0)^T (\mathbf{g} \mathbf{V}^k)^T \begin{bmatrix} \mathbf{Z}^{Q,\beta_x} & \mathbf{O} \\ \mathbf{O} & \mathbf{Z}^{Q,\beta_y} \end{bmatrix} \begin{bmatrix} \mathbb{P} \mathbf{V}^k \\ \mathbb{P} \mathbf{V}^k \end{bmatrix},$$

where $\mathbf{Z}^{Q,\beta_*} \in \mathbb{R}^{N^Q \times N^Q} \forall * \in \{x, y\}$ is the diagonal matrix whose non zero entries are defined as

$$(\mathbf{Z}^{Q,\beta_*})_{ii} = \beta_*(\mathbf{x}_i^Q) \omega_i^Q, \quad \forall i = 1, \dots, N^Q.$$

5.5 | The reaction term $(\gamma p_h, q_h)_{0,E}$

We define the local reaction matrix $\mathbb{P} \mathbf{H}^{k,\gamma} \in \mathbb{R}^{n_k \times n_k}$ as the matrix that collects the terms

$$\mathbb{P} \mathbf{H}_{\alpha\beta}^{k,\gamma} = \int_E \gamma \mathbb{P}_\alpha^k \mathbb{P}_\beta^k \quad \forall \alpha, \beta = 1, \dots, n_k.$$

We observe that this matrix is the mass matrix $\mathbb{P} \mathbf{H}^k$ related to the polynomial basis $\{\mathbb{P}_\alpha^k\}_{\alpha=1}^{n_k}$ of degree k weighted by the reaction coefficient γ . Thus, in matrix form, we have

$$\mathbb{P} \mathbf{H}^{k,\gamma} = (\mathbb{P} \mathbf{V}^k)^T \mathbf{Z}^{Q,\gamma} \mathbb{P} \mathbf{V}^k,$$

where

$$\mathbf{Z}_{ij}^{Q,\gamma} = \gamma(\mathbf{x}_i^Q) \omega_i^Q \delta_{ij}, \quad \forall i, j = 1, \dots, N^Q.$$

In conclusion, the local system matrix $\mathbf{K}^E \in \mathbb{R}^{N_E^{\text{dof}} \times N_E^{\text{dof}}}$ for the discrete problem (27) is given by

$$\mathbf{K}^E = \begin{bmatrix} \mathbf{K}_C^a + \mathbf{K}_S^a & -\mathbf{W}^T + \mathbf{A}^\beta \\ \mathbf{W} & \mathbb{P} \mathbf{H}^{k,\gamma} \end{bmatrix}.$$

6 | NUMERICAL EXPERIMENTS

In this section, we propose some numerical experiments to show the performance of the aforementioned approaches in terms of the following error norms:

$$p_{\text{err}} = \left(\sum_{E \in \mathcal{T}_h} \|p - p_h\|_{0,E}^2 \right)^{\frac{1}{2}}, \quad (59)$$

$$\mathbf{u}_{\text{err}} = \left(\sum_{E \in \mathcal{T}_h} \|\mathbf{u} - \Pi_0^k \mathbf{u}_h\|_{0,E}^2 \right)^{\frac{1}{2}}, \quad (60)$$

$$p_{I,\text{err}} = \left(\sum_{E \in \mathcal{T}_h} \|p_I - p_h\|_{0,E}^2 \right)^{\frac{1}{2}}, \quad (61)$$

where the interpolant $p_I \in Q_{h,k}$ of p is computed locally as

$$p_I = \sum_{\alpha=1}^{n_k} c_\alpha \mathbb{P}_\alpha^k.$$

The vector $\mathbf{c} \in \mathbb{R}^{n_k}$ of coefficients is the solution of the following linear least squares problem

$$\min \|\mathbb{P}^k \mathbf{c} - \mathbf{y}\|,$$

where $\mathbf{y} \in \mathbb{R}^{N^Q}$ is the vector whose entries are the evaluation of p at the internal quadrature points $\{\mathbf{x}_i^Q\}_{i=1}^{N^Q}$.

We will also analyze the condition number of the main local matrices that allow assembling the local system matrix, namely \mathbf{G}^k , \mathbf{W} , \mathbf{B} , Π_k^0 and \mathbf{D} . The condition number of a matrix is computed as the ratio between its largest and its smallest singular values.

For the first two numerical examples, we consider problem (5) on the unit square domain $\Omega = (0, 1)^2$ with

$$\boldsymbol{\kappa} = \begin{bmatrix} y^2 + 1 & -xy \\ -xy & x^2 + 1 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} x \\ y \end{bmatrix}, \quad \gamma = x^2 + y^3.$$

The forcing term, the Dirichlet and the Neumann boundary conditions are set in such a way the exact solution is

$$p(x, y) = x^2 y + \sin(2\pi x) \sin(2\pi y) + 2, \quad (62)$$

and the set Γ_N coincides with the edge of the domain of interest on the x -axis.

In particular, in the first experiment (Test1), we validate the proposed methods showing the order of convergence of the three monomial, partial and full orthogonal approaches against rising polynomial degree k , for various mesh refinement levels, on meshes made of squared elements. Then, in the second test (Test2), we show the performance of the approaches in terms of matrix condition number and error convergence trends in presence of collapsing polygons in the meshes.

Finally, the last example (Test3) proposes the application of the method to flow simulations in Discrete Fracture Networks (DFNs). DFN simulations are, indeed, a typical example where highly distorted mesh elements might appear in the process to generate a conforming mesh [8]. A simple network is here considered, for which an analytic solution is known, in order to compare the convergence curves of the three different approaches.

6.1 | Test1: Convergence rates

In this first experiment, we validate the new approaches by showing that the computational orders of convergence match the theoretical ones for different values of the polynomial degree k on regular meshes. At this aim, convergence tests are performed on a sequence of squared meshes that decompose the domain of interest in 25, 100, 400 and 1600 identical squares, respectively.

Obtained convergence rates are shown in Table 1, while convergence curves of the errors are shown in Figure 1. In this figure, each row reports the graphs of the three error norms at varying k on a fixed mesh. The three rows correspond to the last three considered refinement levels, that is, the 100, 400 and 1600 element meshes, respectively. In the figure, the upper bound of y -axis is fixed to $1.0e + 1$.

The obtained results show that the rates of convergence related to all three approaches match the theoretical ones up to polynomial degree $k = 5$. For higher orders, errors in the monomial approach start to increase due to ill-conditioning, while both the partial and the full orthonormal approaches still provide the expected results in terms of errors, up to stagnation due to finite precision arithmetic. We highlight that, in the partial and full orthonormal approach, the local projection matrices and the global system matrix are well-conditioned, as we will see in the next tests.

6.2 | Test2: Collapsing polygons

Now, we analyze the condition number of local matrices in case of collapsing polygons. For this purpose, we use three meshes of rectangular elements, with elements having aspect ratios of 10, 50, and 100, respectively. We remark that the aspect ratio of an element is here defined as the ratio between the maximum and the minimum length of its edges. These meshes are built starting from a mesh made up of 100 identical squares, then subdivided into rectangles with length fixed to the value 0.1, equal to that of the original square, and a height computed according to the desired aspect ratio value. The mesh with aspect ratio of 10 originated from this process is shown in Figure 2.

Figures 3, 4 and 5 report the highest value among mesh elements of the condition number of elemental matrices \mathbf{G}^k , \mathbf{W} , \mathbf{B} , $\mathbf{\Pi}_0^k$ and \mathbf{D} , on the three considered meshes. In these figures, we set the upper bound of y -axis to $1.0e + 20$. As shown in these figures, mass matrix \mathbf{G}^k , resulting from the full orthonormal approach, is perfectly well-conditioned, since it corresponds to the identity matrix, up to machine precision. Moreover, an algebraic growth of the condition number of the mass matrix obtained with the partial orthonormal approach is observed, whereas conditioning grows exponentially for the monomial one. More generally, we can observe that the condition numbers of the local matrices vary in a very limited range in the partial and full orthonormal approaches, as k increases. Instead,

TABLE 1 Test1: Convergence rates on squared mesh.

	k	0	1	2	3	4	5	6	7	8
Monomial	p_{err}	1.1956	2.4760	3.7464	4.8565	5.9153	6.9456	5.8756	2.4067	1.5749
	u_{err}	0.9947	2.0937	3.1263	4.0877	5.1110	5.9812	3.2093	0.3267	-0.0167
	$p_{l,\text{err}}$	1.9417	3.0056	3.9580	4.9588	5.9529	6.9680	5.8761	2.4071	1.5749
Partial	p_{err}	1.1956	2.4760	3.7464	4.8565	5.9153	6.9456	7.6064	6.8582	5.1747
	u_{err}	0.9947	2.0937	3.1263	4.0877	5.1110	6.0812	5.4221	3.8693	2.1663
	$p_{l,\text{err}}$	1.9417	3.0056	3.9580	4.9588	5.9529	6.9680	7.6098	6.8590	5.1751
Ortho	p_{err}	1.1956	2.4760	3.7464	4.8565	5.9153	6.9456	7.9637	8.5476	7.2031
	u_{err}	0.9947	2.0937	3.1263	4.0877	5.1110	6.1058	7.0385	6.0692	4.3635
	$p_{l,\text{err}}$	1.9417	3.0056	3.9580	4.9588	5.9529	6.9680	7.9773	8.5484	7.2035

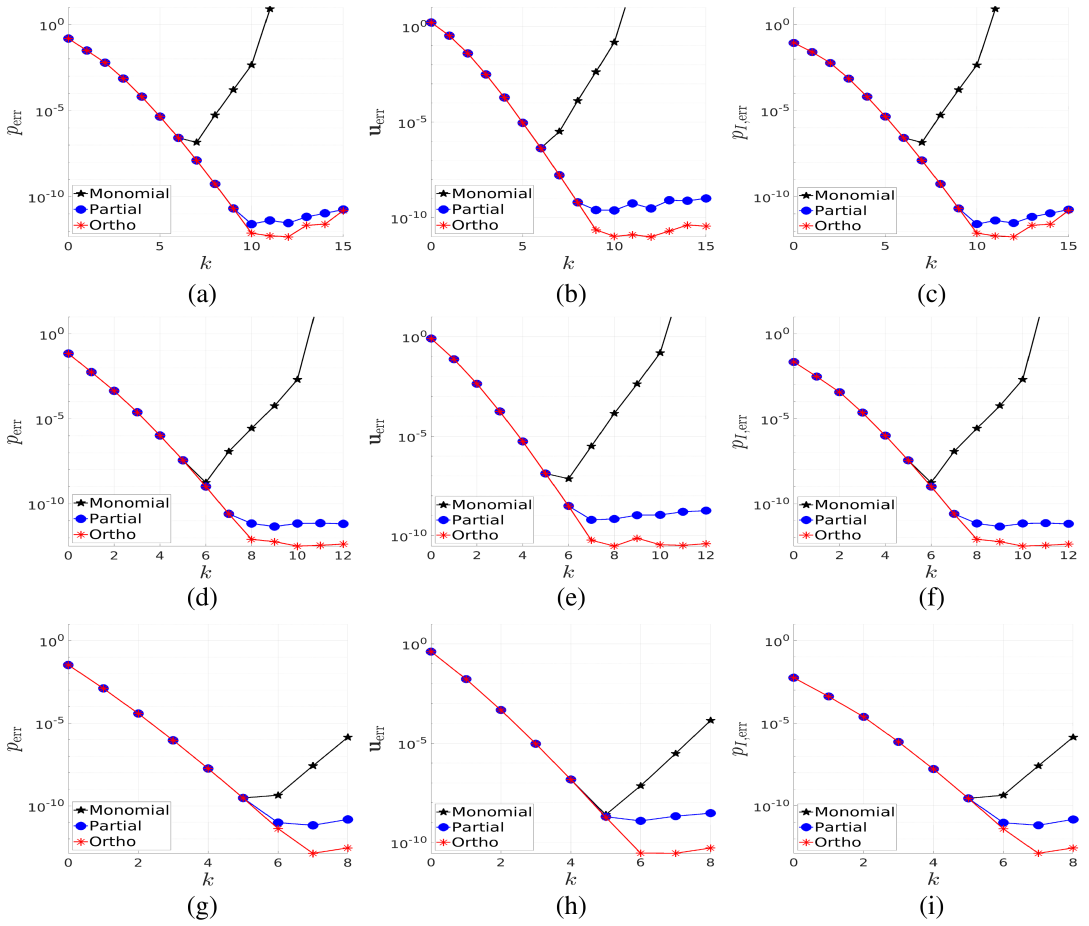


FIGURE 1 Test1: Behaviour of errors (59), (60) and (61) at varying k on square meshes. Pictures on each row refer to a different mesh refinement level: 100, 400 and 1600 element meshes from top to bottom.

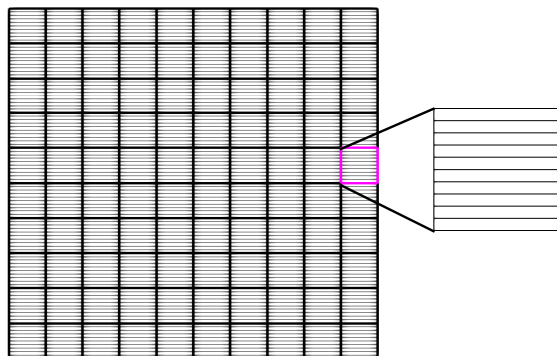


FIGURE 2 Test2: Mesh with aspect ratio 10.

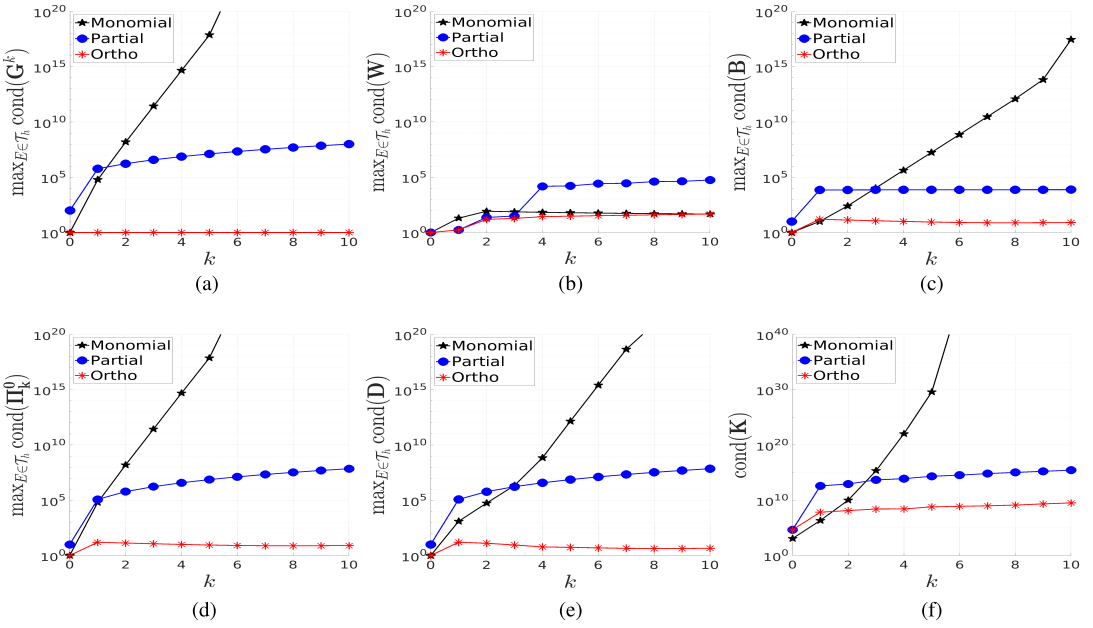


FIGURE 3 Test2: Figure 3a–e show the maximum condition number of local matrices among elements, at varying k . Figure 3f reports the behaviour of the condition number of the global system matrix \mathbf{K} , at varying k . Mesh with aspect ratio 10.

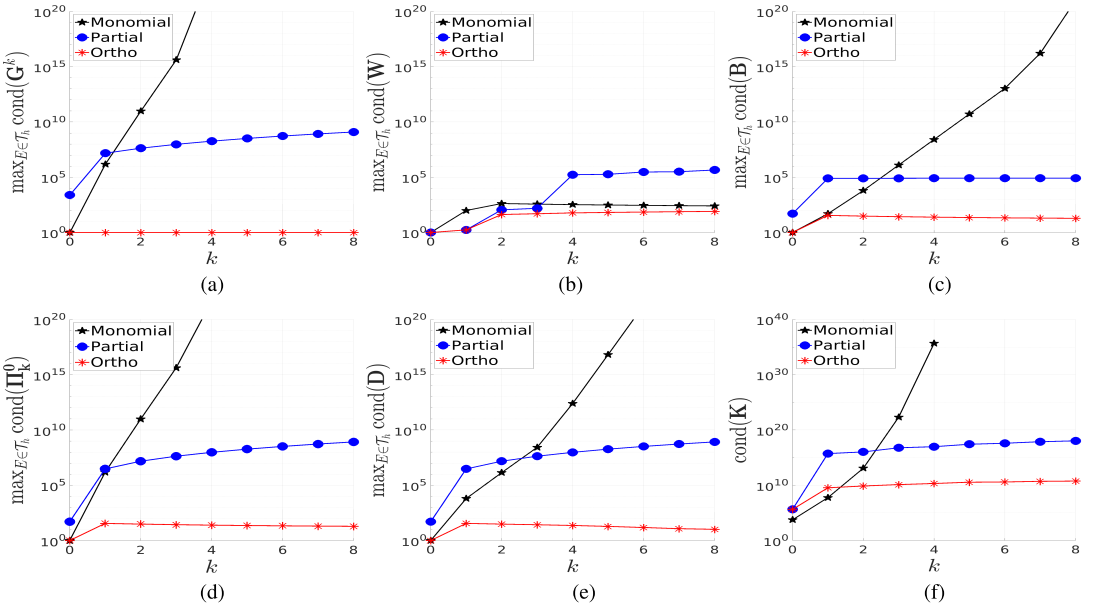


FIGURE 4 Test2: Figure 4a–e show the maximum condition number of local matrices among elements, at varying k . Figure 4f reports the behaviour of the condition number of the global system matrix \mathbf{K} , at varying k . Mesh with aspect ratio 50.

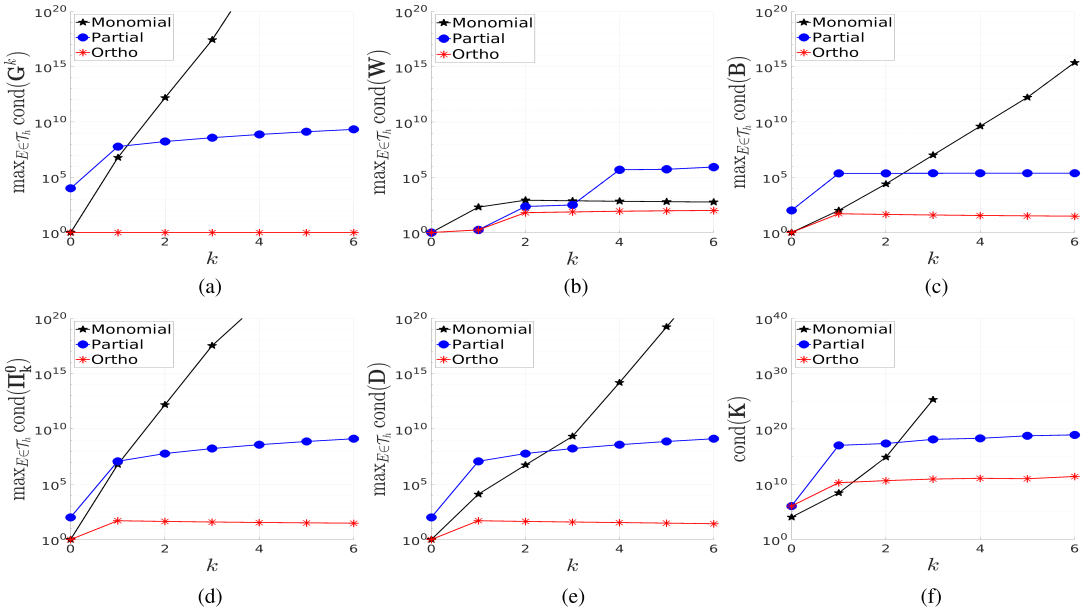


FIGURE 5 Test2: Figure 5a–e show the maximum condition number of local matrices among elements, at varying k . Figure 5f reports the behaviour of the condition number of the global system matrix \mathbf{K} , at varying k . Mesh with aspect ratio 100.

the condition number of local matrices grows exponentially in the monomial approach, with the only exception of matrices \mathbf{W} , which exhibit a good conditioning regardless of the case.

Conditioning of matrices obtained with the partial-orthonormal bases appear to be slightly more affected by increasing aspect ratio values than the corresponding matrices with the full-orthonormal approach, however, still showing much lower values than the matrices given by the monomial basis.

Figure 6 reports convergence curves of the error against growing polynomial accuracy k , for the three considered meshes: top row for the mesh with aspect ratio 10, middle row for aspect ratio 50 and bottom row for aspect ratio 100. Also in this case, the upper bound of y-axis is fixed to $1.0e + 1$. At low values of k , the curves corresponding to the monomial approach are well overlapped with the curves of partial and full-orthonormal approaches. Furthermore, the maximum value of polynomial accuracy k for which the monomial approach provides errors in line with the other approaches reduces as the aspect ratio of mesh elements increases. Finally, error curves relating to the monomial approach are interrupted at values of $k \leq 6$ due to failure of linear algebra libraries (the *SparseLU* solver of Eigen) in computing a solution due to ill-conditioning of the global system matrix. Small differences are instead noticed between the curves obtained with the partial and full orthonormal approaches for all the considered values of k . We remark that these behaviours are coherent with the corresponding trends of the condition number of the global system matrices, which are shown in Figures 3f, 4f and 5f for the three considered meshes. In these figures, the upper bound of y-axis is set to $1.0e + 40$. In the monomial approach, the condition number of the global system matrix \mathbf{K} grows exponentially as the polynomial degree k increases, while the growth in the partial and full orthonormal approaches is linear. The partial-orthonormal approach however provides condition numbers significantly higher than the full approach. In Figures 4f and 5f, the curves reporting the behaviour of the global condition numbers for the monomial approach are interrupted due to the failure of both the MATLAB R2023b routine *condst* and the C++ SuiteSparse v7.7.0 *klu_condst*. This behaviour confirms that using orthonormal

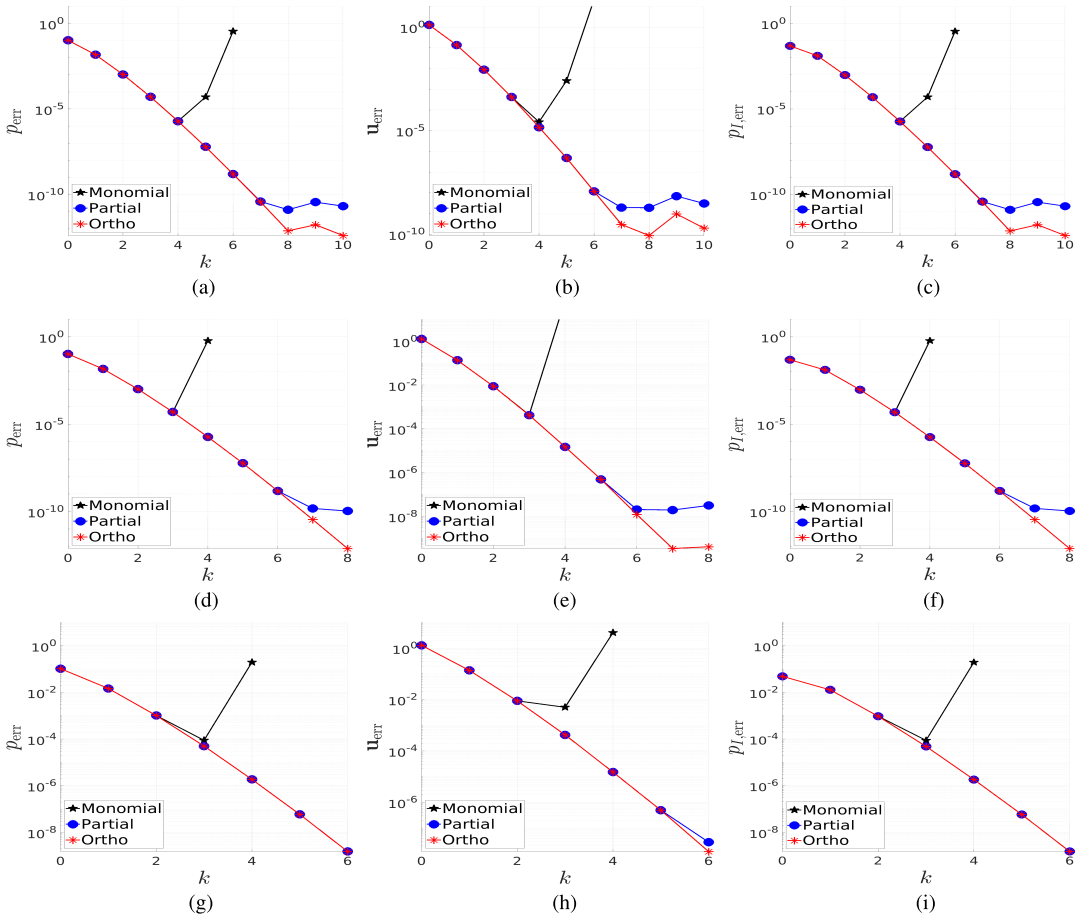


FIGURE 6 Test2: Behaviour of errors (59), (60) and (61) at varying k on rectangular meshes. Each row represents a different mesh: 10, 50, 100 from top to bottom.

polynomial basis to define the internal degrees of freedom allows to improve the conditioning of the global system matrix and to obtain more accurate solutions, as noted also in [15].

6.3 | Test3: Simulations in discrete fracture networks

In this last example, the application of the proposed orthonormal bases to Discrete Fracture Network (DFN) problems is presented. Discrete Fracture Networks are obtained as the union of planar polygonal domains with arbitrary orientations in the 3D space and are used to model the fractures in a porous medium [1]. As fracture thickness is typically orders of magnitude smaller than the other dimensions, fractures are geometrically reduced to 2D domains, and suitable equations, averaged across fracture thickness, are derived to describe the phenomena occurring in such domains [14]. Interface equations are then added at fracture intersections. A major complexity in DFN simulations consists in the generation of a conforming mesh, for realistic configurations characterized by intricate networks with a large number of fractures and fracture intersections. A possible strategy for DFN meshing is proposed, for example, in [8], based on the use of mixed virtual elements: first a triangular mesh is constructed on each fracture domain, independently of the intersections with the other fractures; then, the elements of these meshes are cut according to the interfaces, and hanging nodes are added where needed, to obtain

a fully conforming mesh of the whole domain (see [8] for more details). Highly elongated elements are likely to be generated in this process, such that the standard choice of VEM basis function might yield badly conditioned problems for high order approximations [12].

An advection-diffusion-reaction problem on a network made up by the union of three fractures F_i with three intersections Γ_i , is here considered, namely

$$\begin{aligned}
 F_1 &= \{(x, y, z) \in \mathbb{R}^3 : -1 \leq x \leq 1, -1 \leq y \leq 1, z = 0\}, \\
 F_2 &= \{(x, y, z) \in \mathbb{R}^3 : -1 \leq z \leq 1, -1 \leq x \leq 1, y = 0\}, \\
 F_3 &= \{(x, y, z) \in \mathbb{R}^3 : -1 \leq y \leq 1, -1 \leq z \leq 1, x = 0\}, \\
 \Gamma_1 &= \{(x, y, z) \in \mathbb{R}^3 : -1 \leq x \leq 1, y = 0, z = 0\}, \\
 \Gamma_2 &= \{(x, y, z) \in \mathbb{R}^3 : -1 \leq y \leq 1, x = 0, z = 0\}, \\
 \Gamma_3 &= \{(x, y, z) \in \mathbb{R}^3 : -1 \leq z \leq 1, x = 0, y = 0\}.
 \end{aligned}$$

On each fracture, we choose

$$\kappa_i(\hat{x}, \hat{y}) = \begin{bmatrix} 1 + \hat{y}^2 & -\frac{\hat{x}\hat{y}}{2} \\ -\frac{\hat{x}\hat{y}}{2} & 1 + \hat{x}^2 \end{bmatrix}, \quad \mathbf{b}_i(\hat{x}, \hat{y}) = \begin{bmatrix} \hat{x} - \hat{y} \\ \hat{y} - 1 \end{bmatrix}, \quad \gamma_i(\hat{x}, \hat{y}) = \hat{x}^3 + \hat{y}.$$

where (\hat{x}, \hat{y}) is a proper fracture-local reference system, and a forcing term and Neumann boundary conditions are defined in such a way the exact solution on each fracture is:

$$\begin{aligned}
 h_1(x, y) &= -|x|(1+x)(1-x)y(1+y)(1-y), \\
 h_2(z, x) &= -z(1+z)(1-z)x(1+x)(1-x), \\
 h_3(y, z) &= y(1+y)(1-y)|z|(1+z)(1-z).
 \end{aligned}$$

The same problem is considered in [8] up to order 5. The network and the exact solution are shown in Figure 7.

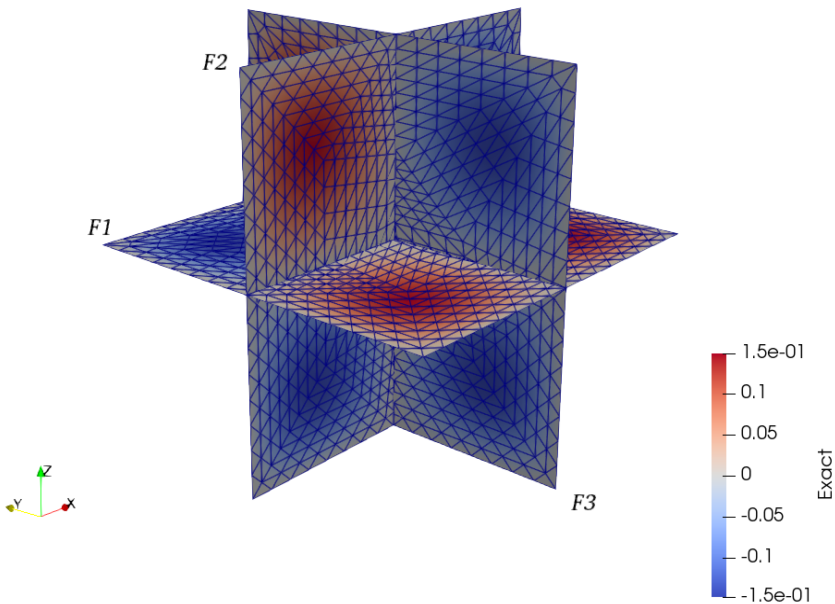


FIGURE 7 Test3: Exact solution DFN benchmark problem.

TABLE 2 Test3: Number of cells (Num Cells), minimum aspect ratio (Min AR) and maximum aspect ratio (Max AR) of mesh cells for the four refinement levels (R0 to R3).

	Num cells	Min AR	Max AR
R0	75	1.41	38.43
R1	246	1.19	62.11
R2	882	1.19	152.11
R3	3294	1.14	215.13

TABLE 3 Test3: Convergence rates on conforming mesh.

	k	0	1	2	3	4	5	6	7
Monomial	p_{err}	1.1984	2.3292	3.5494	4.6377	5.6597	2.8423	1.8227	-4.7912
	u_{err}	1.0354	1.8430	2.8787	3.9900	4.9328	1.1135	0.6710	-6.5224
	$p_{l,err}$	1.7581	2.6572	3.8542	4.8566	5.8788	2.8333	1.8227	-4.7912
Partial	p_{err}	1.1984	2.3292	3.5494	4.6377	5.6633	6.8989	5.7860	-2.8821
	u_{err}	1.0354	1.8430	2.8787	3.9900	5.0341	4.7002	2.8836	-4.0248
	$p_{l,err}$	1.7581	2.6572	3.8542	4.8566	5.8885	6.9844	5.7860	-2.8821
Ortho	p_{err}	1.1984	2.3292	3.5494	4.6377	5.6633	6.9237	8.0005	-0.3784
	u_{err}	1.0354	1.8430	2.8787	3.9900	5.0345	6.0422	6.2590	-2.7340
	$p_{l,err}$	1.7581	2.6572	3.8542	4.8566	5.8885	7.0225	8.0005	-0.3784

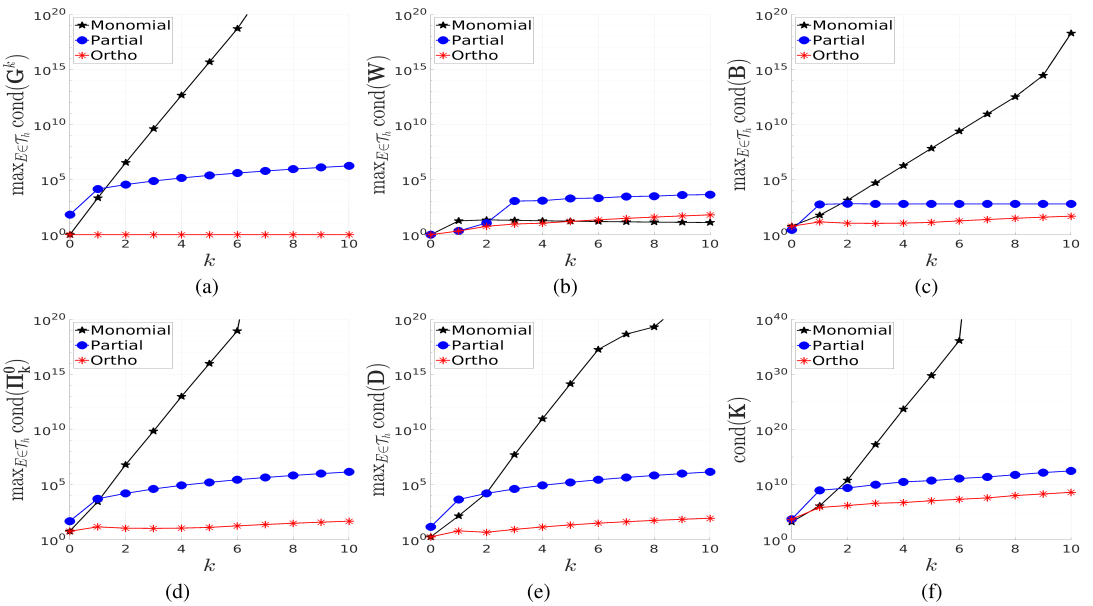


FIGURE 8 Test3: Figure 8a–e show the maximum condition number of local matrices among elements, at varying k . Figure 8f reports the behaviour of the condition number of the global system matrix \mathbf{K} , at varying k . Coarsest mesh.

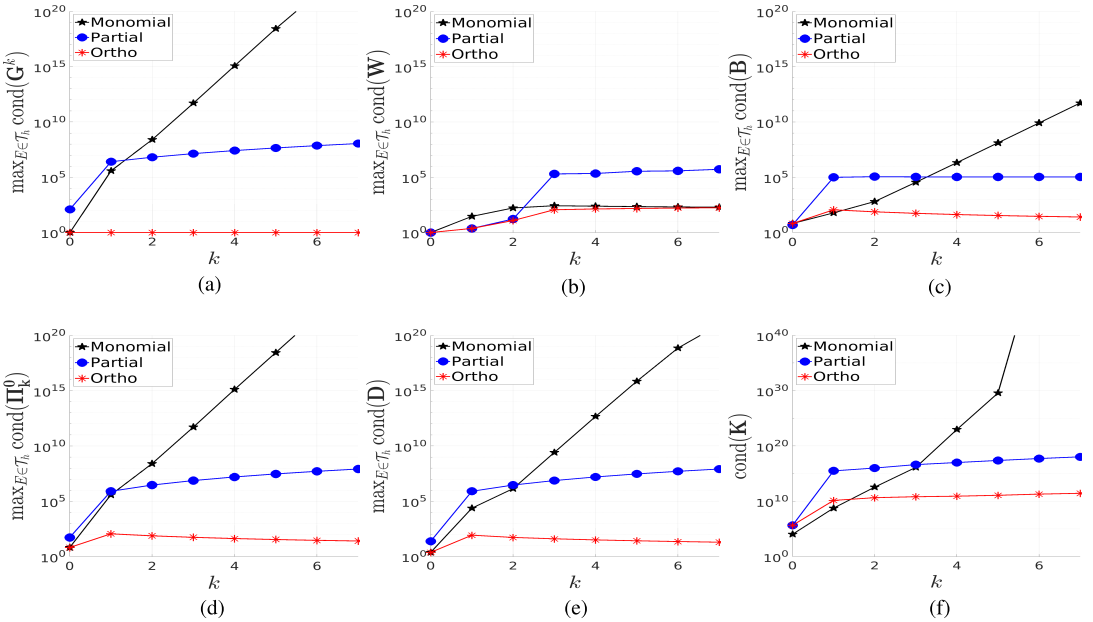


FIGURE 9 Test3: Figure 9a–c show the maximum condition number of local matrices among elements, at varying k . Figure 9f reports the behaviour of the condition number of the global system matrix K , at varying k . Finest mesh.

In this numerical test, we use four refinements of an initially triangular mesh, modified, as mentioned above, in such a way that the final polygonal meshes are conforming at the traces. Table 2 reports the number, the minimum and maximum aspect ratio of mesh elements for the four refinements (R0 to R3).

Table 3 shows the computed convergence rates of errors (59) and (61), for all the three tested approaches. In each sub-domain, the velocity field is a vector of polynomials of degree 7, such that, for $k \geq 7$ only errors related to floating point arithmetic computations are to be expected.

Figures 8 and 9 report, as previously, the maximum condition number across mesh elements of the computed local matrices and the trend of the condition number of the global system matrix as k varies, on the coarsest and finest considered meshes. We can observe that these data show the same behaviour as in the previous tests. Figure 10 shows error convergence curves against polynomial accuracy k , for the four considered mesh refinement levels. Pictures on each row correspond to the same mesh. It can be seen that the curves obtained with the three approaches are almost indistinguishable for k up to 4 – 5. For higher values errors given by the monomial approach start growing rapidly due to ill-conditioning, whereas errors given by partial and full orthonormal approaches still decrease up to stagnation due to finite precision arithmetic.

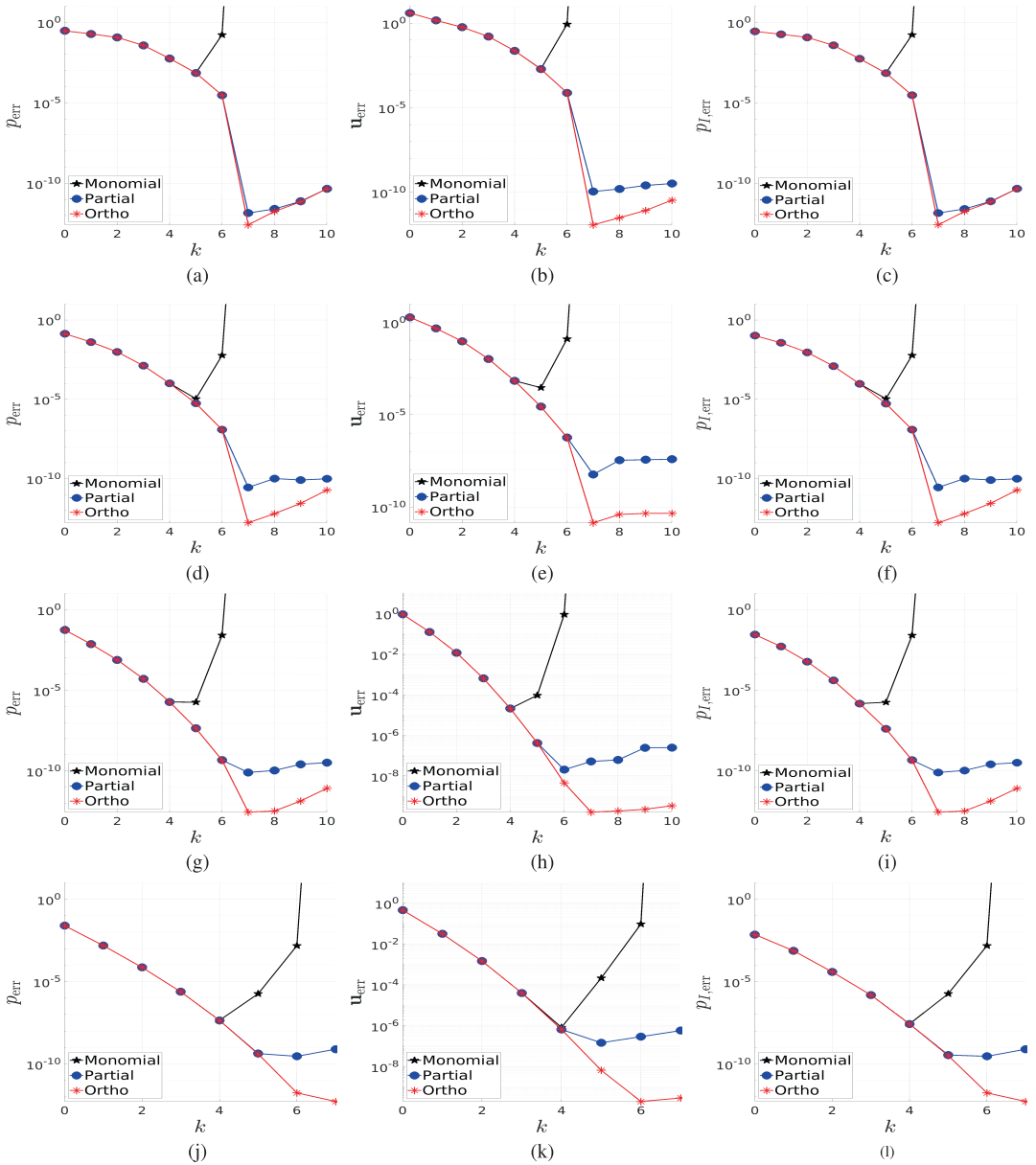


FIGURE 10 Test3: Behaviour of errors (59), (60) and (61) at varying k on conforming meshes. Each row represents a different refinement, from the coarsest mesh on top to the finest mesh at the bottom.

7 | CONCLUSIONS

In this paper, we presented a possible solution to cure the ill-conditioning of system matrix in the mixed formulation of the Virtual Element Method.

Since in the mixed formulation we need to introduce a discrete local space for both the pressure and the velocity variable, we have first introduced an orthonormal scalar-polynomial basis in the pressure space and then we have also orthonormalized the vector-polynomial basis used in the definition of the degrees of freedom related to the velocity variable.

Numerical experiments suggest that the introduction of orthonormal polynomial basis in both spaces allows to improve stability of mixed Virtual Elements for high order applications on distorted elements.

The additional computational cost for the Ortho and the Partial approaches with respect to the standard, monomial approach stems from the application of the Gram-Schmidt algorithm. This needs to be performed twice for the Partial approach and three times for the full-orthonormal approach. However, its cost is only associated with local quantities, such as the local polynomial degree k and the cardinality of the employed quadrature formula. Additionally, in the full-orthonormal approach, this cost is at least partially mitigated by the elimination of the two linear system resolutions, namely (58) and (55). It is also to remark that, overall, the leading cost is the one related to the resolution of the global system, that does not increase after the use of the proposed polynomial bases, and may even decrease due to improvements in the condition number.

It is worth to mention that the methods here suggested to build orthonormal polynomial bases improve the conditioning of the Vandermonde matrix defined with respect to the quadrature formulas in the interior of each element E . However, in general, this does not guarantee an improvement in the conditioning of the Vandermonde matrix defined with respect to quadrature formulas on the boundary of the elements [10]. Nonetheless, this appears to be sufficient to recover optimal convergence trends in the considered cases.

ACKNOWLEDGMENTS

The author S.B. kindly acknowledges partial financial support provided by PRIN project “Advanced polyhedral discretisations of heterogeneous PDEs for multiphysics problems” (No. 20204LN5N5_003) and by PNRR M4C2 project of CN00000013 National Centre for HPC, Big Data and Quantum Computing (HPC) (CUP: E13C22000990001). The author S.S. kindly acknowledges partial financial support provided by INdAM-GNCS through project “Sviluppo ed analisi di Metodi agli Elementi Virtuali per processi accoppiati su geometrie complesse” and that this publication is part of the project NODES which has received funding from the MUR-M4C2 1.5 of PNRR with grant agreement no. ECS00000036. The author G.T. kindly acknowledges financial support provided by the MIUR programme “Programma Operativo Nazionale Ricerca e Innovazione 2014–2020” (CUP: E11B21006490005). Computational resources are partially supported by SmartData@polito. Open access publishing facilitated by Politecnico di Torino, as part of the Wiley - CRUI-CARE agreement.

DATA AVAILABILITY STATEMENT

The data that support the findings of this study are available from the corresponding author upon reasonable request.

ORCID

Gioana Teora  <https://orcid.org/0000-0002-8540-3639>

REFERENCES

- [1] P. M. Adler, J. F. Thovert, and V. V. Mourzenko, *Fractured porous media*, Oxford University Press, Oxford, UK, 2012.
- [2] F. Bassi, L. Botti, A. Colombo, D. Di Pietro, and P. Tesini, *On the flexibility of agglomeration based physical space discontinuous Galerkin discretizations*, *J. Comput. Phys.* 231 (2012), no. 1, 45–65. <https://doi.org/10.1016/j.jcp.2011.08.018>.

- [3] L. Beirão da Veiga, F. Brezzi, A. Cangiani, G. Manzini, L. D. Marini, and A. Russo, *Basic principles of virtual element methods*, Math Models Methods Appl Sci 23 (2013), no. 1, 199–214. <https://doi.org/10.1142/S0218202512500492>.
- [4] L. Beirão da Veiga, F. Brezzi, L. D. Marini, and A. Russo, *The Hitchhiker's guide to the virtual element method*, Math Models Methods Appl Sci 24 (2014), no. 8, 1541–1573. <https://doi.org/10.1142/S021820251440003X>.
- [5] L. Beirão da Veiga, F. Brezzi, L. D. Marini, and A. Russo, H(div) and H(curl)-conforming VEM. 2016a. <https://doi.org/10.1007/s00211-015-0746-1>.
- [6] L. Beirão da Veiga, F. Brezzi, L. D. Marini, and A. Russo, *Mixed virtual element methods for general second order elliptic problems on polygonal meshes*, Esaim: M2AN 50 (2016b), no. 3, 727–747. <https://doi.org/10.1051/m2an/2015067>.
- [7] L. Beirão da Veiga, F. Brezzi, L. D. Marini, and A. Russo, *Virtual element implementation for general elliptic equations*, Springer International Publishing, Cham, 2016c, 39–71. https://doi.org/10.1007/978-3-319-41640-3_2.
- [8] M. F. Benedetto, A. Borio, and S. Scialò, *Mixed virtual elements for discrete fracture network simulations*, Finite Elem. Anal. Des. 134 (2017), 55–67. <https://doi.org/10.1016/j.finel.2017.05.011>.
- [9] S. Berrone and A. Borio, *Orthogonal polynomials in badly shaped polygonal elements for the virtual element method*, Finite Elem. Anal. Des. 129 (2017), 14–31. <https://doi.org/10.1016/j.finel.2017.01.006>.
- [10] S. Berrone, S. Scialò, and G. Teora, *The mixed virtual element discretization for highly-anisotropic problems: The role of the boundary degrees of freedom*, Math Eng 5 (2023), no. 6, 1–32. <https://doi.org/10.3934/mine.2023099>.
- [11] F. Brezzi, R. S. Falk, and L. Donatella Marini, *Basic principles of mixed virtual element methods*. ESAIM, Math Model Numer Anal 48 (2014), no. 4, 1227–1240. <https://doi.org/10.1051/m2an/2013138>.
- [12] F. Dassi, A. Fumagalli, D. Losapio, S. Scialò, A. Scotti, and G. Vacca, *The Mixed Virtual Element Method for Grids with Curved Interfaces in Single-Phase Flow Problems*. 2021. <https://doi.org/10.2118/203998-MS>.
- [13] L. Giraud, J. Langou, and M. Rozloznik, *The loss of orthogonality in the gram-Schmidt orthogonalization process*, Comput Math Appl 50 (2005), no. 7, 1069–1075. <https://doi.org/10.1016/j.camwa.2005.08.009>.
- [14] V. Martin, J. Jaffré, and J. E. Roberts, *Modeling fractures and barriers as interfaces for flow in porous media*, SIAM J. Sci. Comput. 26 (2005), no. 5, 1667–1691. <https://doi.org/10.1137/S1064827503429363>.
- [15] L. Mascotto, *Ill-conditioning in the virtual element method: Stabilizations and bases*, Numer Methods Partial Differ Equ 34 (2018), no. 4, 1258–1281. <https://doi.org/10.1002/num.22257>.

How to cite this article: S. Berrone, S. Scialò, and G. Teora, *Orthogonal polynomial bases in the mixed virtual element method*, Numer. Methods Partial Differ. Eq. (2024), e23144. <https://doi.org/10.1002/num.23144>