# Fundamental Data Structures for Matrix-Free Finite Elements on Hybrid Tetrahedral Grids

Nils Kohl[*,§], Daniel Bauer[*,†], Fabian Böhm[*,†], and Ulrich Rüde[*,‡]

[*]Friedrich-Alexander-Universität Erlangen-Nürnberg (FAU), Erlangen, Germany
[†]Erlangen National High Performance Computing Center (NHR@FAU), Erlangen, Germany
[‡]Centre Européen de Recherche et de Formation Avancée en Calcul Scientifique (CERFACS), Toulouse, France
[§]Contact: Nils Kohl (nils.kohl@fau.de)

### Abstract

This paper presents efficient data structures for the implementation of matrix-free finite element methods on block-structured, hybrid tetrahedral grids. It provides a complete categorization of all geometric sub-objects that emerge from the regular refinement of the unstructured, tetrahedral coarse grid and describes efficient iteration patterns and analytical linearization functions for the mapping of coefficients to memory addresses. This foundation enables the implementation of fast, extreme-scalable, matrix-free, iterative solvers, and in particular geometric multigrid methods by design. Their application to the variable-coefficient Stokes system subject to an enriched Galerkin discretization and to the curl-curl problem discretized with Nédélec edge elements showcases the flexibility of the implementation. Eventually, the solution of a curl-curl problem with $1.6 \cdot 10^{11}$ (more than one hundred billion) unknowns on more than $32\,000$ processes with a matrix-free full multigrid solver demonstrates its extreme-scalability.

## 1 Introduction

The numerical approximation of partial differential equations (PDEs) and their solutions at the extreme scale is a challenging task that requires the design and implementation of efficient and scalable parallel algorithms and data structures. Linear systems that arise from discretizations with very fine spatial and temporal resolution can be subject to billions ($10^9$) and even trillions ($10^{12}$) of unknowns [8]. In such cases, the solution vector alone requires the allocation of multiple terabytes of main memory. It therefore renders the explicit storage of the system matrix impossible since it typically requires of the order of a hundred times more memory [27, 40].

*Matrix-free* methods are crucial to solving such systems [13, 57, 44, 29, 28, 45, 5, 46, 40]. Since most iterative solvers only require the result of matrix-vector operations but no explicit access to the matrix entries, the assembly of the matrix and the evaluation of matrix-vector operations are fused and performed on-the-fly. On top of the reduced need for storage space, such approaches reduce the pressure on the memory bandwidth while the arithmetic intensity of the relevant compute kernels increases, favoring the machine balance of current cache-based architectures [60, 31].

This paper introduces fundamental data structures for the implementation of massively parallel, matrix-free finite element methods on block-structured triangular and tetrahedral grids. The concepts herein are implemented in the Hybrid Tetrahedral Grids (HyTeG) finite element framework, building upon what has been described in [41].[1]

HyTeG is a re-implementation of the Hierarchical Hybrid Grids (HHG) prototype software [16, 15], intending to cast the general concepts into a sustainable, extensible software framework. With extreme-scale applications in mind, the fundamental idea of HHG and HyTeG is to exploit the block-structured grids to implement fast, matrix-free linear solvers. Block-structured grids enable the construction of *implicit* index mappings. This is a critical advantage since it avoids the need for any bookkeeping data structures or indirect array accesses. Due to the regular refinement algorithm applied to the coarse grid elements, a grid hierarchy that enables straightforward implementation of geometric multigrid solvers is embedded into the framework by design. The parallel data structures combined with solvers of optimal asymptotic complexity using efficient matrix-free linear algebra enable the solution of systems with trillions ($\mathcal{O}(10^{12})$) of unknowns on hundreds of thousands of parallel processes [27, 40].
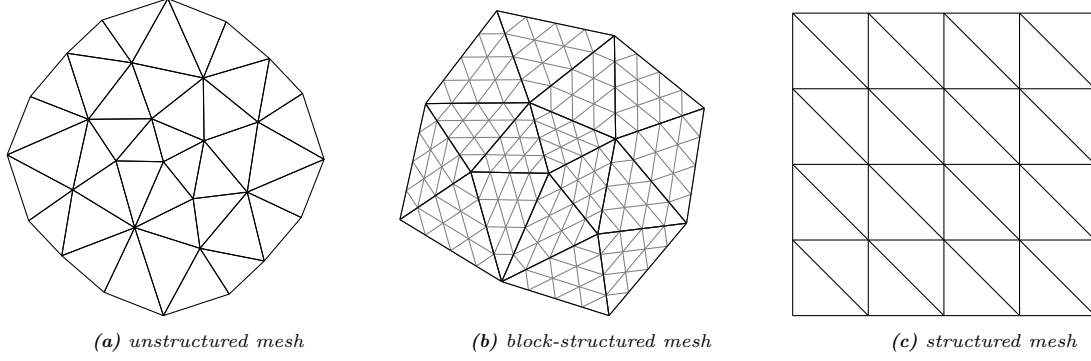
---

[1]https://i10git.cs.fau.de/hyteg/hyteg

**(a)** *unstructured mesh*    **(b)** *block-structured mesh*    **(c)** *structured mesh*

**Figure 1**  *Examples of unstructured, block-structured, and structured triangular meshes.*

The HHG prototype implementation has successfully demonstrated the performance and extreme scalability of this approach in a sequence of articles [16, 14, 13, 15, 28, 27, 11, 9, 12, 35, 29, 58, 36, 59]. The new HyTeG framework is a much more powerful reimplementation with more general data structures enabling much more flexible simulations on block-structured tetrahedral grids [41, 21, 8, 7, 43, 40, 39, 26].

This article provides a description of data structures that enable the implementation of *arbitrary* finite element discretizations for matrix-free computations on block-structured tetrahedral grids. Section 2 introduces the underlying block-structured mesh. Section 3 describes the regular refinement algorithm and categorizes the evolving mesh structure. Section 4 defines an indexing scheme that supports the implementation of efficient data structures for the implementation of matrix-free compute kernels. Section 5 discusses the design of the actual finite element spaces and matrix-free operators. Finally, Section 6 demonstrates the flexibility of the concept with two illustrative example applications.

## 2  Domain partitioning

The mesh that approximates the domain determines the requirements for the data structures of a simulation code. Unstructured meshes (example in Figure 1a) with varying element shapes throughout the grid are well-suited to approximate complex domain geometries. This comes at a cost since the data structures must incorporate the necessary bookkeeping to keep track of the connectivity and mesh shapes over the entire computational domain. Uniformly structured grids (example in Figure 1c) allow for efficient iteration patterns and data structures that can be optimized for low memory consumption since their geometry and connectivity are usually implicitly defined.

Uniform patterns of mesh elements provide crucial optimization potential for finite element codes, e.g., since the element integration is invariant under translation (given constant PDE coefficients) [15], and structured iteration patterns favor the implementation of efficient compute kernels on current cache-based architectures due to predictable memory accesses. Block-structured meshes (example in Figure 1b) are a tradeoff between adaptability and performance. The coarse elements provide a rough approximation of the domain. They are then uniformly refined to exhibit a local structure that is exploited for efficient compute kernels and data structures. HHG and HyTeG are based on this idea and employ block-structured triangular or tetrahedral elements. Tetrahedral meshes provide more flexibility than hexahedral meshes due to their superior geometric adaptability but require more complicated data structures.

This paper only deals with triangular and tetrahedral meshes that approximate the physical domain $\Omega$. The triangulation $\mathcal{T}_h(\Omega)$ of $\Omega \subset \mathbb{R}^d$ is defined as a partition of $\Omega$ into triangles (in two dimensions) or tetrahedrons (in three dimensions) $T_i$. In particular, $\mathcal{T}_h(\Omega) \coloneqq \{T_1, \ldots, T_{n_E}\}$, $T_i \subset \Omega$ open, such that

$$\overline{\Omega} = \bigcup_{T_i \in \mathcal{T}_h(\Omega)} \overline{T}_i, \tag{2.1}$$

$$T_i \cap T_j = \emptyset \text{ for } i \neq j, \text{ and} \tag{2.2}$$

$$\overline{T}_i \cap \overline{T}_j, \ i \neq j \text{ is a common} \begin{cases} \text{vertex or edge} & \text{if } d = 2, \\ \text{vertex, edge, or face} & \text{if } d = 3. \end{cases} \tag{2.3}$$

The mesh size $h$ is indicated by the maximum of the diameters $h_i$ of all elements.

**(a)** *unstructured mesh*   **(b)** *macro-primitives*   **(c)** *graph representation*
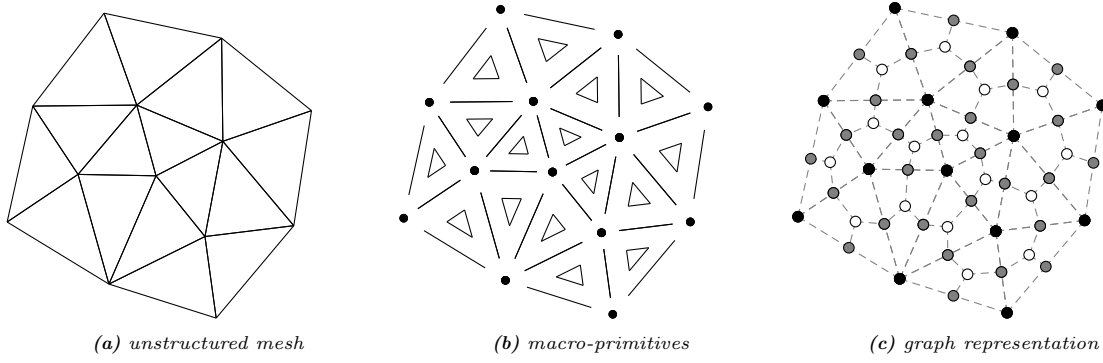
***Figure 2***  *Unstructured mesh (Figure 2a), the corresponding macro-primitives (Figure 2b), and a possible graph representation (Figure 2c). Each macro-primitive can be interpreted as a node of a graph. The connectivity of the graph defines the pairs of macro-primitives that can exchange information with each other. It is usually discretization-dependent and can be configured as required. Figure 2c illustrates a configuration with (undirected) (graph-)edges that connect macro-vertices (black nodes) to macro-edges (gray nodes), and macro-edges to macro-faces (white nodes).*

*Remark* 1. Domains with curved boundaries are handled via projection of elements using suitable mapping techniques, e.g., those described in [30]. In HyTeG, a surrogate approach is used to efficiently approximate arising integrals. Details are found in, e.g., [11, 10, 21]. See Section 6.2 for a demonstration.

In HyTeG the unstructured coarse mesh is represented by an undirected graph of so-called *macro-primitives*. See Figure 2 for a simple 2D example. For each cell, face, edge, and vertex of the coarse mesh, a macro-cell, macro-face, macro-edge, and macro-vertex data structure is allocated, respectively. Macro-primitives correspond to the nodes of a graph, with (graph-)edges between neighboring primitives. It may depend on the particular application which primitives are considered neighbors. For instance, each macro-cell usually references the neighboring macro-faces, macro-edges, and macro-vertices. However, for certain discretizations (e.g., discontinuous Galerkin (DG) schemes), references from macro-cells to neighboring macro-cells are additionally required. Details are found in [41].

On distributed parallel systems, each macro-primitive is assigned to one of the available processes. A load balancing algorithm determines the distribution via weights for each of the graph's nodes and edges. Each process only stores their assigned macro-primitives and metadata about those that are connected to the local primitives.

Spatially-dependent simulation data is allocated locally on each macro-primitive but can be migrated dynamically from one process to another one during run time. The migration and parallel re-partitioning are vital to load balancing, coarse grid agglomeration [19], and in-memory checkpointing [38].

# 3 Refinement

This section covers the refinement process of the block-structured grid and introduces an indexing scheme to organize the emerging geometric patterns. All considerations herein build on the refinement algorithm due to Bey [17].

*Remark* 2 (Macro- and micro-primitives). To avoid confusion, we refer to the cells, faces, edges, and vertices that are part of the unstructured coarse mesh as *macro*-cells, -faces, -edges, and -vertices and those created during the refinement of a macro-primitive as *micro*-cells, -faces, -edges, and -vertices.

*Remark* 3 (Two dimensions). All following structures and indexing schemes apply to two-dimensional settings by restriction to the $xy$-plane.

## 3.1 Regular refinement

The refinement process of an individual tetrahedron is described for a reference tetrahedron. For an illustration of the structure see Figure 3a. Any tetrahedron from the unstructured coarse grid can be mapped to the reference tetrahedron.

To create a mesh hierarchy, the volume-primitives (macro-faces in 2D, macro-cells in 3D) are regularly refined according to [17]. Given the vertices $T = [v_0, v_1, v_2, v_3]$ of a tetrahedron, and the midpoints $v_{ij}$ of
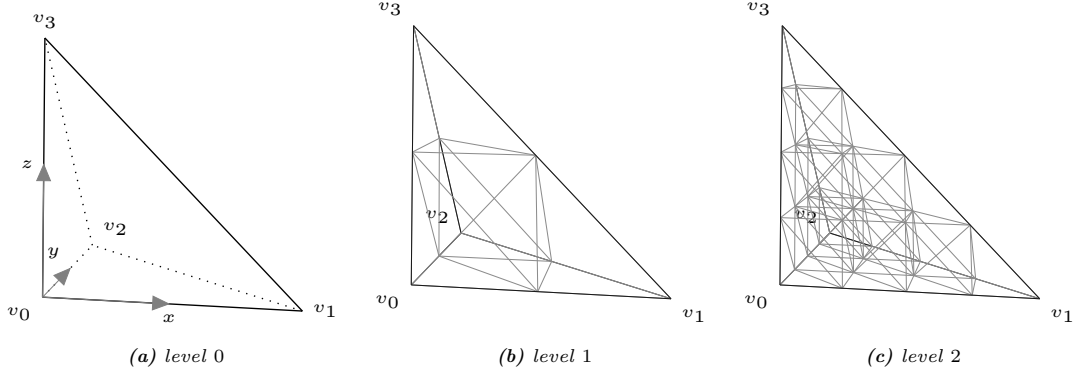
**Figure 3** *Uniform refinement of the reference tetrahedron according to [17] (i.e., (3.1)). Figure 3a also shows the reference coordinate system.*
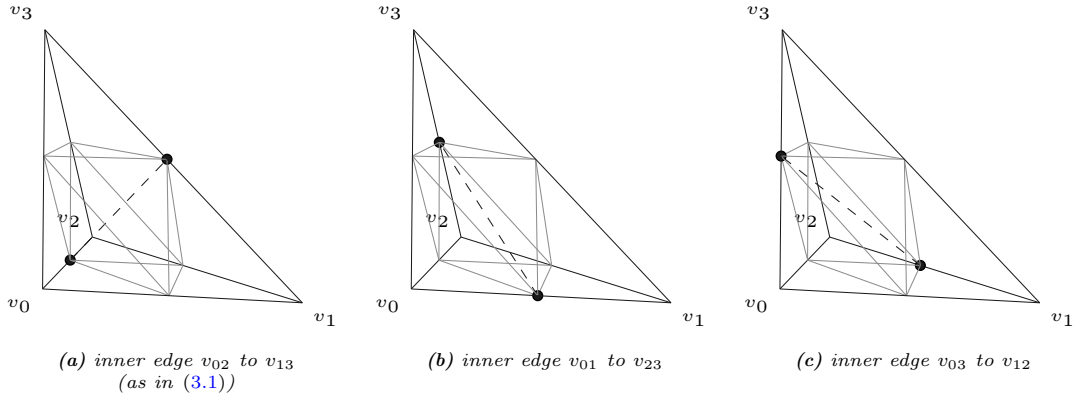


**(a)** *inner edge $v_{02}$ to $v_{13}$ (as in (3.1))*  **(b)** *inner edge $v_{01}$ to $v_{23}$*  **(c)** *inner edge $v_{03}$ to $v_{12}$*

**Figure 4** *Different ways of choosing the inner edge (dashed) during the refinement process of a tetrahedron.*

the edges connecting $v_i$ and $v_j$, $0 \leq i, j \leq 3$, $i \neq j$, $T$ is divided into the eight subtetrahedra $T_i$, $1 \leq i \leq 8$:

$$
\begin{aligned}
T_1 &:= [\ v_0,\ v_{01},\ v_{02},\ v_{03}\ ], \\
T_2 &:= [\ v_{01},\ v_1,\ v_{12},\ v_{13}\ ], \\
T_3 &:= [\ v_{02},\ v_{12},\ v_2,\ v_{23}\ ], \\
T_4 &:= [\ v_{03},\ v_{13},\ v_{23},\ v_3\ ], \\
T_5 &:= [\ v_{01},\ v_{02},\ v_{03},\ v_{13}\ ], \\
T_6 &:= [\ v_{01},\ v_{02},\ v_{12},\ v_{13}\ ], \\
T_7 &:= [\ v_{02},\ v_{03},\ v_{13},\ v_{23}\ ], \\
T_8 &:= [\ v_{02},\ v_{12},\ v_{13},\ v_{23}\ ].
\end{aligned}
\tag{3.1}
$$

This process is recursively applied to the subtetrahedra $T_1, \ldots, T_8$. The meshes of each refinement iteration are identified by levels. Level 0 refers to the original, coarse grid tetrahedron $T$ (or the entire macro-primitive coarse grid without refinement), level 1 to the mesh after one refinement iteration, etc. Figure 3 illustrates the refinement of a single tetrahedron. Note that this algorithm produces a globally conforming mesh, i.e., no hanging nodes are introduced at the intersection of two coarse grid tetrahedra.

The rule (3.1) produces different grids depending on the permutation of the vertices $v_0, \ldots, v_3$ of the original tetrahedron. Three different possible grids may occur. Starting from the reference tetrahedron, they are identified by the orientation of the inner edge, as illustrated in Figure 4. The permutation of the local vertex indices of a macro-element enables the optimization of the refinement process towards well-shaped elements and simplifies the implementation since the inner edge orientation can be fixed relative to the reference tetrahedron (in HyTeG, the inner edge always connects $v_{02}$ and $v_{13}$ as in (3.1)).

## 3.2  Micro-primitives

The recursive refinement of the reference tetrahedron, according to [17], introduces a fixed number of *subgroups* of micro-primitives. Regardless of the refinement level, each micro-primitive belongs to one
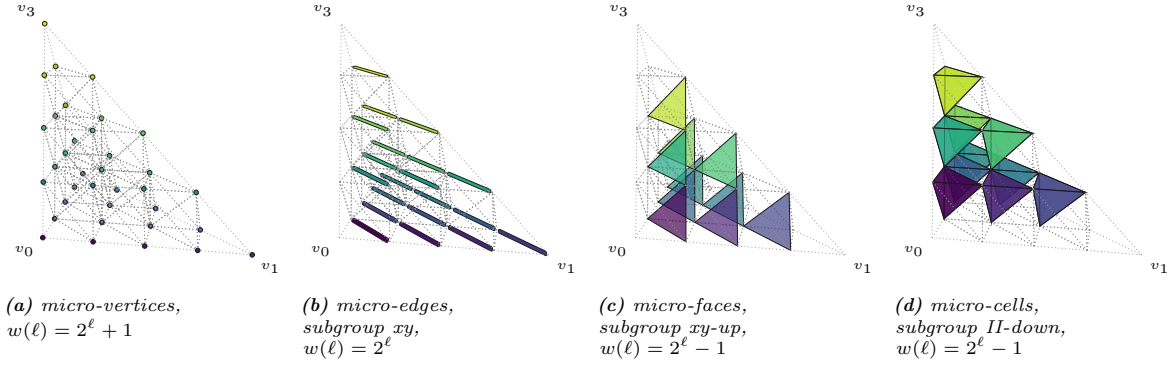
**(a)** *micro-vertices,*
$w(\ell) = 2^\ell + 1$

**(b)** *micro-edges,*
*subgroup xy,*
$w(\ell) = 2^\ell$

**(c)** *micro-faces,*
*subgroup xy-up,*
$w(\ell) = 2^\ell - 1$

**(d)** *micro-cells,*
*subgroup II-down,*
$w(\ell) = 2^\ell - 1$

***Figure 5*** *Illustration of one example subgroup for each micro-primitive type and the corresponding polytope side lengths* $w(\ell)$ *as functions of the refinement level* $\ell \geq 2$ *($\ell = 2$ in the figures). Illustrations and the corresponding polytope side lengths* $w(\ell)$ *for all subgroups are found in Appendix A, Figures 15 to 17 and Table 1.*

of the subgroups. Within each subgroup, micro-primitives are identical up to translation on the same refinement level and identical up to translation and scaling on different refinement levels.

This categorization yields

- 1 micro-vertex subgroup,

- 7 micro-edge subgroups,

- 12 micro-face subgroups,

- 6 micro-cell subgroups.

Figure 5 illustrates one subgroup of each micro-primitive type after two refinement iterations on the reference tetrahedron. Figures 15 to 17 (Appendix A) give a complete overview of all subgroups.

*Remark* 4 (Minumum refinement level). All considerations here are only valid on refinement levels $\ell \geq 2$. On level $\ell = 0$, for instance, there is only one micro-cell subgroup. Generally, on levels 0 and 1, the refined tetrahedron does not exhibit micro-primitives of all subgroups.

*Remark* 5 (Similar subgroups). The six subgroups of micro-cells are arranged into pairs, that are identical up to translation, reflection, and rotation. The naming scheme reflects this categorization (-up and -down types). Similar properties apply to the micro-face subgroups. See Figures 16 and 17.

The categorization allows for a unique mapping of each element of a subgroup to the set of indices

$$I_{\text{tet}}(w) = \{(i, j, k) : i + j + k < w, \ i, j, k \in \mathbb{Z}_{\geq 0}\}, \tag{3.2}$$

where $w \in \mathbb{Z}_{>0}$ is the side length of a corresponding tetrahedral polytope [48]. The side length $w$ is equal along all macro-edges, and depends on the refinement level $\ell$. That dependence may be different for each subgroup. The triangular polytope is defined by

$$I_{\text{tri}}(w) = \{(i, j) : i + j < w, \ i, j \in \mathbb{Z}_{\geq 0}\}. \tag{3.3}$$

The cardinalities of $I_{\text{tet}}(w)$ and $I_{\text{tri}}(w)$ are equal to the $w^{\text{th}}$ tetrahedral number

$$N_{\text{tet}}(w) = \sum_{k=1}^{w} \left( \sum_{i=1}^{k} i \right) = \frac{w(w+1)(w+2)}{6}, \tag{3.4}$$

and the $w^{\text{th}}$ triangular number

$$N_{\text{tri}}(w) = \sum_{k=1}^{w} k = \frac{w(w+1)}{2}. \tag{3.5}$$

The concrete values of $w$ of each subgroup illustrated in Figure 5 are given in the respective captions. A complete list of values for all subgroups is provided in Table 1.

The mapping of subgroups of micro-primitives to (3.2) facilitates the construction of loop nests and corresponding index linearization functions, which will be discussed next, in Section 4.

5

| micro-primitive | subgroup | $w(\ell)$ | illustration |
|---|---|---|---|
| vertex | - | $2^\ell + 1$ | Figure 15a |
| edge | x | $2^\ell$ | Figure 15b |
| | y | $2^\ell$ | Figure 15c |
| | z | $2^\ell$ | Figure 15d |
| | xy | $2^\ell$ | Figure 15e |
| | xz | $2^\ell$ | Figure 15f |
| | yz | $2^\ell$ | Figure 15g |
| | xyz | $2^\ell - 1$ | Figure 15h |
| face | z-up | $2^\ell$ | Figure 16a |
| | z-down | $2^\ell - 1$ | Figure 16b |
| | y-up | $2^\ell$ | Figure 16c |
| | y-down | $2^\ell - 1$ | Figure 16d |
| | x-up | $2^\ell$ | Figure 16e |
| | x-down | $2^\ell - 1$ | Figure 16f |
| | xyz-up | $2^\ell$ | Figure 16g |
| | xyz-down | $2^\ell - 2$ | Figure 16h |
| | xy-up | $2^\ell - 1$ | Figure 16i |
| | xy-down | $2^\ell - 1$ | Figure 16j |
| | yz-up | $2^\ell - 1$ | Figure 16k |
| | yz-down | $2^\ell - 1$ | Figure 16l |
| cell | I-up | $2^\ell$ | Figure 17a |
| | I-down | $2^\ell - 2$ | Figure 17b |
| | II-up | $2^\ell - 1$ | Figure 17c |
| | II-down | $2^\ell - 1$ | Figure 17d |
| | III-up | $2^\ell - 1$ | Figure 17e |
| | III-down | $2^\ell - 1$ | Figure 17f |

**Table 1** *Overview of all micro-primitive subgroups and the side lengths $w(\ell)$ of the corresponding tetrahedral polytopes, as functions of the refinement level $\ell \geq 2$.*

## 4 Indexing

In the context of finite element discretizations, degrees of freedom (DoFs) are usually associated with the underlying mesh's vertices, edges, faces, or volume elements. One or more scalars (e.g., for vector-valued discretizations or when a finite element discretization associates multiple DoFs with a certain micro-primitive) are then allocated for each micro-primitive of a specific type. For instance, a $\mathbb{P}_2$ finite element discretization typically associates one DoF per vertex and one per edge.

The coefficients are stored in arrays, and a linearization function has to be defined that uniquely maps each DoF to a corresponding memory address. Due to the association of DoFs and micro-primitives, this is accomplished via the construction of a unique mapping of each micro-primitive to an integer. Such a mapping can generally not be defined analytically on unstructured meshes, and therefore requires an additional data structure. The critical advantage of (block-) structured meshes is that they enable the construction of *implicit* index mappings and avoid the need for any bookkeeping data structures or indirect array accesses.

The regular structure of each micro-primitive subgroup enables the definition of linearization functions that only depend on the side length $w$ of the respective polytope. A bijective map $t_w : I_{\text{tet}}(w) \rightarrow \{0, \ldots, N_{\text{tet}}(w) - 1\}$ from the tetrahedral index set (3.2) to an array index can be defined as

$$t_w(i, j, k) \coloneqq N_{\text{tet}}(w) - N_{\text{tet}}(w - k) + N_{\text{tri}}(w - k) - N_{\text{tri}}(w - k - j) + i, \tag{4.1}$$

using the cardinalities (3.4) and (3.5) of the index sets (3.2) and (3.3). Figure 6 illustrates the map (4.1) for the z-down micro-face subgroup. Obviously, other mappings can be constructed.

If the number of DoFs $m \in \mathbb{Z}_{>0}$ per micro-primitive is larger than 1, (4.1) cannot directly be used as a linearization function. Two standard layouts to arrange the corresponding sequences of data sets in memory are referred to as array of structures (AoS) and structure or arrays (SoA). Here, a structure is a tuple of $m$ DoFs.

The memory layouts differ in how the structures interleave. The SoA layout enumerates the first element of each structure sequentially, followed by the second element of each structure, etc., i.e., the first $N_{\text{tet}}$ array elements correspond to the first DoF of the $N_{\text{tet}}$ micro-primitives. In contrast, the AoS layout
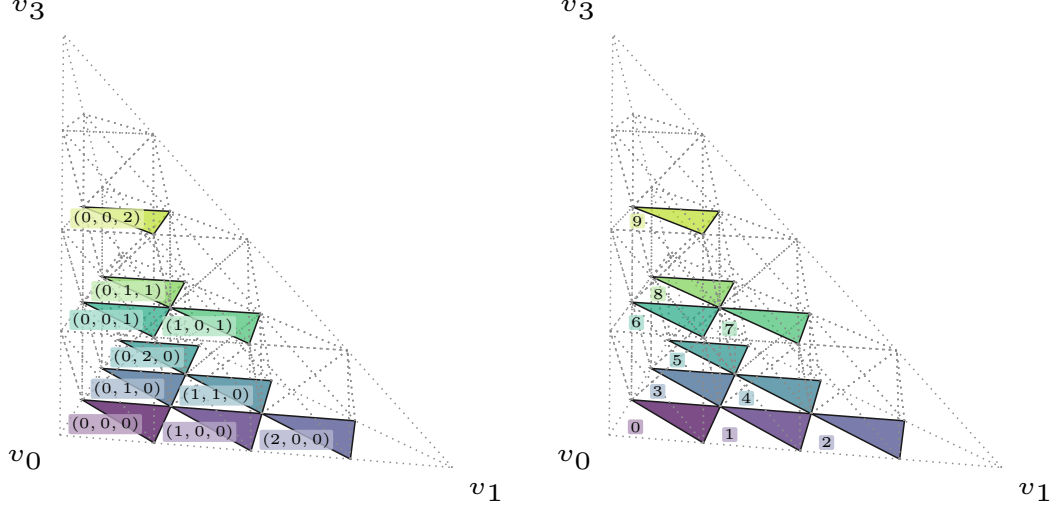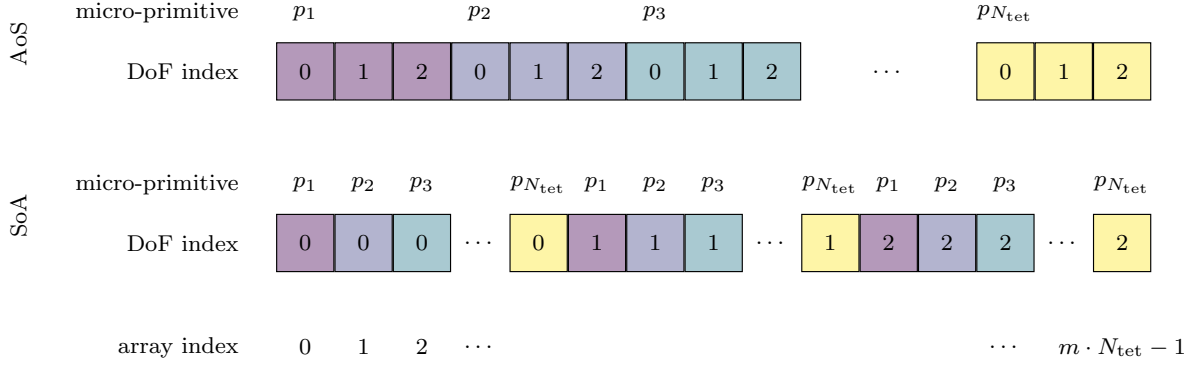
**Figure 6** *Mapping of polytope coordinates $(i, j, k)$ (left) to consecutive integers (right) via the linearization function $t_w(i, j, k)$ defined in (4.1) for micro-faces of the subgroup z-down on level $\ell = 2$ ($w = 3$).*



**Figure 7** *Illustration of the AoS (top) and SoA (bottom) memory layouts for $m = 3$ DoFs per micro-primitive of one subgroup. The micro-primitive indices are denoted by $p_\xi = (i_\xi, j_\xi, k_\xi)$. Array elements of the same color refer to DoFs on the same micro-primitive, i.e., to DoFs in the same structure.*

linearizes each entire structure after the other, i.e., the first $m$ array elements correspond to the $m$ DoFs of the first micro-primitive. Using (4.1) for the linearization of the micro-primitive indices, the AoS and SoA linearization functions are defined as

$$t_{w,m}^{\mathrm{AoS}}(i, j, k, d) = m \cdot t_w(i, j, k) + d, \tag{4.2}$$

and

$$t_{w,m}^{\mathrm{SoA}}(i, j, k, d) = d \cdot N_{\mathrm{tet}}(w) + t_w(i, j, k), \tag{4.3}$$

respectively, where $d \in \{0, \ldots, m - 1\}$ is the index of the DoF on the micro-primitive with index $(i, j, k)$. Both layouts are illustrated in Figure 7.

The concrete selection of the iteration pattern and linearization function are crucial for the performance of the algorithm. Both must be designed to preserve cache locality and to enable efficient vectorization. However, in general, the iteration pattern cannot be chosen arbitrarily since the access order modification may influence the underlying algorithm's properties. For instance, the convergence properties of a Gauss-Seidel iteration depend on the update pattern [56]. Code Listing 1 lists two example loop nests that result in entirely consecutive access patterns for the AoS and SoA linearization functions (4.2) and (4.3). The three loops with counters $i$, $j$, and $k$ correspond to an iteration along the Cartesian coordinates $x$, $y$, and $z$ (c.f., Figure 3a). The micro-primitives are traversed first in $x$-direction, then in $y$-direction, and finally in $z$-direction. The remaining loops in Code Listing 1 iterate over all $m$ allocated scalars on each micro-primitive.

```
1   // array of structures (AoS)
2   for ( int k = 0; k < w; k++ )
3     for ( int j = 0; j < w - k; j++ )
4       for ( int i = 0; i < w - k - j; i++ )
5         for ( int d = 0; d < m; d++ ) // innermost loop over structure
6           int lin_idx = t_aos( w, i, j, k, d );
7
8   // structure of arrays (SoA)
9   for ( int d = 0; d < m; d++ )        // outermost loop over structure
10    for ( int k = 0; k < w; k++ )
11      for ( int j = 0; j < w - k; j++ )
12        for ( int i = 0; i < w - k - j; i++ )
13          int lin_idx = t_soa( w, i, j, k, d );
```

**Code Listing 1**  *Loop nests for the iteration over the index set $I_{\text{tet}}(w)$ with $m$ DoFs per micro-primitive (see (3.2)). Each loop nest is structured to preserve consecutive array accesses for the AoS and SoA memory layouts respectively.*



**(a)** *micro-vertices, $m = 1$*    **(b)** *micro-edges, $m = 2$*    **(c)** *micro-faces, $m = 3$*    **(d)** *micro-volumes, $m = 4$*

**Figure 8**  *Visualization of the association of DoFs to micro-primitives. The integer $m$ denotes the number of DoFs that are associated with one micro-primitive, as discussed in Section 4. It is randomly chosen here for illustrative purposes. An overview of various types of finite element spaces and corresponding DoFs layouts can be found in [3].*

# 5   Matrix-free finite elements

The present section describes the construction of data structures for finite element discretizations based on the indexing schemes introduced in Sections 3 and 4 and discusses the implementation of matrix-free operators that act on the corresponding coefficient vectors.

## 5.1   Function spaces

Finite element functions can be represented by the coefficients of their basis functions. Depending on the type of function space, one or multiple coefficients (or basis functions, respectively) are logically associated to respective micro-primitives. Figure 8 illustrates some examples.

HyTeG provides the `Function` data structure for those constructions. `Function`s represent coefficient vectors and the corresponding finite element functions simultaneously. Essential functionality is provided through various methods, for example, for the computation of linear combinations of coefficient vectors, for the evaluation of the finite element function anywhere in the domain, or for the calculation of scalar products of the coefficient vectors.

`Function`s can be composed to construct additional function spaces. For instance, multiple scalar `Function`s can be combined to represent a vector-valued function using the `VectorFunction` class. An example composition for the $\mathbb{P}_2$ space, with vertex and edge DoFs, is illustrated in Figure 9.

Since HyTeG focuses on geometric multigrid methods, the underlying coefficient vectors are usually allocated on several levels of the mesh hierarchy. The range of levels is specified during allocation, and virtually all methods require the user to specify the refinement level to operate on. This embeds the notion of a grid hierarchy into the framework and allows for the convenient construction of multigrid methods by design.

## 5.2   Operator design

The majority of iterative solvers and preconditioners only require the result of the action of an operator (e.g., matrix-vector multiplication) to a vector. However, they do not need explicit access to the matrix entries. Since the arithmetic intensity of matrix-vector operations using standard sparse matrix formats
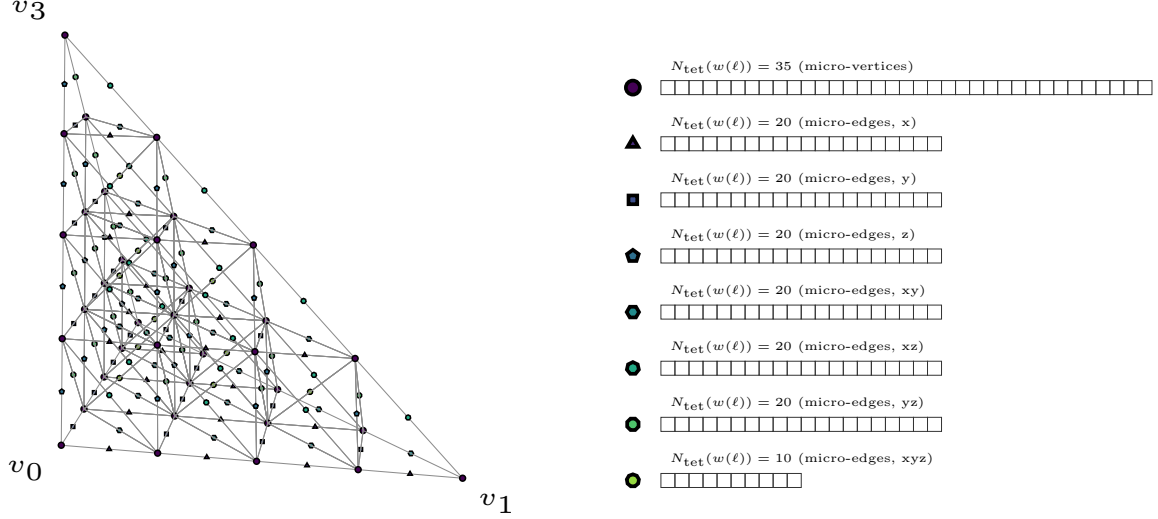
**Figure 9** *Illustration of the composition of DoFs of several micro-primitive subgroups into one function space data structure. In this example, one DoF is allocated for each micro-vertex and each micro-edge, corresponding to a $\mathbb{P}_2$ discretization. The respective arrays for the individual micro-primitive subgroups are displayed on the right. DoFs that belong to the same subgroup are allocated contiguously.*

is low and does not favor modern computer architectures, and the amount of memory required to store even sparse matrices for extreme-scale applications is not always available, *matrix-free* methods are critical to the efficient solution of extreme-scale problems [13, 57, 44, 29, 28, 45, 5, 46, 40]. The general idea is to (re-)compute the relevant matrix entries on-the-fly during the operator application instead of precomputing and storing the entire system matrix. This reduces storage space and bandwidth pressure since the matrix entries do not have to be loaded from memory at the cost of additional arithmetic operations [31].

In HyTeG, linear operators are therefore designed with a focus on *matrix-free* computations. The underlying block-structured grids enable crucial optimization techniques, outlined below and studied in, e.g., [16, 27, 9, 40].

The general interface of all operators requires the implementation of an `apply()` method that applies the operator to a vector (i.e., to a `Function` object) and writes the result to a target vector. The `apply()` method of a concrete operator in HyTeG comprises two main ingredients:

- The specification of the approximation to the bilinear form on a single element and

- the iteration pattern over the structured mesh on a macro-primitive.

The former defines a routine that computes local element matrices. It is derived from the weak formulation of the underlying PDE, the finite element function spaces, and the order of the quadrature formula.

The latter specifies how the element matrices are applied to a vector. While different iteration patterns are often mathematically equivalent, they must be carefully selected to achieve optimal performance. The performance of an iteration pattern may depend on the memory layout, the type of update rule (matrix-vector multiplication, relaxation, etc.), the function spaces, the PDE, and the underlying hardware. To achieve flexibility and maintainability, implementations should be composable, such that, for instance, optimizations to iteration patterns are independent of bilinear forms.

Two typical iteration patterns for finite element discretizations are illustrated in Figure 10 and described in the following.

### 5.2.1 Element-wise

Directly inferred from the standard derivation of the finite element method (FEM), *element-wise* update routines iterate over the grid elements. The weak form integrals are evaluated (or approximated) on each element to construct the local stiffness matrix. To perform a matrix-vector multiplication, this matrix is multiplied with the DoFs associated with the local element, and the result is added to the destination vector.
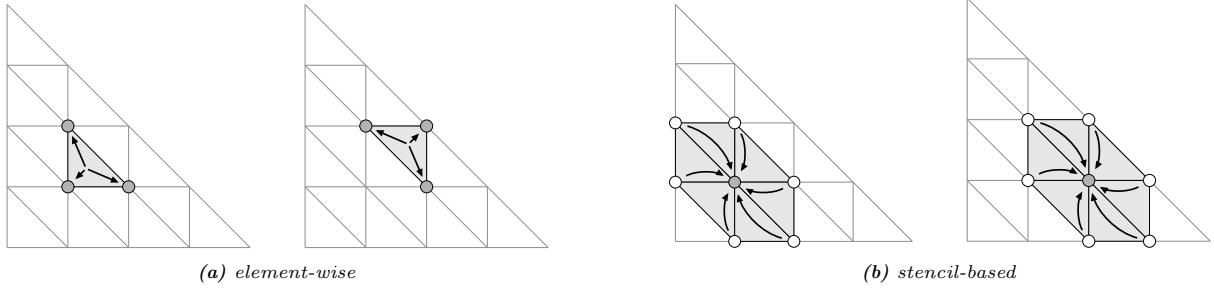
**(a)** *element-wise*        **(b)** *stencil-based*

***Figure 10*** *Illustrations of the element-wise and stencil-based update patterns. Each cartoon shows one update step for a discretization with vertex DoFs as an example. In Figure 10a (element-wise pattern), two elements are traversed, with read-write access to all associated DoFs. In Figure 10b (stencil-based pattern) two DoFs are traversed, with read-only access to DoFs marked with empty nodes, and read-write access to the center DoF, marked by a filled node. The iteration order is not fixed and could be chosen differently than indicated in the figures. To achieve the best possible performance, it must be selected in accordance with the underlying memory layout.*

As an example, consider the standard weak formulation of the Poisson equation. The (global) stiffness matrix is computed via

$$a_{ij} = \int_\Omega \nabla \phi_i \cdot \nabla \phi_j \, \mathrm{d}x = \sum_{T \in \mathcal{T}_h(\Omega)} \int_T \nabla \phi_i \cdot \nabla \phi_j \, \mathrm{d}x, \tag{5.1}$$

with basis functions $\phi_i$, and $\phi_j$. The entry $a_{kl}^{(T)}$ of the local stiffness matrix $A^{(T)}$ of the element $T$ is computed as

$$a_{kl}^{(T)} = \int_T \nabla \phi_{m_T(k)} \cdot \nabla \phi_{m_T(l)} \, \mathrm{d}x, \tag{5.2}$$

where $k$ and $l$ are indices of DoFs associated with the element $T$. $m_T(\cdot)$ maps the element-local indices of the basis functions to their global indices. The computation of the corresponding entry $a_{ij}$ of the global system matrix may require contributions from elements other than $T$. So instead of computing $a_{ij}$, the result of the matrix-vector product is computed by adding the local element contributions into the destination vector:

$$y_{m_T(k)} \leftarrow y_{m_T(k)} + \sum_l a_{kl}^{(T)} x_{m_T(l)}. \tag{5.3}$$

Figure 10a illustrates this pattern. Obviously, **y** has to be appropriately initialized before the iteration.

The advantage of this approach is that operations for the integration only have to be performed once per element. This includes, for example, calculating the transformation to the reference element.

On the downside, the actual entries $a_{ij}$ of the global system matrix are not computed explicitly. This means that certain matrix-vector operations cannot be performed in this way. In particular, many preconditioners require access to the diagonal entries of the system matrix. Those usually have to be precomputed and stored in a vector. More complicated access patterns, such as those required by Gauss-Seidel-type smoothers, are not suited for strictly element-wise iterations in general. Popular smoothers that only require matrix-vector products and diagonal entries include the Chebyshev accelerated Jacobi iteration [1]. Also, element-wise patterns involve multiple store operations per vector entry when the contributions from multiple elements are summed up. This induces more memory traffic than is theoretically necessary.

A comprehensive discussion of this element-wise approach can be found in [20, 44, 45], although the approach in [20] is not matrix-free (the local element matrices are precomputed and stored).

### 5.2.2 Stencil-based

Instead of iterating over the elements of the grid, *stencil-based* update patterns iterate over the DoFs of the destination vector. We can interpret a stencil $S_i$ as a bipartite graph that relates an entry $y_i$ of the destination vector to those entries $x_j$ of the source vector that correspond to the non-zero entries in row $i$ of the system matrix:

$$S_i := \{(j, a_{ij}) : a_{ij} \neq 0\}. \tag{5.4}$$

```
1   // Class template for a stencil-based operator
2   // that exploits the absence of space-dependent coefficients.
3   template< typename P1Form > P1ConstantOperator { ... };
4
5   // The template can be used to implement different operators via a Form object.
6
7   // Laplace operator
8   using P1ConstantLaplaceOperator =
9       P1ConstantOperator< forms::p1_diffusion_affine_q2 >;
10  // Mass operator ('qe' indicates analytical integration)
11  using P1ConstantMassOperator = P1ConstantOperator< forms::p1_mass_affine_qe >;
12
13  // For operators with variable coefficients,
14  // these optimizations cannot be exploited.
15  // A different implementation must be used.
16  template< typename P1Form > P1ElementwiseOperator { ... };
17
18  // div( k(x) grad ) operator with a space-dependent coefficient k.
19  using P1DivKGradElementwiseOperator =
20      P1ElementwiseOperator< forms::p1_div_k_grad_affine_q3 >;
```

**Code Listing 2** *Examples of linear operators in* HyTeG. *The implementation is opaque to the user, and all listed operators implement the* **apply()** *method to perform the corresponding matrix-vector multiplication. Depending on the implementation, different additional functionality may be supported, e.g., stencil-based operators can be employed as a Gauss-Seidel smoother.*

The computation of $S_i$ requires the evaluation of integrals over all elements that are associated with $y_i$. Using the same example as for the element-wise case, we get for each row $i$:

$$a_{ij} = \sum_{T \in \mathcal{T}_h(\Omega)} \int_T \nabla \phi_i \cdot \nabla \phi_j \, \mathrm{d}x, \quad \text{for all } j. \tag{5.5}$$

The result of a matrix-vector product can then be computed for a single row $i$ by setting

$$y_i \leftarrow \sum_j a_{ij} x_j. \tag{5.6}$$

An illustration is shown in Figure 10b.

The advantage of this method is that the entries of the system matrix are computed explicitly, enabling the implementation of a variety of matrix-vector operations that are not feasible with element-wise patterns (such as Gauss-Seidel type smoothers). Also, each destination vector entry is written to only once.

On the downside, each update requires the evaluation of integrals on different elements, and elements are traversed multiple times. Although only a subset of the entries of the local stiffness matrix needs to be computed on each element, some terms depend on the element itself and are therefore computed repeatedly.

The performance of stencil-based update patterns has been studied extensively, e.g., in [16, 29, 40].

Neither of the approaches is generally superior. It depends on many factors, such as the underlying grid structure, the finite element discretization, the desired solver, and the properties of the employed hardware. In some instances, the element-wise and stencil-based approaches are equivalent. Consider, for example, standard DG discretizations. Since each DoF is only associated with a single element, there is no difference between the element-wise iteration and the stencil-based pattern.

### 5.2.3 Optimizations

Apart from reducing bandwidth pressure and memory consumption, one of the decisive advantages of matrix-free approaches compared to pure sparse matrix computations is that they can exploit the underlying properties of the discretization (grid structure, basis functions, etc.) to perform heavy optimizations. An overview of standard techniques is given below, with examples of operator definitions implemented in HyTeG in Code Listing 2.

**(Block-)Structured grids** If the underlying grid is structured, matrix-free routines benefit in multiple ways. Most importantly, the alignment of the iteration pattern and the memory layout is possible and enables the exploitation of data locality strategies of modern cache-based processors. In particular, this avoids scattered memory accesses in favor of contiguous accesses to the linearized data and enables vectorization. Bandwidth pressure further reduces due to the on-the-fly computation of element vertex coordinates. On unstructured grids, coordinates must be loaded from memory.

Eventually, the evaluation of the integrals can be optimized since the element geometry is known a priori. The transformations from the computational elements to the reference element can be precomputed for a small set of different elements and possibly even simplified for specific grids (axis-aligned hexahedral elements only need to be scaled and translated).

Optimizations that exploit block-structured grids are central to the design and performance of matrix-free methods in HHG and HyTeG, and are covered in detail in, e.g., [16, 28, 40, 41].

**Constant coefficients**  Additional optimizations is possible depending on the underlying PDE. If the operator does not involve a variable coefficient, the integration can be simplified or even computed analytically.

If both the underlying grid is structured and the operator does not involve a variable coefficient, the element matrices and stencils for each type of element can be precomputed, which translates to only storing one (or several) rows of the sparse system matrix on each macro-primitive, *regardless* of the refinement level. This is arguably one of the most crucial optimizations that can be applied in matrix-free implementations and is exploited heavily in, e.g., [27, 40].

**Surrogates**  If the grid exhibits a specific structure, but the operator involves a space-dependent coefficient, the local element matrices and stencil entries are not constant either. A key observation is that, for sufficiently smooth coefficients, the entries of the local stiffness matrix and the entries of the stencil also vary smoothly over the domain. Those coefficients can then be approximated, for instance, via polynomials. Instead of evaluating the weak form integrals to compute a coefficient, the precomputed surrogate polynomial is evaluated. Obviously, this introduces errors, but if implemented carefully may drastically reduce the computational effort to obtain the entries of either the local element matrices or the stencil and still grants a sufficient approximation of the operator. A series of articles emerged from this idea  [11, 9, 10, 12, 21, 22, 23, 24, 25].

**Grid-aware precomputation**  Depending on the discretization and hardware, matrix-free approaches are not always superior to standard sparse-matrix implementations. As a rule of thumb, matrix-free implementations for low-order discretizations still exhibit significant pressure on the bandwidth, so it is unclear a priori, whether matrix-free approaches outperform standard sparse matrix computations [47]. Still, if memory *consumption*, i.e., storage space, is the limiting resource, matrix-free methods are generally necessary.

However, the structured access patterns can still be exploited in the non-matrix-free case. Instead of using a sparse format such as compressed row storage (CRS), the local stiffness matrices or stencils can be stored in a linearized layout corresponding to the iteration pattern. While this still requires loading the coefficients from memory, it avoids the frequent indirections associated with standard sparse formats. This approach is discussed, for example, in [20, 9, 55], and is implemented in HyTeG.

# 6  Demonstration

This section presents results on the numerical approximation of solutions to two model problems to demonstrate the flexibility and performance of the presented data structures in nontrivial use cases. The performance and scalability of both demonstrator applications are subject to ongoing research and are expected to be analyzed in forthcoming articles. For extreme-scale studies with up to more than a trillion ($10^{12}$) DoFs on more than $140\,000$ processes, we refer to [40, 39].

## 6.1  Enriched Galerkin elements for the Stokes system

Let $\Omega \in \mathbb{R}^d$ be a bounded domain, $\vec{f}$ a given forcing term, and $\mu$ a given space-dependent viscosity. We consider the Stokes system

$$
\begin{aligned}
-\nabla \cdot (2\mu\epsilon(\vec{u}) - pI) &= \vec{f} \quad \text{in } \Omega, \\
\nabla \cdot \vec{u} &= 0 \quad \text{in } \Omega, \\
\vec{u} &= \vec{g} \quad \text{on } \partial\Omega,
\end{aligned}
\tag{6.1}
$$

where $\vec{u}$ and $p$ are the velocity and pressure solutions, $\epsilon(\vec{u}) = \frac{1}{2}\left(\nabla\vec{u} + (\nabla\vec{u})^T\right)$ is the symmetric gradient, and $\vec{g}$ the given velocity boundary condition.

**(a)** *EG velocity element*    **(b)** *EG pressure element*    **(c)** *Nédélec element*
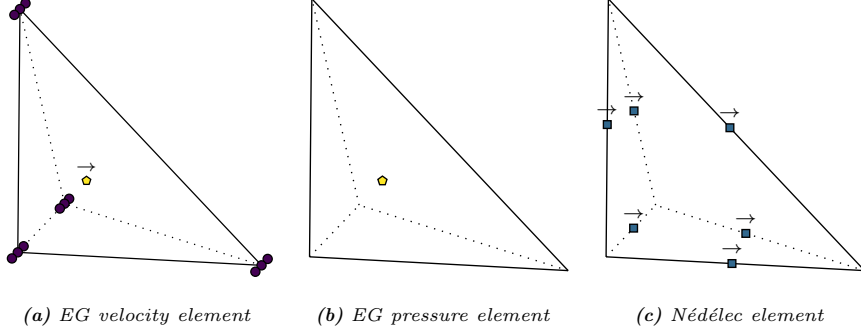
**Figure 11**    *Velocity (Figure 11a) and pressure (Figure 11b) elements as employed for the EG discretization of the Stokes system, and Nédélec element (Figure 11c) used for the discretization of the curl-curl problem. The arrows indicate that the respective DoFs correspond to vector-valued basis functions.*

The Stokes system (6.1) is discretized via the enriched Galerkin (EG) approach described in [61, 62]. This discretization employs a piecewise linear conforming velocity space that is enriched by a vector-valued, globally discontinuous component:

$$\mathbf{X}_h := \mathbb{P}_1 \oplus \left\{ \phi \in [L^2(\Omega)]^d : \phi|_T = c_T(\mathbf{x} - \mathbf{x}_T),\ c_T \in \mathbb{R},\ T \in \mathcal{T}_h \right\}, \tag{6.2}$$

where $\mathbb{P}_1$ is the standard continuous Galerkin space of piecewise linear elements, and $\mathbf{x}_T$ is the centroid of the element $T \in \mathcal{T}_h$. For the pressure, the $\mathbb{P}_0$ piecewise constant space is selected. Figures 11a and 11b illustrate the corresponding elements and DoFs. A single DoF in the tetrahedron's volume represents vector-valued enrichment. The discontinuity in the velocity (which is only in $L^2(\Omega)$) enforces a DG weak formulation and requires the computation of jumps and averages, and edge and face integrals [53].

The discretization is inf-sup stable [61]. Since the enrichment only adds a single additional DoF per element to the piecewise linear velocity space, it requires far fewer DoFs per element than other inf-sup stable combinations like the traditional Taylor-Hood pair ($\mathbb{P}_2$-$\mathbb{P}_1$) and is thus an interesting choice for extreme-scale applications with billions or trillions of DoFs.

The EG discretization requires the construction of composite spaces, combining scalar vertex DoFs and vectorial DoFs within the volumes for the velocity, with scalar volume DoFs for the pressure. The matrix-free operators map between the different spaces, e.g., from $\mathbb{P}_1$ to the enrichment space and back. Furthermore, the DG weak formulation entails integrals over lower-dimensional geometries like edges in 2D and triangles in 3D. Those integrals also lead to dependencies across edges and faces and complicate the update patterns during the operator application. See [18] for implementation details and further analysis. Figure 12a illustrates the vertex-centered EG stencil that couples the velocity with itself in 3D.

We showcase the discretization for a Stokes problem with strong viscosity variations, as frequently encountered in Earth mantle convection models [52, 32]. One test case that mimics this problem characteristic is the so-called *multi-sinker benchmark* [49, 54]. A variable number of spherical high-viscosity inclusions (the sinkers) are randomly placed in a 3D domain. The viscosity rises exponentially at the boundaries of the sinkers, leading to ill-conditioned problems. Figure 12 visualizes the viscosity inclusions and the velocity solution computed by HyTeG's EG implementation.
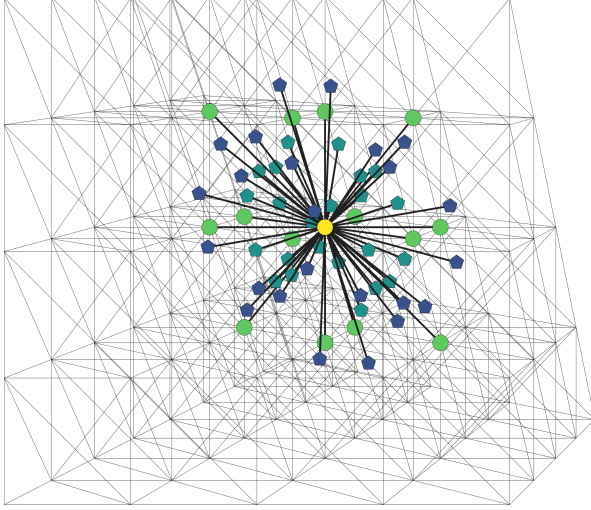
## 6.2    Nédélec elements for the curl-curl problem
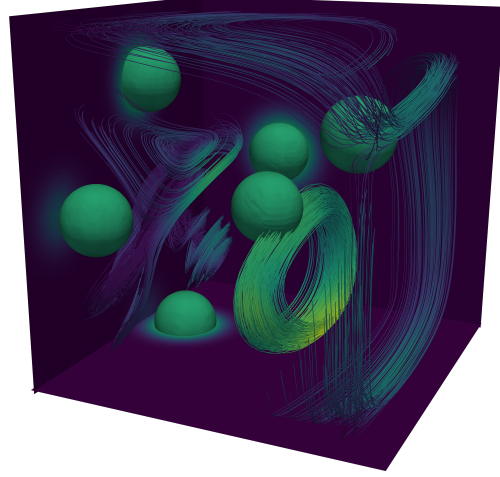
The (homogeneous) curl-curl problem

$$\begin{aligned} \alpha \,\mathbf{curl}\,\mathbf{curl}\,\vec{u} + \beta\vec{u} &= \vec{f} \quad \text{in } \Omega, \\ \vec{u} \times \vec{n} &= 0 \quad \text{on } \partial\Omega, \end{aligned} \tag{6.3}$$

in three dimensions ($d = 3$), with given $\vec{f}$, and $\alpha,\ \beta \in \mathbb{R}^+$, arises from Maxwell's equations [34]. Because standard piecewise Lagrangian elements are not suited for the discretization of (6.3), linear *Nédélec* elements of the first kind [50] are used to approximate $\vec{u}$ and $\vec{f}$. Figure 11c illustrates the corresponding vector-valued edge elements. Positioning the DoFs on the edges of the mesh enforces continuity of tangential components while components normal to cell faces are discontinuous. This matches the amount of continuity as required for conformity in $\mathbf{H}(\mathbf{curl})$ [2]. See [6] for implementation details.

We pick the solid torus as the domain $\Omega$ and discretize it with $65\,280$ macro-cells. Because the triangulation is only a rough approximation of the curved geometry, a curvilinear mapping is applied [30]. Figure 13a shows the resulting curvilinear mesh on the coarsest level.
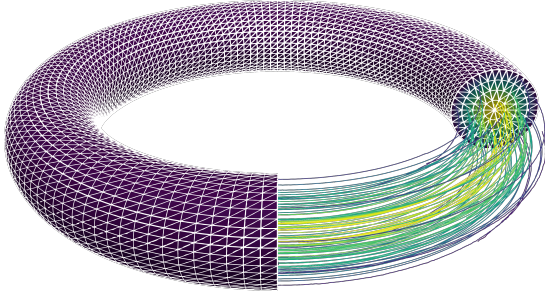
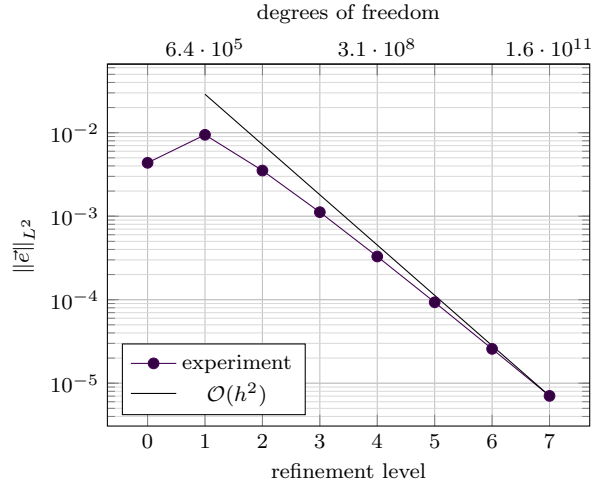**(a)** *vertex-centered EG velocity-velocity stencil in 3D*



**(b)** *sinker benchmark*

**Figure 12** *Illustration of a 3D stencil of the EG discretization and a visualization of the computed solution of the sinker benchmark. Figure 12a: The vertex-centered stencil represents the coupling of the velocity components at a vertex DoF. Circles correspond to vertex DoFs; pentagons correspond to volume DoFs. As shown in Figure 11a, 3 DoFs are allocated at each micro-vertex, and a vectorial DoF is placed in each micro-volume. Figure 12b: Spherical viscosity inclusions and velocity streamlines of a computed solution of the sinker benchmark. Six randomly placed sinkers exhibit an exponential viscosity decay of a factor of $10^3$. A multigrid preconditioned MINRES solver has been applied to solve the saddle point system with roughly $5 \cdot 10^5$ DoFs.*



**(a)** *coarse mesh and electric field lines of the solution on the solid torus*



**(b)** *$L^2$-error grid convergence*

**Figure 13** *Error convergence experiment using linear Nédélec elements of the first kind to solve the curl-curl problem (6.3) on the toroidal solid, discretized with $65\,280$ curvilinear macro-tetrahedra. (6.3) is solved using a matrix-free full multigrid (FMG) solver with 5 V(1,1) cycles per level. The total number of DoFs on the finest level is approximately $1.6 \cdot 10^{11}$.*

To evaluate the grid convergence of our discretization, we construct an analytic solution with homogeneous tangential boundary conditions and determine the right-hand-side $\vec{f}$ from it. For the sake of simplicity, we set $\alpha = \beta = 1$ in (6.3).

Numerically, we solve the PDE with a matrix-free FMG solver performing five V(1,1) cycles on each refinement level. Inside the V-cycles, Hiptmair's hybrid smoother [34] is used. Due to the non-elliptic nature of (6.3), standard smoothers can only reduce error components orthogonal to the nullspace of the curl-operator. Remaining error components in the nullspace of the curl-operator must be handled separately to obtain an effective multigrid scheme. To that end, the residual remaining after smoothing is determined, lifted to the space of scalar potentials, and discretized by $\mathbb{P}_1$ elements. Smoothing again in potential space removes the remaining oscillating error components. The smoothed $\mathbb{P}_1$ vector must then be transformed back to the space of Nédélec elements, where it is added to the current iterate. We
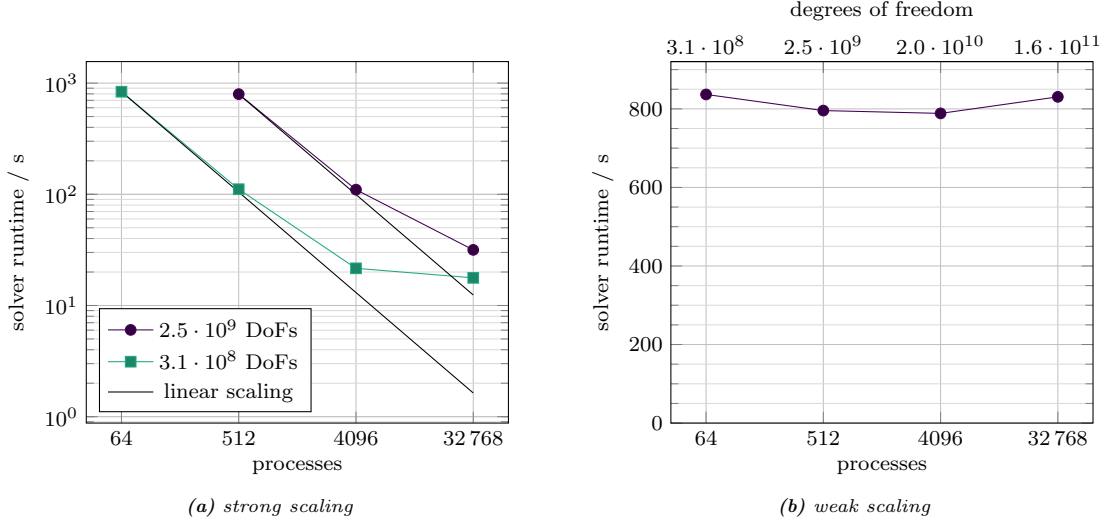
14

**(a)** *strong scaling*

**(b)** *weak scaling*

***Figure 14*** *Strong and weak scalability of the FMG solver for the curl-curl problem* (6.3). (6.3) *is solved on the toroidal solid (*65 280 *curvilinear macro-tetrahedra) on different refinement levels and with varying number of processes. The largest run comprises roughly* $1.6 \cdot 10^{11}$ *DoFs and is executed on* 32 768 *cores.*

choose Chebyshev smoothers of order 2 [1, 4] in both spaces. This means that one hybrid smoothing step requires in total three matrix-vector products in the Nédélec space, two matrix-vector products in $\mathbb{P}_1$, and two transfer operations between the spaces. Note that two additional $\mathbb{P}_1$ vectors must be allocated.

The system is solved up to refinement level 7, which comprises roughly $1.6 \cdot 10^{11}$ DoFs on 32 768 processes. The numerical solution is compared against the known analytic solution and the error is measured in the $L^2$-norm. The $L^2$-error is expected to reduce quadratically [33, Theorem 5.8, Remark 18]. According to Figure 13b, our convergence results agree with the theory.

For comparison, the same equation with jumping coefficients has recently been solved in [37] with an algebraic multigrid (AMG) preconditioned conjugate gradient (CG) solver implemented in the ParELAG miniapplication [51] in MFEM [42]. The authors scaled their system up to $1.4 \cdot 10^9$ DoFs on 4608 processes and report a total solve time of around three minutes. Furthermore, they used a graded hexahedral mesh containing distorted elements but no curvilinear transformation.

Next, we assess the strong and weak scalability of our solver. To that end, we solve the same system on 32 768, 4096, 512, and 64 processes. This corresponds to 256, 32, 4, and 1 compute nodes, respectively. The results are summarized in Figure 14.

We examine strong scalability for two problem sizes: $2.5 \cdot 10^9$ DoFs (level 5) and $3.1 \cdot 10^8$ DoFs (level 4). In the larger setup, when increasing the number of processes from 512 to 32 768, a parallel efficiency of 39% is observed. In the smaller case, increasing the processing cores from 64 to 4096 results in a high parallel efficiency of 61%. Further, solving the small problem ($3.1 \cdot 10^8$ DoFs) with 32 768 processes only yields a parallel efficiency of 9%. This result is the expected strong scaling behavior, given that solving the system with 4096 processes is a matter of a few seconds. Overall, better strong scalability is mainly hindered by the runtime on the coarser levels, where the arithmetic workload is very low compared to the amount of inter-process communication. Note that only a few ($< 10$) DoFs are allocated per process on the coarsest grids in both problem setups. Solving the system directly on a higher level (i.e., using a finer coarse mesh for the FMG solver) might have a positive impact on the scaling results.

On the other hand, weak scaling is not impacted by this effect. This shows in nearly constant runtimes over the range from 64 to 32 768 processes. Here, the scalability of the coarse grid solver is of less significance since 85% of the overall runtime is spent on the finest level (level 7).

# 7 Conclusion

This paper has provided a systematic approach to the design of data structures for arbitrary finite element discretizations on hybrid tetrahedral grids. It has categorized all geometric structures that emerge from regular refinement and provided associated illustrations. Particular focus has been put on the exploitation of the grid structure to enable efficient matrix-free kernels and corresponding memory layouts for the coefficient vectors, regardless of the element type. The flexibility of the implementation in the

finite element framework HyTeG have been showcased via the variable viscosity Stokes system and the curl-curl problem. Extreme-scalability has been demonstrated via the solution of curl-curl systems with approximately $1.6 \cdot 10^{11}$ unknowns on more than 32000 processes with a matrix-free full multigrid method. The presented considerations lay the ground for the formalization, generalization, and implementation of efficient compute kernels for large-scale, matrix-free finite element methods.

# Acknowledgments

# References

[1] M. ADAMS, M. BREZINA, J. HU, and R. TUMINARO, *Parallel Multigrid Smoothing: Polynomial versus Gauss–Seidel*, Journal of Computational Physics 188.2 (July 2003), pp. 593–610, ISSN: 00219991, DOI: `10.1016/S0021-9991(03)00194-3`.

[2] D. N. ARNOLD, *Finite Element Exterior Calculus*, Society for Industrial and Applied Mathematics, Philadelphia, PA, Dec. 2018, ISBN: 978-1-61197-553-6 978-1-61197-554-3, DOI: `10.1137/1.9781611975543`.

[3] D. N. ARNOLD and A. LOGG, *Periodic Table of the Finite Elements*, Siam News 47.9 (2014), p. 212.

[4] A. H. BAKER, R. D. FALGOUT, T. V. KOLEV, and U. M. YANG, *Multigrid Smoothers for Ultraparallel Computing*, SIAM J. Sci. Comput. 33.5 (Jan. 2011), pp. 2864–2887, ISSN: 1064-8275, 1095-7197, DOI: `10.1137/100798806`.

[5] P. BASTIAN, E. H. MÜLLER, S. MÜTHING, and M. PIATKOWSKI, *Matrix-Free Multigrid Block-Preconditioners for Higher Order Discontinuous Galerkin Discretisations*, Journal of Computational Physics 394 (Oct. 2019), pp. 417–439, ISSN: 00219991, DOI: `10.1016/j.jcp.2019.06.001`.

[6] D. BAUER, *Multigrid in H(curl) on Hybrid Tetrahedral Grids*, MA thesis, Friedrich-Alexander-Universität Erlangen-Nürnberg (FAU), 2023.

[7] M. BAUER, S. EIBL, C. GODENSCHWAGER, N. KOHL, M. KURON, C. RETTINGER, F. SCHORNBAUM, C. SCHWARZMEIER, D. THÖNNES, H. KÖSTLER, and U. RÜDE, *waLBerla: A Block-Structured High-Performance Framework for Multiphysics Simulations*, Computers & Mathematics with Applications 81 (Jan. 2021), pp. 478–501, ISSN: 08981221, DOI: `10.1016/j.camwa.2020.01.007`.

[8] S. BAUER, H.-P. BUNGE, D. DRZISGA, S. GHELICHKHAN, M. HUBER, N. KOHL, M. MOHR, U. RÜDE, D. THÖNNES, and B. WOHLMUTH, *TerraNeo — Mantle Convection Beyond a Trillion Degrees of Freedom*, in: *Softw. Exascale Comput. - SPPEXA 2016-2019*, ed. by H.-J. BUNGARTZ, S. REIZ, B. UEKERMANN, P. NEUMANN, and W. NAGEL, vol. 136, Lecture Notes in Computational Science and Engineering, Springer, 2020, pp. 569–610, DOI: `10.1007/978-3-030-47956-5_19`.

[9] S. BAUER, D. DRZISGA, M. MOHR, U. RÜDE, C. WALUGA, and B. WOHLMUTH, *A Stencil Scaling Approach for Accelerating Matrix-Free Finite Element Implementations*, SIAM J. Sci. Comput. 40.6 (Jan. 2018), pp. C748–C778, ISSN: 1064-8275, 1095-7197, DOI: `10.1137/17M1148384`.

[10] S. BAUER, M. HUBER, S. GHELICHKHAN, M. MOHR, U. RÜDE, and B. WOHLMUTH, *Large-Scale Simulation of Mantle Convection Based on a New Matrix-Free Approach*, Journal of Computational Science 31 (Feb. 2019), pp. 60–76, ISSN: 18777503, DOI: `10.1016/j.jocs.2018.12.006`.

[11] S. BAUER, M. MOHR, U. RÜDE, J. WEISMÜLLER, M. WITTMANN, and B. WOHLMUTH, *A Two-Scale Approach for Efficient on-the-Fly Operator Assembly in Massively Parallel High Performance Multigrid Codes*, Applied Numerical Mathematics 122 (Dec. 2017), pp. 14–38, ISSN: 01689274, DOI: `10.1016/j.apnum.2017.07.006`.

[12] S. BAUER, M. HUBER, M. MOHR, U. RÜDE, and B. WOHLMUTH, *A New Matrix-Free Approach for Large-Scale Geodynamic Simulations and Its Performance*, in: *Computational Science – ICCS 2018*, ed. by Y. SHI, H. FU, Y. TIAN, V. V. KRZHIZHANOVSKAYA, M. H. LEES, J. DONGARRA, and P. M. A. SLOOT, vol. 10861, Springer International Publishing, Cham, 2018, pp. 17–30, ISBN: 978-3-319-93700-7 978-3-319-93701-4, DOI: `10.1007/978-3-319-93701-4_2`.

[13] B. BERGEN, T. GRADL, F. HULSEMANN, and U. RUDE, *A Massively Parallel Multigrid Method for Finite Elements*, Comput. Sci. Eng. 8.6 (Nov. 2006), pp. 56–62, ISSN: 1521-9615, DOI: `10.1109/MCSE.2006.102`.

[14] B. BERGEN, F. HULSEMANN, and U. RUDE, *Is $1.7 \times 10^{10}$ Unknowns the Largest Finite Element System that Can Be Solved Today?*, in: *ACMIEEE SC 2005 Conf. SC05*, IEEE, Seattle, WA, USA, 2005, pp. 5–5, ISBN: 978-1-59593-061-3, DOI: `10.1109/SC.2005.38`.

---

[2] `https://gauss-allianz.de/en/project/title/CoMPS`
[3] `https://www.nhr-verein.de`

[15] B. BERGEN, *Hierarchical Hybrid Grids: Data Structures and Core Algorithms for Efficient Finite Element Simulations on Supercomputers*, Advances in Simulation, SCS Publishing House, 2006.

[16] B. K. BERGEN and F. HÜLSEMANN, *Hierarchical Hybrid Grids: Data Structures and Core Algorithms for Multigrid*, Numer. Linear Algebra Appl. 11.23 (Mar. 2004), pp. 279–291, ISSN: 1070-5325, 1099-1506, DOI: 10.1002/nla.382.

[17] J. BEY, *Tetrahedral Grid Refinement*, Computing 55.4 (Dec. 1995), pp. 355–378, ISSN: 0010-485X, 1436-5057, DOI: 10.1007/BF02238487.

[18] F. BÖHM, *Matrix-Free Implementation and Evaluation of the Enriched Galerkin Finite Element Method for the Stokes Problem with Varying Viscosity*, MA thesis, Friedrich-Alexander-Universität Erlangen-Nürnberg (FAU), 2023.

[19] A. BUTTARI, M. HUBER, P. LELEUX, T. MARY, U. RÜDE, and B. WOHLMUTH, *Block Low-rank Single Precision Coarse Grid Solvers for Extreme Scale Multigrid Methods*, Numer Linear Algebra Appl 29.1 (Jan. 2022), ISSN: 1070-5325, 1099-1506, DOI: 10.1002/nla.2407.

[20] G. F. CAREY, E. BARRAGY, R. MCLAY, and M. SHARMA, *Element-by-Element Vector and Parallel Computations*, Commun. appl. numer. methods 4.3 (May 1988), pp. 299–307, ISSN: 0748-8025, 1555-2047, DOI: 10.1002/cnm.1630040303.

[21] D. DRZISGA, B. KEITH, and B. WOHLMUTH, *The Surrogate Matrix Methodology: A Priori Error Estimation*, SIAM J. Sci. Comput. 41.6 (Jan. 2019), A3806–A3838, ISSN: 1064-8275, 1095-7197, DOI: 10.1137/18M1226580.

[22] D. DRZISGA, B. KEITH, and B. WOHLMUTH, *The Surrogate Matrix Methodology: A Reference Implementation for Low-Cost Assembly in Isogeometric Analysis*, MethodsX 7 (2020), p. 100813, ISSN: 22150161, DOI: 10.1016/j.mex.2020.100813.

[23] D. DRZISGA, B. KEITH, and B. WOHLMUTH, *The Surrogate Matrix Methodology: Accelerating Isogeometric Analysis of Waves*, Computer Methods in Applied Mechanics and Engineering 372 (Dec. 2020), p. 113322, ISSN: 00457825, DOI: 10.1016/j.cma.2020.113322.

[24] D. DRZISGA, B. KEITH, and B. WOHLMUTH, *The Surrogate Matrix Methodology: Low-cost Assembly for Isogeometric Analysis*, Computer Methods in Applied Mechanics and Engineering 361 (Apr. 2020), p. 112776, ISSN: 00457825, DOI: 10.1016/j.cma.2019.112776.

[25] D. DRZISGA, U. RÜDE, and B. WOHLMUTH, *Stencil Scaling for Vector-Valued PDEs on Hybrid Grids With Applications to Generalized Newtonian Fluids*, SIAM J. Sci. Comput. 42.6 (Jan. 2020), B1429–B1461, ISSN: 1064-8275, 1095-7197, DOI: 10.1137/19M1267891.

[26] D. DRZISGA, A. WAGNER, and B. WOHLMUTH, *A Matrix-Free ILU Realization Based on Surrogates*, ArXiv Prepr. (2022), DOI: 10.48550/ARXIV.2210.15280.

[27] B. GMEINER, M. HUBER, L. JOHN, U. RÜDE, and B. WOHLMUTH, *A Quantitative Performance Study for Stokes Solvers at the Extreme Scale*, Journal of Computational Science 17 (Nov. 2016), pp. 509–521, ISSN: 18777503, DOI: 10.1016/j.jocs.2016.06.006.

[28] B. GMEINER, U. RÜDE, H. STENGEL, C. WALUGA, and B. WOHLMUTH, *Performance and Scalability of Hierarchical Hybrid Multigrid Solvers for Stokes Systems*, SIAM J. Sci. Comput. 37.2 (Jan. 2015), pp. C143–C168, ISSN: 1064-8275, 1095-7197, DOI: 10.1137/130941353.

[29] B. GMEINER, U. RÜDE, H. STENGEL, C. WALUGA, and B. WOHLMUTH, *Towards Textbook Efficiency for Parallel Multigrid*, Numer. math.: theory methods appl. 8.1 (Feb. 2015), pp. 22–46, ISSN: 1004-8979, 2079-7338, DOI: 10.4208/nmtma.2015.w10si.

[30] W. J. GORDON and C. A. HALL, *Transfinite Element Methods: Blending-function Interpolation over Arbitrary Curved Element Domains*, Numer. Math. 21.2 (Apr. 1973), pp. 109–129, ISSN: 0029-599X, 0945-3245, DOI: 10.1007/BF01436298.

[31] G. HAGER and G. WELLEIN, *Introduction to High Performance Computing for Scientists and Engineers*, 0th ed., CRC Press, July 2010, ISBN: 978-1-4398-1193-1, DOI: 10.1201/EBK1439811924.

[32] T. HEISTER, J. DANNBERG, R. GASSMÖLLER, and W. BANGERTH, *High Accuracy Mantle Convection Simulation through Modern Numerical Methods – II: Realistic Models and Problems*, Geophys. J. Int. 210.2 (Aug. 2017), pp. 833–851, ISSN: 0956-540X, 1365-246X, DOI: 10.1093/gji/ggx195.

[33] R. HIPTMAIR, *Finite Elements in Computational Electromagnetism*, Acta Numerica 11 (Jan. 2002), pp. 237–339, ISSN: 0962-4929, 1474-0508, DOI: 10.1017/S0962492902000041.

[34] R. HIPTMAIR, *Multigrid Method for Maxwell's Equations*, SIAM J. Numer. Anal. 36.1 (Jan. 1998), pp. 204–225, ISSN: 0036-1429, 1095-7170, DOI: 10.1137/S0036142997326203.

[35] M. HUBER, B. GMEINER, U. RÜDE, and B. WOHLMUTH, *Resilience for Massively Parallel Multigrid Solvers*, SIAM J. Sci. Comput. 38.5 (2016), S217–S239, DOI: 10.1137/15M1026122.

[36] M. HUBER, U. RÜDE, C. WALUGA, and B. WOHLMUTH, *Surface Couplings for Subdomain-Wise Isoviscous Gradient Based Stokes Finite Element Discretizations*, J Sci Comput 74.2 (Feb. 2018), pp. 895–919, ISSN: 0885-7474, 1573-7691, DOI: 10.1007/s10915-017-0470-3.

[37] D. Z. KALCHEV, P. S. VASSILEVSKI, and U. VILLA, *Parallel Element-Based Algebraic Multigrid for $\boldsymbol{H}(\underline{curl})$ and $\boldsymbol{H}(div)$ Problems Using the ParELAG Library*, SIAM Journal on Scientific Computing 45.3 (June 30, 2023), S371–S400, ISSN: 1064-8275, 1095-7197, DOI: 10.1137/21M1433253.

[38] N. KOHL, J. HÖTZER, F. SCHORNBAUM, M. BAUER, C. GODENSCHWAGER, H. KÖSTLER, B. NESTLER, and U. RÜDE, *A Scalable and Extensible Checkpointing Scheme for Massively Parallel Simulations*, The International Journal of High Performance Computing Applications 33.4 (July 2019), pp. 571–589, ISSN: 1094-3420, 1741-2846, DOI: 10.1177/1094342018767736.

[39] N. KOHL, M. MOHR, S. EIBL, and U. RÜDE, *A Massively Parallel Eulerian-Lagrangian Method for Advection-Dominated Transport in Viscous Fluids*, SIAM J. Sci. Comput. 44.3 (June 2022), pp. C260–C285, ISSN: 1064-8275, 1095-7197, DOI: 10.1137/21M1402510.

[40] N. Kohl and U. Rüde, *Textbook Efficiency: Massively Parallel Matrix-Free Multigrid for the Stokes System*, SIAM J. Sci. Comput. 44.2 (Apr. 2022), pp. C124–C155, issn: 1064-8275, 1095-7197, doi: 10.1137/20M1376005.

[41] N. Kohl, D. Thönnes, D. Drzisga, D. Bartuschat, and U. Rüde, *The* HyTeG *Finite-Element Software Framework for Scalable Multigrid Solvers*, International Journal of Parallel, Emergent and Distributed Systems 34.5 (Sept. 2019), pp. 477–496, issn: 1744-5760, 1744-5779, doi: 10.1080/17445760.2018.1506453.

[42] T. Kolev and V. Dobrev, *Modular Finite Element Methods (MFEM)*, Language: en, 2010, doi: 10.11578/DC.20171025.1248.

[43] H. Köstler, M. Heisig, N. Kohl, S. Kuckuk, M. Bauer, and U. Rüde, *Code Generation Approaches for Parallel Geometric Multigrid Solvers*, Analele Univ. Ovidius Constanta - Ser. Mat. 28.3 (Dec. 2020), pp. 123–152, issn: 1844-0835, doi: 10.2478/auom-2020-0038.

[44] M. Kronbichler and K. Kormann, *A Generic Interface for Parallel Cell-Based Finite Element Operator Application*, Computers & Fluids 63 (June 2012), pp. 135–147, issn: 00457930, doi: 10.1016/j.compfluid.2012.04.012.

[45] M. Kronbichler and K. Kormann, *Fast Matrix-Free Evaluation of Discontinuous Galerkin Finite Element Operators*, ACM Trans. Math. Softw. 45.3 (Sept. 2019), pp. 1–40, issn: 0098-3500, 1557-7295, doi: 10.1145/3325864.

[46] M. Kronbichler, D. Sashko, and P. Munch, *Enhancing Data Locality of the Conjugate Gradient Method for High-Order Matrix-Free Finite-Element Implementations*, The International Journal of High Performance Computing Applications (July 2022), p. 109434202211078, issn: 1094-3420, 1741-2846, doi: 10.1177/10943420221107880.

[47] M. Kronbichler and W. A. Wall, *A Performance Comparison of Continuous and Discontinuous Galerkin Methods with Fast Multigrid Solvers*, SIAM J. Sci. Comput. 40.5 (Jan. 2018), A3423–A3448, issn: 1064-8275, 1095-7197, doi: 10.1137/16M110455X.

[48] C. Lengauer, *Loop Parallelization in the Polytope Model*, in: *CONCUR'93*, ed. by G. Goos, J. Hartmanis, and E. Best, vol. 715, Springer Berlin Heidelberg, Berlin, Heidelberg, 1993, pp. 398–416, isbn: 978-3-540-57208-4 978-3-540-47968-0, doi: 10.1007/3-540-57208-2_28.

[49] D. A. May, J. Brown, and L. L. Pourhiet, *pTatin3D: High-Performance Methods for Long-Term Lithospheric Dynamics*, in: *SC14 Int. Conf. High Perform. Comput. Netw. Storage Anal.* IEEE, New Orleans, LA, USA, Nov. 2014, pp. 274–284, isbn: 978-1-4799-5500-8 978-1-4799-5499-5, doi: 10.1109/SC.2014.28.

[50] J. C. Nedelec, *Mixed finite elements in* $\mathbb{R}^3$, Numer. Math. 35.3 (Sept. 1980), pp. 315–341, issn: 0029-599X, 0945-3245, doi: 10.1007/BF01396415.

[51] *ParELAG Mini Applications in MFEM*, url: http://github.com/mfem/mfem/tree/master/miniapps/parelag.

[52] Y. Ricard, *Physics of Mantle Convection*, in: *Mantle Dynamics*, ed. by D. Bercovici, vol. 7, Treatise on Geophysics, Elsevier, 2007, pp. 31–89, doi: 10.1016/B978-0-444-53802-4.00127-5.

[53] B. Rivière, *Discontinuous Galerkin Methods for Solving Elliptic and Parabolic Equations: Theory and Implementation*, Society for Industrial and Applied Mathematics, Jan. 2008, isbn: 978-0-89871-656-6 978-0-89871-744-0, doi: 10.1137/1.9780898717440.

[54] J. Rudi, G. Stadler, and O. Ghattas, *Weighted BFBT Preconditioner for Stokes Flow Problems with Highly Heterogeneous Viscosity*, SIAM J. Sci. Comput. 39.5 (Jan. 2017), S272–S297, issn: 1064-8275, 1095-7197, doi: 10.1137/16M108450X.

[55] H. D. Tran, M. Fernando, K. Saurabh, B. Ganapathysubramanian, R. M. Kirby, and H. Sundar, *A Scalable Adaptive-Matrix SPMV for Heterogeneous Architectures*, in: *2022 IEEE Int. Parallel Distrib. Process. Symp. IPDPS*, IEEE, Lyon, France, May 2022, pp. 13–24, isbn: 978-1-66548-106-9, doi: 10.1109/IPDPS53621.2022.00011.

[56] U. Trottenberg, C. W. Oosterlee, and A. Schüller, *Multigrid*, Academic Press, San Diego, 2001, isbn: 978-0-12-701070-0.

[57] P. E. Vos, S. J. Sherwin, and R. M. Kirby, *From h to p Efficiently: Implementing Finite and Spectral/Hp Element Methods to Achieve Optimal Performance for Low- and High-Order Discretisations*, Journal of Computational Physics 229.13 (July 2010), pp. 5161–5181, issn: 00219991, doi: 10.1016/j.jcp.2010.03.031.

[58] C. Waluga, B. Wohlmuth, and U. Rüde, *Mass-Corrections for the Conservative Coupling of Flow and Transport on Collocated Meshes*, Journal of Computational Physics 305 (Jan. 2016), pp. 319–332, issn: 00219991, doi: 10.1016/j.jcp.2015.10.044.

[59] J. Weismüller, B. Gmeiner, S. Ghelichkhan, M. Huber, L. John, B. Wohlmuth, U. Rüde, and H.-P. Bunge, *Fast Asthenosphere Motion in High-resolution Global Mantle Flow Models*, Geophys. Res. Lett. 42.18 (Sept. 2015), pp. 7429–7435, issn: 0094-8276, 1944-8007, doi: 10.1002/2015GL063727.

[60] S. Williams, A. Waterman, and D. Patterson, *Roofline: An Insightful Visual Performance Model for Multicore Architectures*, Commun. ACM 52.4 (Apr. 2009), pp. 65–76, issn: 0001-0782, 1557-7317, doi: 10.1145/1498765.1498785.

[61] S.-Y. Yi, X. Hu, S. Lee, and J. H. Adler, *An Enriched Galerkin Method for the Stokes Equations*, Computers & Mathematics with Applications 120 (Aug. 2022), pp. 115–131, issn: 08981221, doi: 10.1016/j.camwa.2022.06.018.

[62] S.-Y. Yi, S. Lee, and L. Zikatanov, *Locking-Free Enriched Galerkin Method for Linear Elasticity*, SIAM J. Numer. Anal. 60.1 (Feb. 2022), pp. 52–75, issn: 0036-1429, 1095-7170, doi: 10.1137/21M1391353.

# A    Illustrations of all micro-primitive subgroups



**(a)** *micro-vertices*

**(b)** *micro-edges, subgroup x*

**(c)** *micro-edges, subgroup y*

**(d)** *micro-edges, subgroup z*

**(e)** *micro-edges, subgroup xy*

**(f)** *micro-edges, subgroup xz*

**(g)** *micro-edges, subgroup yz*

**(h)** *micro-edges, subgroup xyz*

**Figure 15**   *Micro-vertices, and all 7 types of micro-edges, illustrated on refinement level 2.*

*(a) subgroup z-up*

*(b) subgroup z-down*

*(c) subgroup y-up*

*(d) subgroup y-down*

*(e) subgroup x-up*

*(f) subgroup x-down*

*(g) subgroup xyz-up*

*(h) subgroup xyz-down*

*(i) subgroup xy-up*

*(j) subgroup xy-down*

*(k) subgroup yz-up*

*(l) subgroup yz-down*

**Figure 16**  *All 12 types of micro-faces, illustrated on refinement level 2.*

**(a)** subgroup I-up

**(b)** subgroup I-down

**(c)** subgroup II-up

**(d)** subgroup II-down
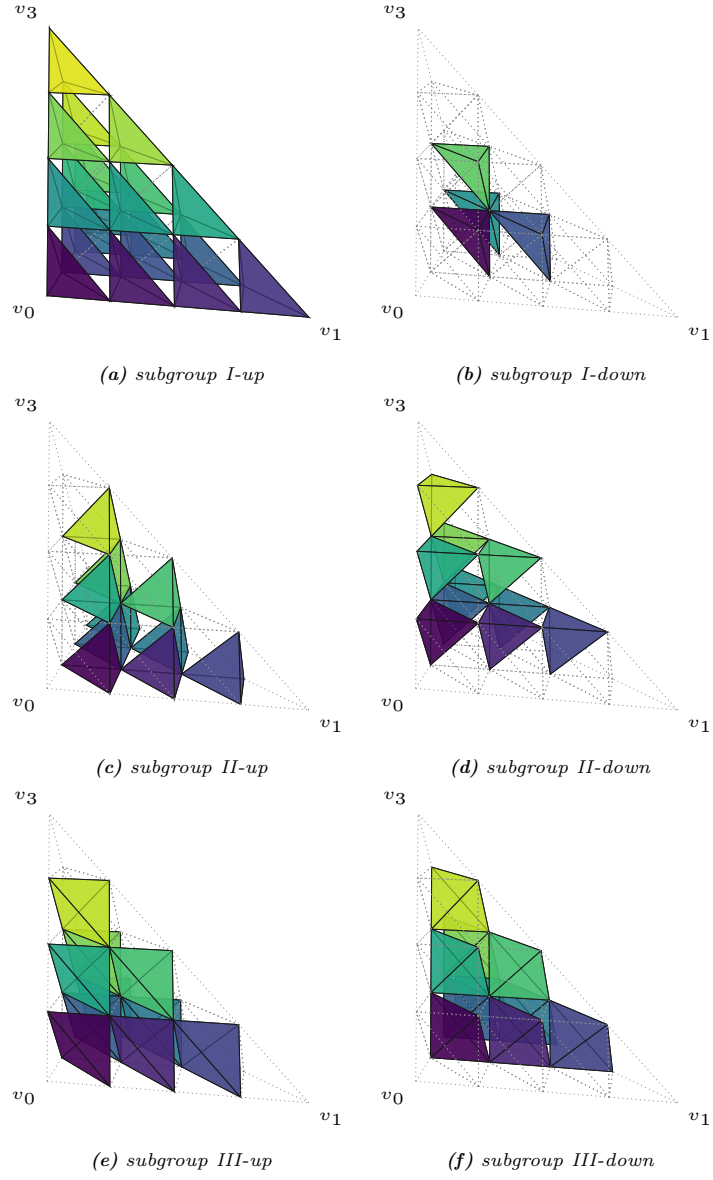
**(e)** subgroup III-up

**(f)** subgroup III-down

**Figure 17**  *All 6 types of micro-cells, illustrated on refinement level 2.*