

1 **FULLY PARALLEL MESH I/O USING PETSC DMPLEX WITH AN**
2 **APPLICATION TO WAVEFORM MODELING***

3 VACLAV HAPLA[†], MATTHEW G. KNEPLEY[‡], MICHAEL AFANASIEV[†], CHRISTIAN
4 BOEHM[†], MARTIN VAN DRIEL[†], LION KRISCHER[†], AND ANDREAS FICHTNER[†]

5 **Abstract.** Large-scale PDE simulations using high-order finite-element methods on unstruc-
6 tured meshes are an indispensable tool in science and engineering. The widely used open-source
7 PETSc library offers an efficient representation of generic unstructured meshes within its DMPlex
8 module. This paper details our recent implementation of parallel mesh reading and topological in-
9 terpolation (computation of intermediate codimensions from a cell-vertex mesh) into DMPlex. We
10 apply these developments to seismic wave propagation scenarios on Mars as a sample application.
11 The principal motivation is to overcome single-node memory limits and reach mesh sizes which were
12 impossible before. Moreover, we demonstrate that scalability of I/O and topological interpolation
13 goes beyond 12'000 cores, and memory imposed limits on maximum mesh size vanish.

14 **Key words.** unstructured mesh, directed acyclic graph, partitioning, topological interpolation,
15 parallel I/O, PETSc, DMPlex, seismic waveform modeling, Spectral Element Method

16 **AMS subject classifications.** 65-04, 65Y05, 65M50, 05C90, 35L05

17 **1. Introduction.** Finite-element methods (FEM) [37, 24, 51] are widely used
18 in science and engineering to simulate complex physical systems. Many applications
19 require FEM discretization with high polynomial order on large unstructured meshes
20 requiring distributed memory computer architectures [1, 14, 46, 34, 2]. Here, two
21 complications may arise:

22 (1) The distributed memory mesh representation should explicitly include mesh
23 entities of all codimensions (vertices, edges, faces, cells). However, mesh datafiles
24 typically store only vertices and cells to avoid redundant storage, disk operations,
25 and to conform to widely used formats. Edges and faces can be computed in runtime
26 using a process which will be called here *topological interpolation*.

27 (2) Meshing tools typically create a single datafile, but loading the whole mesh
28 onto one processor and distributing onto all remaining processors becomes a bottle-
29 neck for a sufficiently large mesh. Instead, we need to load different portions of the
30 mesh file directly onto target processors and maintain a distributed mesh representa-
31 tion right from beginning. Additional load balancing and redistribution can be done
32 using a parallel partitioner.

33 PETSc DMPlex (section 2) is a flexible mesh implementation which addresses
34 (1) by design, i.e., it can represent any number of codimensions and implements
35 topological interpolation. Moreover, it is completely agnostic to the mesh shape, di-
36 mensionality and cell types. However, only with developments presented in this paper,
37 also (2) has been addressed. This in turn brought another challenge to implement the
38 topological interpolation from (1) in parallel.

39 In this paper, we do not deal with hierarchical mesh representations such as oc-
40 tree data structures [20, 47]. They are advantageous for certain data access patterns,
41 especially in the context of Adaptive Mesh Refinement (AMR) [20, 12], at the ex-
42 pense of generality; DMPlex supports, for instance, arbitrary mesh partitions and
43 extraction of arbitrary subsets of cells (or facets) as submeshes: features which are

*Submitted to the editors on April 18, 2020.

[†]ETH Zurich, Switzerland (vaclav.hapla@erdw.ethz.ch, andreas.fichtner@erdw.ethz.ch).

[‡]University at Buffalo, NY (knepley@buffalo.edu).

44 typically missing from hierarchical meshing frameworks [26]. Note that PETSc offers
 45 the DMForest wrapper of p4est [13, 25, 46], and conversion between DMPlex and
 46 DMForest [26].

47 We demonstrate the efficiency and potential of our new DMPlex functionality on
 48 a real world application in the context of the full waveform modeling. The spectral-
 49 element method on unstructured conforming hexahedral meshes has become the de-
 50 facto standard for global-scale simulations of seismic waves [1, 18, 19, 44]. We apply
 51 it to simulate full 3D high-frequency wave propagation on Mars, based on data from
 52 the NASA InSight mission [6]. This consists in solving a coupled system of the elastic
 53 and acoustic wave equations. To accurately model these data in the desired frequency
 54 band, large scale simulations are required. In the presented simulation, more than
 55 100 million 4-th order hexahedral mesh elements were used.

56 Before the developments presented in this manuscript, with mesh reading and
 57 topological interpolation being serial, mesh sizes were limited by the single node
 58 memory. On our testing platform, Piz Daint supercomputer at the Swiss National
 59 Supercomputing Center, the mesh size limit was approximately 16 million elements,
 60 rendering such simulations impossible.

61 The manuscript is organized as follows. First, we describe abstractions for mesh
 62 data management of unstructured meshes in high-order finite-element discretizations
 63 using PETSc DMPlex. Then we explain our new strategies for the parallel simulation
 64 startup (mesh reading, topological interpolation and redistribution) on distributed
 65 memory HPC architectures. We continue with a brief introduction of the spectral-
 66 element method and its implementation which uses DMPlex. Next, waveform mod-
 67 eling benchmarks demonstrate scalable performance of the parallel startup for up to
 68 256 million hexahedral mesh elements, running on up to 1024 Cray XC50 nodes of Piz
 69 Daint. Finally, the mentioned Mars seismic wave propagation simulation is presented
 70 as a practical application.

71 **2. DMPlex.** PETSc [3, 4, 5, 10] is a well-known library for numerical meth-
 72 ods, used by the scientific and engineering computing communities. It provides par-
 73 allel data management, structured and unstructured meshes, linear and nonlinear
 74 algebraic solvers and preconditioners, time integrators, optimization algorithms and
 75 others. Many of these methods (such as geometric multigrid and domain decomposi-
 76 tion solvers) can take advantage of the geometric/topological setting of a discretized
 77 problem, i.e., mesh information.

78 DMPlex is a PETSc module for generic unstructured mesh storage and operations.
 79 It decouples user applications from the implementation details of common mesh-
 80 related utility tasks, such as file I/O and mesh partitioning. It represents the mesh
 81 topology in a flexible way, providing topological connectivity of mesh entities at all
 82 codimensions (vertices, edges, faces and cells), crucial for high-order finite-element
 83 method (FEM) simulations, and provides a wide range of common mesh management
 84 functionalities.

85 PETSc’s interface for serving mesh data to numerical algorithms is the DM object.
 86 PETSc has several DM implementations. The native implementations of structured
 87 grids (DMDA), staggered grids DMStag, and unstructured meshes (DMPlex) have the
 88 most complete coverage of the DM API, and are developed most actively. Besides
 89 these two, PETSc has several DM implementations that wrap external libraries, such
 90 as DMMOAB for MOAB [48] and DMFOREST for p4est [13]. Here we will focus on DMPlex
 91 which proved to be most relevant for the discussed waveform modeling applications.

92 DMPlex encapsulates the topology of unstructured grids and provides a wide range

93 of common mesh management functionalities to application programmers [45, 36, 35,
 94 7]. It provides a domain topology abstraction that decouples user applications from
 95 the implementation details of common mesh-related utility tasks, such as file I/O,
 96 mesh partitioning, and parallel load balancing [33]. It aims to increase extensibility
 97 and interoperability between scientific applications through librarization [5, 10].

98 **2.1. Mesh representation and basic operations.** DMPlex uses an abstract
 99 representation of the unstructured mesh topology, where the connectivity of topo-
 100 logical entities (vertices, edges, faces, cells) is stored as a directed acyclic graph
 101 (DAG) [38, 32], also referred to as a Hasse Diagram in topology. Let us call this
 102 representation a *plex* and refer to the topological mesh entities as *points*. The plex
 103 is constructed of clearly defined layers (*strata*) that enable access to mesh entities
 104 by their topological codimension without explicit reference to the overall topological
 105 dimension d of the mesh, see Figure 1(b). Note also that d can be reconstructed in
 106 linear time from an unlabeled plex using Depth First Search. Let us denote a set of
 107 points at the same stratum as $\text{Stratum}(h)$, where h is an integer *height*, $h \leq d + 1$.

108 All points in the plex share a single consecutive numbering, emphasizing that
 109 each point is treated equally no matter its shape or dimension, and allowing DMPlex
 110 to store the graph connectivity in a single array where dimensional layers are defined
 111 as consecutively numbered subranges. The directional connectivity of the plex is
 112 defined by the covering relation called *cone*, denoted here as $C(p)$, yielding a set of all
 113 points directly connected to p in the next codimension. The transitive closure of the
 114 cone operation shall be denoted by $C^+(p)$. Both $C(p)$ and $C^+(p)$ are illustrated in
 115 Figure 1(d). The dual operation called *support* and denoted $S(p)$, and its transitive
 116 closure $S^+(p)$ are shown in Figure 1(e).

117 In addition to the abstract topology data, PETSc provides two utility objects to
 118 describe the parallel data layout: a `PetscSection` object maps the graph-based topol-
 119 ogy information to discretized solution data through an offset mapping, and a star for-
 120 est (`PetscSF`, see subsection 3.2.1) object holds a one-sided description of shared data
 121 in parallel. These data layout mappings allow DMPlex to manage distributed solution
 122 data by automating the preallocation of distributed vector and matrix data structures
 123 and performing halo data exchanges. Moreover, by storing grid topology alongside
 124 discretized solution data, DMPlex is able to provide the mappings required for sophis-
 125 ticated preconditioning algorithms, such as geometric multigrid methods [11, 17] and
 126 multiblock, or “fieldsplit” preconditioning for multiphysics problems [9].

127 **2.2. Topological interpolation.** For high order methods we are interested in,
 128 we need mesh entities of all codimensions (vertices, edges, faces, cells) be present
 129 in the memory mesh representation. However, usual mesh generation algorithms or
 130 mesh file readers result in a mesh with cells and vertices only, while edges and faces
 131 need to be inferred at runtime. This is accomplished by the so-called topological
 132 interpolation.

133 Topological Interpolation is the process of constructing intermediate levels of the
 134 ranked poset describing a mesh, given information at bracketing levels. For example,
 135 if we receive triangles and their covering vertices, as in Figure 2, interpolation will
 136 construct edges. The first algorithm for interpolation on the Hasse diagram was
 137 published in [38], but this version is only appropriate for simplices, ignores orientation
 138 of the mesh points, and did not give a complexity bound.

139 The interpolation procedure selects a given point stratum as cells, for which it will
 140 construct facets. It iterates over the cells, whose cones are an oriented set of vertices.
 141 The two essential operations are to extract an oriented facet from lower dimensional

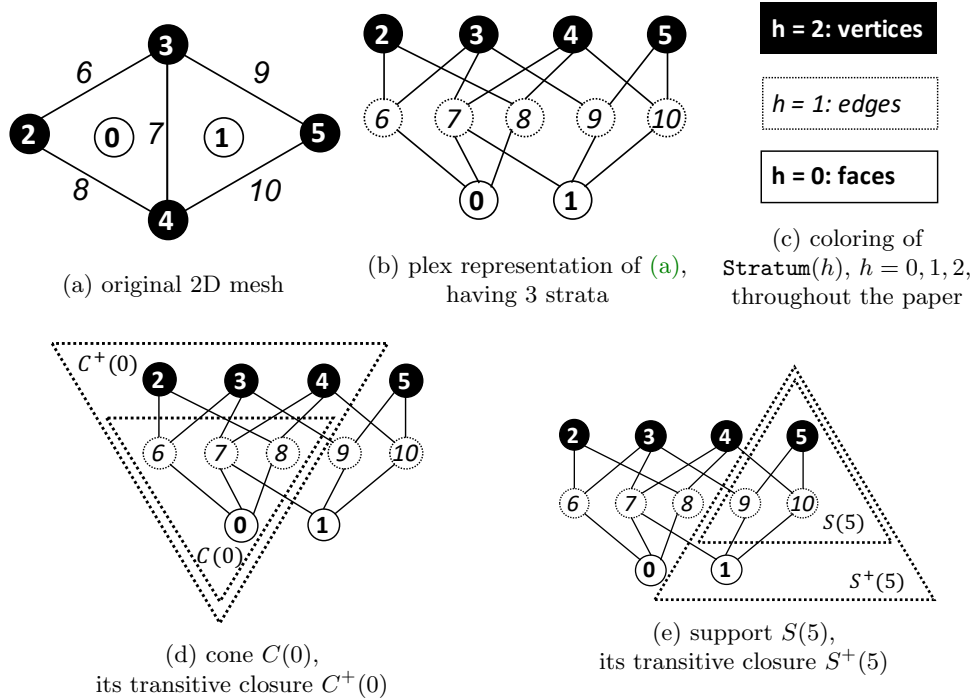


Fig. 1: DMPlex mesh representation and basic relations. A 1D mesh would have 2 strata: $h = 0$ edges, $h = 1$ vertices. A 3D mesh would have 4 strata: $h = 0$ cells, $h = 1$ faces, $h = 2$ edges, $h = 3$ vertices. Note we can also refer to the $h = 1$ entities commonly, in a dimension-independent way, as *facets*, and $h = 0$ as *cells*. The entities at $h = d$, where d is the topological dimension, are always vertices.

142 points, and then attach it to a cell with the correct orientation. Orientation of mesh
 143 points is detailed in [subsection 2.3](#). In order to enumerate the facets for a given
 144 cell type, we have `DMPlexGetRawFaces_Internal()` which returns an oriented list of
 145 vertices for each facet of the input cell.

146 An initial iteration over cells constructs all facets, and enters them into a hash
 147 table, where the hash key is the sorted list of vertices in each face. We need one pass
 148 for each type of face. Once the hash table is constructed, we know the number of
 149 new facets to be inserted, and can allocate a new plex. The new plex is identical to
 150 the old, except that it has a new face stratum, and the cone sizes of cells need to be
 151 calculated anew. For example, hexahedral cells have 8 vertices, but 6 facets, so that
 152 cone size would change from 8 to 6.

153 A second iteration over cells inserts the facets. We clear the hash table and repeat
 154 the face extraction above. If a face is missing from the table, we insert it into the table,
 155 record its cone, and also insert it into the cell cone with default orientation. If instead
 156 it is present in the table, we insert it into the cell cone with orientation computed from
 157 comparing the face cone with that returned from `DMPlexGetRawFaces_Internal()`.
 158 Again, we need one pass for each face type.

159 The complexity to interpolate a given stratum is in $\mathcal{O}(N_C N_F N_T)$, where N_C is the

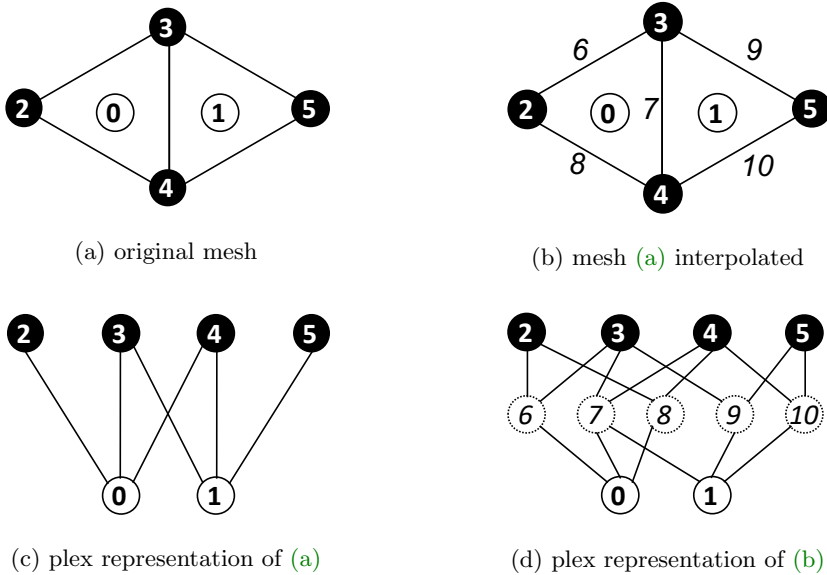


Fig. 2: Sequential topological interpolation: original mesh and interpolated mesh in classical and plex representation.

160 number of cells, N_F is the maximum number of faces per cell, and N_T is the number
 161 of used face types. However, N_F and N_T are obviously constant, so $\mathcal{O}(N_C N_F N_T)$
 162 $= \mathcal{O}(N_C)$. The \mathcal{O} -complexity remains the same even if we sum all strata because
 163 the size of (number of points in) stratum $h = k + 1$ is a certain multiple of the size
 164 of stratum $h = k$ for any feasible height k . We can conclude that the complexity is
 165 linear with respect to the mesh size.

166 **2.3. Orienting edges and faces.** Let us focus on two particular points in
 167 the interpolated mesh in [Figure 2\(b\)](#): cell 0, and edge $7 \in C(0)$. Their cones are
 168 $C(0) = \{6, 7, 8\}$ and $C(7) = \{3, 4\}$, respectively. We have so far spoken about *which*
 169 points form $C(p)$ for given p . However, the orientation of p is also important; for
 170 instance, to have a well-defined direction of an outer normal, or to assign field values
 171 correctly during simulation. The orientation of p is represented as the order of points
 172 in $C(p)$. Hence, e.g. $C(0)$ is rather a tuple, $C(0) = (6, 7, 8)$. This is how the array
 173 implementing $C(0)$ is stored. Let us from now denote by $C(p, c)$ the c -th point in this
 174 tuple, $c = 0, 1, \dots, n - 1$, where $n = \text{size}(C(p))$.

175 For consistency, the orientation of points in plex should be in line with the ori-
 176 entation of their supporting points. There is never a problem for the lowermost and
 177 uppermost stratum; cells have no supporting points, and vertices have no orientation.
 178 However, for the intermediate strata (i.e. edges and faces computed by interpolation),
 179 we can get a conflicting situation as depicted in [Figure 3](#). Edge 7 points against the
 180 direction of cell 1 but if we flip it, it will point against the direction of cell 0. Thus
 181 we need a mechanism to allow for this.

182 In general, suppose point $p \in \text{Stratum}(0)$, its cone point $q = C(p, c) \in C(p) \subset$
 183 $\text{Stratum}(1)$, and $C(q) \subset \text{Stratum}(2)$. For example, in [Figures 3\(b\)](#) and [3\(d\)](#), $p = 1$,
 184 $c = 0$, $q = 7$, $C(q) = (3, 4)$. To compensate the given (stored) order of $C(q)$, additional

185 information about the proper interpretation order needs to be defined. This informa-
 186 tion must be attached to the edge (p, q) in the DAG because it varies for different
 187 choices of p even for the same q . Naturally, the order is not completely arbitrary; it
 188 can be described by (1) the starting point index $S(p, c) \geq 0$ in 0-based local numbering
 189 with respect to q , and (2) the direction $D(p, c) \in \{-1, 1\}$ (reverse/forward). These
 190 two can be represented by a single signed integer $O(p, c)$: $S(p, c)$ by its magnitude,
 191 and $D(p, c)$ by its sign. Since the sign is undefined for 0, the negative values are
 192 shifted by -1. To summarize,

$$193 \quad (2.1) \quad O(p, c) = \begin{cases} S(p, c), & D(p, c) = 1, \\ -S(p, c) - 1, & D(p, c) = -1, \end{cases}$$

194 and the other way around,

$$195 \quad (2.2) \quad D(p, c) = \begin{cases} 1, & O(p, c) \geq 0, \\ -1, & O(p, c) < 0, \end{cases}$$

$$196 \quad (2.3) \quad S(p, c) = \begin{cases} O(p, c), & O(p, c) \geq 0, \\ -O(p, c) - 1, & O(p, c) < 0. \end{cases}$$

197
 198 Note that for $C(p, c)$ being an edge, flipping the orientation of the edge implies chang-
 199 ing the starting point, so $O(p, c) \in \{0, -2\}$ only. We can also define the whole tuple
 200 of orientations for point p ,

$$201 \quad (2.4) \quad O(p) = (O(p, 0), \dots, O(p, n - 1)),$$

202 where $n = \text{size}(C(p))$. This $O(p)$ is attached to every plex point p in the same way as
 203 $C(p)$. We note that $O(p)$ is just a numbering for the elements of the dihedral group
 204 for a given face.

205 **3. Parallel simulation startup.** Let us call *startup phase* all steps necessary
 206 to load the mesh from disk storage and prepare it for use in the simulation time
 207 loop. It consists of the following steps: (1) raw data loading, (2) plex construction,
 208 (3) topological interpolation, (4) distribution.

209 Before the developments of this paper, these steps were serial and only the last
 210 step, a one-to-all distribution using a serial partitioner such as METIS [28, 30], re-
 211 sulted in the distributed mesh. This approach inevitably led to the upper limit on
 212 the mesh size due to the memory constraints of a single node of a cluster. There-
 213 fore we developed a new, completely parallel startup phase where all four steps are
 214 done in parallel right from the beginning. This parallel startup phase is schematically
 215 depicted in Figure 4. We further show that even for meshes that fit into memory,
 216 parallel startup can bring significant time and energy saving. Let us describe these
 217 stages more in detail in the following subsections.

218 **3.1. Raw data loading.** This stage forms the first part of our mesh reader
 219 implementation, and consists in reading distributed raw topology and geometry data
 220 by generic index set and vector readers, dominated by low level I/O operations.

221 **3.1.1. HDF5.** Hierarchical Data Format 5 (HDF5) [22] is a file format and
 222 library, designed to store and organize large amounts of N-dimensional array data. It is
 223 currently supported by the HDF Group, a not-for-profit corporation whose mission is
 224 to ensure continued development of HDF5 technologies and the continued accessibility
 225 of data stored in HDF.

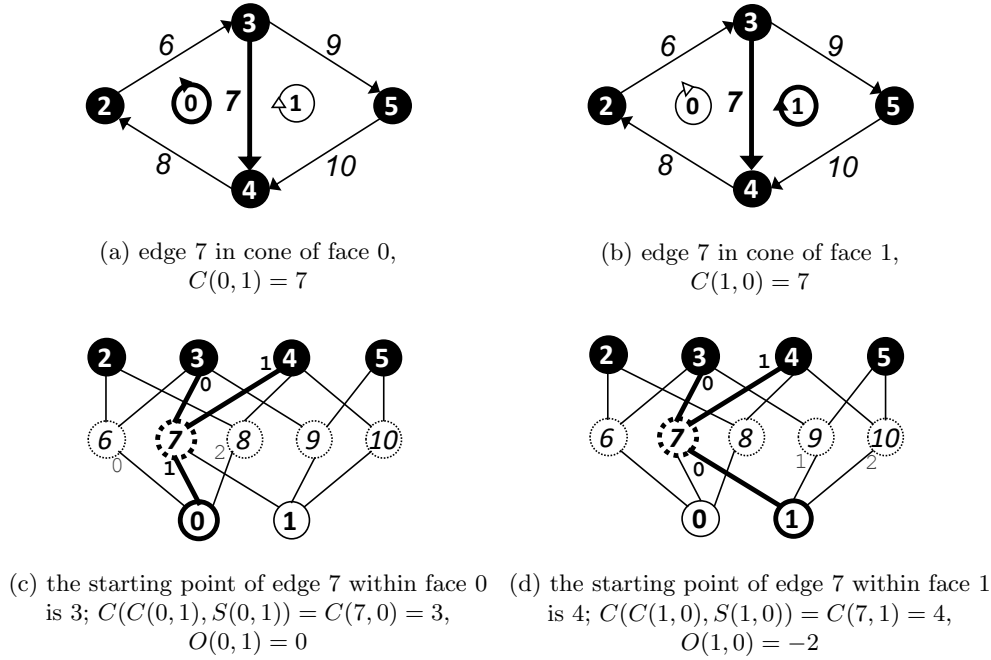


Fig. 3: Mesh from Figures 2(b) and 2(d) with cone points order and orientation. We focus here on the edge 7 within the cones of faces 0 and 1.

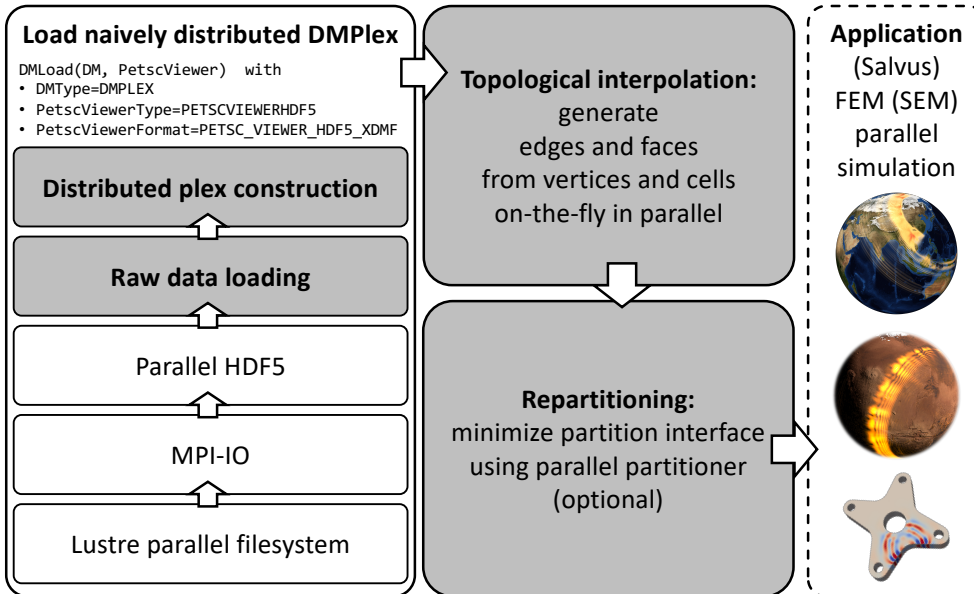


Fig. 4: Schematic diagram of the parallel simulation startup. Grey stages are within PETSc scope.

226 HDF5’s file structure includes two major types of objects: (1) datasets, multidimensional arrays of a homogeneous type; (2) groups, container structures which can
 227 hold datasets and other groups. Every HDF5 file has a root group /, under which one
 228 can add additional groups and datasets. This results in a hierarchical, filesystem-like
 229 data format. Resources in an HDF5 file can be accessed using the POSIX-like syntax
 230 /group1/group2/dataset. Metadata is stored in the form of user-defined, named
 231 attributes attached to groups and datasets [22].

233 HDF5 transparently handles how all the objects map to the actual bytes in the
 234 file. HDF5 actually provides an abstracted filesystem-within-a-file that is portable
 235 to any system with the HDF5 library installed, regardless of the underlying storage
 236 type, filesystem, or endianness. It does automatic conversions between storage
 237 datatypes (dictated by the data file) and runtime memory datatypes (dictated by the
 238 application) [22].

239 HDF5 supports parallel shared-file I/O using MPI-IO [40] capabilities which in
 240 turn provide scalable access to the underlying parallel filesystem such as Lustre [42].
 241 By default, HDF5 provides uniform access to all parts of the file for all processes
 242 of the communicator (passed to HDF5 using `H5Pset_fapl_mpio()`). However, since
 243 PETSc uses data parallelism, it would be very inefficient to load all data to all processes
 244 and then distribute them again. The most important functionality in this
 245 regard are *hyperslabs* which can read or write to a portion of a dataset. A hyperslab
 246 can be a logically contiguous collection of points in a dataspace, or a regular pattern
 247 of points or blocks in a dataspace. The hyperslab can select a separate chunk
 248 of the file for each process individually by means of a rank-dependent offset. The
 249 `H5Sselect_hyperslab()` function is used for this purpose [23].

250 **3.1.2. XDMF.** XDMF (eXtensible Data Model and Format) [50] is a mesh
 251 data file format. It distinguishes the metadata (light data) and the values themselves
 252 (heavy data). Light data and heavy data are stored using XML and HDF5, respectively.
 253 The data format is stored redundantly in both XML and HDF5. There are
 254 two crucial datasets describing the mesh in a minimal sufficient way: (1) `<Geometry>`,
 255 a 2D dataset where each row contains coordinates of a vertex (2 or 3 scalars based
 256 on dimensionality); (2) `<Topology>`, a 2D dataset where each row represents a cell,
 257 listing indices of all incident vertices. Each vertex index in `<Topology>` corresponds
 258 to a row index within `<Geometry>`. Both these datasets can be defined within the
 259 XML file as plain text, or refer to a dataset path within the standalone HDF5 file (e.g.
 260 `MyData.h5:/geometry/vertices`). Both ways can be mixed within the same XDMF
 261 file, and are both supported by widely used visualization programs such as ParaView,
 262 VisIt and EnSight. Nevertheless, the former way is advisable only for small datasets.
 263 We always store all data in the HDF5 file and use the XDMF file only as a descriptor.
 264 We will refer to this as HDF5/XDMF format.

265 **3.1.3. PETSc data loading.** PETSc contains a class, called `PetscViewer`, designated
 266 for all I/O of any `PetscObject` such as a vector, matrix, linear solver. The
 267 source/destination is dictated by the type (`PetscViewerType`) such as `ascii`, `binary`,
 268 `hdf5`, `socket` and others. A more fine-grained control of how the object is read/viewed
 269 is accomplished with the viewer format (`PetscViewerFormat`); e.g. `ascii_info` and
 270 `ascii_info_detail` print plain text information about the object with a different
 271 level of verbosity.

272 Most relevant for this work is that DM supports both reading and writing with
 273 `PetscViewer`, using `DMLoad(DM,PetscViewer)` and `DMView(DM,PetscViewer)`, respectively.
 274 Recently, we have implemented a new `PetscViewerFormat hdf5_xdmf`

275 for the `PetscViewerType hdf5` to enable parallel reading and writing of unstruc-
 276 tured meshes represented in memory as plex (DM with `DMType plex`) and stored in
 277 an HDF5 file with topology and geometry data compatible with XDMF (see [subsec-](#)
 278 [tions 3.1.1](#) and [3.1.2](#)). HDF5/XDMF has become the first widely used mesh format
 279 supported by PETSc which is both readable and writable in parallel.

280 As we are here interested mainly in the simulation startup phase, we will now
 281 focus only on `DMLoad()`, namely its implementation for the combination of `DMType`,
 282 `PetscViewerType` and `PetscViewerFormat` described above. This implementation
 283 relies on lower level readers for integer vectors (index sets, `IS`) and scalar vectors (`Vec`):
 284 `ISLoad()` and `VecLoad()`. They both read datasets from given paths within the
 285 HDF5 file (see [subsection 3.1.1](#)) using the `H5Dread()` routine but with different HDF5
 286 datatypes (`H5T_NATIVE_INT` and `H5T_NATIVE_DOUBLE`, respectively). They make use
 287 of a hyperslab (see [subsection 3.1.1](#)) reflecting the given parallel layout (`PetscLayout`)
 288 to divide the global dataset into local chunks along the first dimension. The layout is
 289 either specified by user, or calculated automatically so that the chunks' lengths differ
 290 by 1 at most. The second dimension of the dataset is interpreted as a block size,
 291 i.e., the resulting vector is divided into equally sized shorter blocks. Blocks can have
 292 various contextual meanings such as DOFs of the same element.

293 Within `DMLoad()`, `VecLoad()` loads the geometry information (`<Geometry>` in
 294 [subsection 3.1.2](#)) into a `Vec`, whose blocks and entries represent vertices and their
 295 coordinates, respectively. The size of all blocks is the same and corresponds to the
 296 spatial dimension of the mesh, and global indices of the blocks form an implicit global
 297 vertex numbering.

298 `ISLoad()` then loads the cell-vertex topology information which is represented in
 299 XDMF by `<Topology>` ([subsection 3.1.2](#)). Each `IS` block corresponds to an element,
 300 and each single entry refers to one of this element's vertices using the implicit global
 301 vertex numbering. Global indices of the blocks form an implicit global cell numbering.

302 When we read such representation in parallel, all processes load approximately
 303 equally sized, contiguously numbered, disjoint portions of vertices and cells in a single
 304 parallel I/O operation. Note that the global vertex and cell numberings are not
 305 affected by parallel loading on any number of processes.

306 **3.2. Distributed plex construction.** The raw topology and geometry data
 307 loaded in [subsection 3.1](#) need to be transformed into a plex representation. This forms
 308 the second part of `DMLoad()`, and is realized by a call to a communication-bounded
 309 operation described in detail in [subsection 3.2.2](#). The resulting `DMPlex` instance is
 310 naively distributed with vertices and cells only. Let us first describe a parallel star
 311 forest graph implementation, which allows gluing together serial plexes across different
 312 processes.

313 **3.2.1. Star Forest.** The Star Forest (SF, in PETSc called `PetscSF` [[4](#), [8](#)]) is a
 314 forest of star graphs. The *root* of each star graph corresponds to something owned
 315 by a process, such as a mesh point or solution degree-of-freedom (DOF). The *leaves*
 316 of each star graph are the shared versions of that point or DOF on other processes,
 317 or what is often called ghost points. In particular, the SF is a map from local points
 318 p to remote points (r, q) , where r is the remote rank and q is the local number for
 319 the point on process r . Since the SF only deals with local numberings, not global
 320 numberings, parallel mesh modification is much easier. Moreover, SF accepts data
 321 buffers with arbitrary MPI types for its communication routines.

322 SF supports a broadcast operation from roots to leaves, as well as a reduction
 323 from leaves to roots. In addition, it supports gather and scatter operations over the

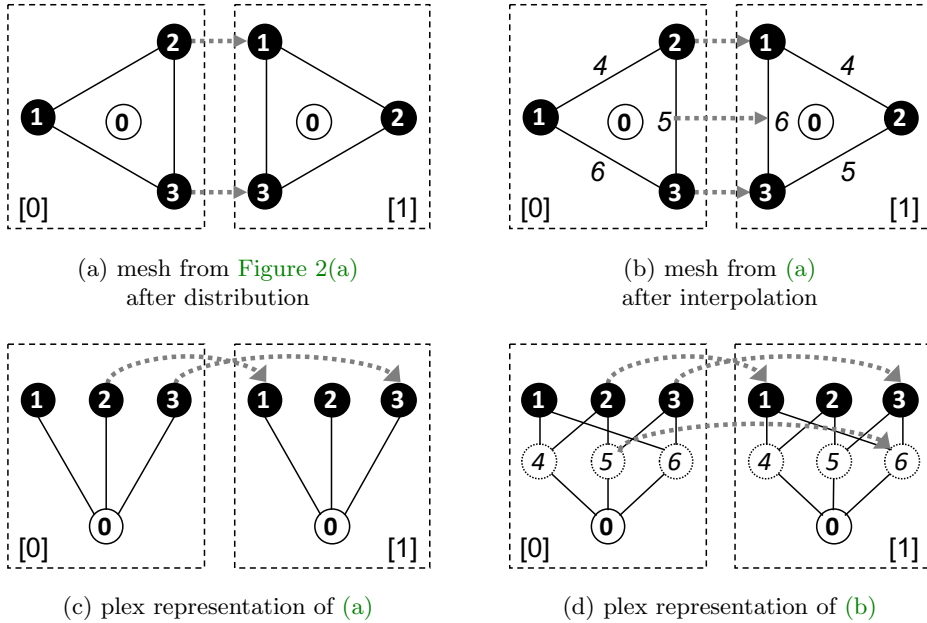


Fig. 5: Parallel DMPlex. Grey dotted arrows denote `sfPoint`.

324 roots, inversion of the SF to get two-sided information, and a fetch-and-op atomic
 325 update. Plex can build a hierarchy of SF objects describing the overlap of partitions,
 326 point set, topology, data layout, and Jacobian information [33].

327 **3.2.2. Construction of distributed plex from raw data.** Once the raw
 328 topology and geometry data is loaded as described in subsection 3.1.3, vertex coordi-
 329 nates are initially distributed independently of the cells, with process p getting some
 330 $N_V^{(p)} \times d$ vertex coordinates. For example, if we have 20 vertices in 3D and 4 pro-
 331 cesses, each process might get 15 doubles, namely the coordinates for 5 contiguous
 332 numbered vertices, regardless of the cells it has been assigned.

333 In order to construct a parallel plex from these vertex and cell portions, we employ
 334 `DMPlexCreateFromCellListParallel()`. From the cell-vertex topology information,
 335 it determines how many unique local vertices a process owns using a hash table of the
 336 vertex numbers. The cone information for a local plex is constructed by translating
 337 the global vertex numbers to local vertex numbers, using the index of each local
 338 vertex in the hash table. Once we have the local cone information, we symmetrize
 339 this to get support information (`DMPlexSymmetrize()`), and construct strata labels
 340 (`DMPlexStratify()`). Both these operations have linear complexity.

341 Given this hash table of local vertices and the global division, we construct an
 342 SF mapping the initial vertex division to the local division, `sfVert`. With this SF,
 343 we broadcast coordinates from the initial distribution roots to the local distribution
 344 leaves, determined by the mesh topology. In addition, we can use the initial SF to
 345 construct an SF describing the sharing of local vertices between processes, `sfPoint`.
 346 We first construct an array which holds (r, l) for each local vertex. Then we reduce this
 347 array to the roots using `sfVert`, taking the maximum over ranks. This information

348 is then broadcast back to the leaves, giving a unique owner for each vertex, allowing
 349 us to construct `sfPoint`.

350 The resulting distributed plex object has partition-wise complete topological in-
 351 formation. Each partition is a serial plex object which includes all incident vertices
 352 of its elements. These serial plex objects are combined together by the parallel SF
 353 object `sfPoint` as shown in [Figure 5](#).

354 **3.3. Parallel topological interpolation.** Since the steps above lead to a dis-
 355 tributed plex, we need a parallel version of the topological interpolation ([subsec-
 356 tion 2.2](#)). This step appeared to be the most challenging one. It consists of the serial
 357 interpolation (in-memory computations) and a small communication.

358 The first step consists in applying the sequential topological interpolation ([subsec-
 359 tion 2.2](#)) and cone orientation ([subsection 2.3](#)) on each rank independently. Then we
 360 must alter the `PetscSF` structure which identifies mesh points which are owned by dif-
 361 ferent processes, or *leaf* points. The SF structure is described in [subsection 3.2.1](#). We
 362 mark all leaf points which are adjacent to another ghost point as candidates. These
 363 candidate points are then gathered to *root* point owners (using `PetscSFbcst()`). For
 364 each candidate, the root checks for each point in the cone that either it owns that
 365 point in the SF or it is a local point. If so, it claims ownership. These claims are
 366 again broadcast, allowing a new SF to be created incorporating the new edges and
 367 faces.

368 The cone orientation has been done on each rank independently, and hence it is
 369 only partition-wise correct. However, we have not yet handled the following assump-
 370 tion: If interface edges/faces owned by different ranks represent the same geometrical
 371 entity, i.e., they are connected by `pointSF` like edges `[0]5` and `[1]6` in [Figure 6](#), they
 372 must have a conforming order of cone points (`[r]` means ownership by rank r). This
 373 requirement can be written more rigorously as an implication

$$374 \quad (3.1) \quad \left. \begin{array}{l} p_0 \quad \rightarrow p_1 \\ C(p_0) = (q_{0,0}, \dots, q_{0,n-1}) \\ C(p_1) = (q_{1,0}, \dots, q_{1,n-1}) \end{array} \right\} \Rightarrow \left\{ \begin{array}{l} q_{0,0} \quad \rightarrow q_{1,0} \\ \dots \\ q_{0,n-1} \rightarrow q_{1,n-1}, \end{array} \right.$$

376 using notation from [subsection 2.3](#) and relation \rightarrow meaning a connection via `pointSF`.

377 In [Figures 6\(a\)](#) and [6\(c\)](#) this assumption is violated for the edges `[0]5` and `[1]6`.
 378 They are flipped to each other, more rigorously speaking `pointSF` connects the edge
 379 and its incident vertices

$$380 \quad [0]5 \rightarrow [1]6, \quad [0]2 \rightarrow [1]1, \quad [0]3 \rightarrow [1]3,$$

382 but the order of cone points does not conform,

$$383 \quad [0]2 = C([0]5, 0) \not\rightarrow C([1]6, 0) = [1]3,$$

$$384 \quad [0]3 = C([0]5, 1) \not\rightarrow C([1]6, 1) = [1]1.$$

386 This would lead to incorrect PDE solution if the used discretization method makes
 387 use of the edge.

388 In order to satisfy this requirement, and additional synchronization of the interface
 389 cones must be carried out. We start by synchronization of the interface cone point
 390 numbering. Let us remind that `pointSF` is a one-sided structure, so only the origins
 391 of the arrows can be found directly.

392 Let us assume rank r_0 , its edge/face $[r_0]p$, and that there is a `pointSF` ar-
 393 row pointing *from* $[r_0]p$ to some $[r_1]p$ If we detect an arrow directed *from* the cone

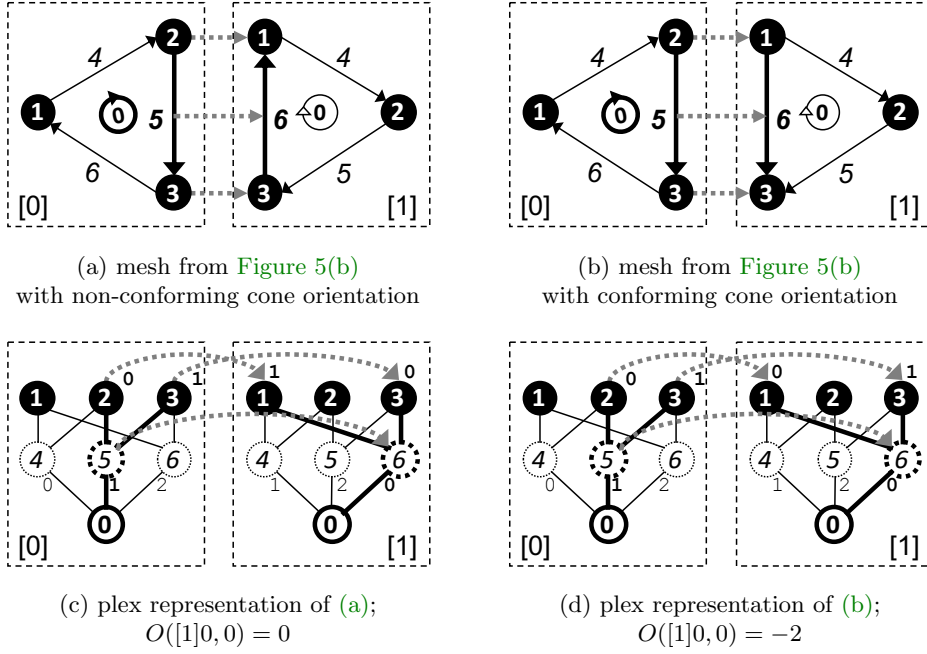


Fig. 6: Parallel DMPlex with cone points order and orientation.

394 point $C([r_0]p, c) \rightarrow [r_1]q_c$, we set $\text{root}([r_0]p, c) = (r_1, q_c)$, otherwise $\text{root}([r_0]p, c) =$
 395 $(r_0, C([r_0]p, c))$. This $\text{root}([r_0]p, c)$ is sent to r_1 using `PetscSFBCastBegin/End()`,
 396 and stored at the destination rank as $\text{leaf}([r_1]p, c)$. This is done for each rank, each
 397 point in $\text{Stratum}(h)$, $h > 0$, and each $c = 0, 1$.

398 Now from the rank r_1 view, it has for $c = 0, 1$ point $[r_1]p$, $\text{root}([r_1]p, c)$ and
 399 the received $\text{leaf}([r_1]p, c)$. If $\text{root}([r_1]p, c) = \text{leaf}([r_1]p, c)$ does not hold for both
 400 $c = 0, 1$, we must rotate and/or flip the cone so that this condition gets satisfied. In
 401 that case we must also update $O(s)$ for all $s \in S([r_1]p)$ accordingly to compensate the
 402 change of cone order.

403 We can see that the orientation synchronization heavily relies on the correct
 404 `pointSF`. This is why it must be processed first.

405 **3.3.1. Redistribution.** We now have a correct distributed DMPlex instance rep-
 406 resenting all codimensions. However, the distribution is naive; it is load balanced with
 407 respect to the size of partitions but the partition shape is not optimized. It can op-
 408 tionally be further improved using a parallel partitioner to minimize the partition
 409 interfaces and hence reduce the halo communication in the subsequent application
 410 computation. PETSc offers interfaces to ParMETIS [29, 30, 31] and PT-Scotch [15].
 411 For this paper, we always used ParMETIS. We will not discuss this stage further be-
 412 cause it is not critically needed to overcome the single node memory limit (our main
 413 motivation), its cost-effectiveness depends on the application, and it was implemented
 414 in PETSc beforehand.

415 **4. Seismic wave propagation modeling.** As a representative use case and
 416 benchmarking tool for the new parallel simulation startup phase described above in

417 [section 3](#), we use an implementation of the spectral-element method (SEM).

418 Although originally developed for applications in fluid dynamics [43], continuous-
419 Galerkin SEM on hexahedral elements has emerged as the de-facto standard for
420 global-scale seismic wave simulations [44, 18, 19]. SEM is a high-order finite-element
421 method with very low dispersion and dissipation errors [18]. The choice of the Gauss-
422 Lobatto-Legendre collocation points for the interpolating Lagrange polynomials nat-
423 urally yields a diagonal mass matrix, which enables the use of explicit time stepping
424 schemes. A second-order Newmark time-stepping additionally allows us to compute
425 coupling terms along any solid-fluid interfaces without the need to solve a linear sys-
426 tem [41]. Furthermore, the tensorized structure of the finite element basis on hexahe-
427 dral elements allows for efficient computations of internal forces. These element-wise
428 operations can be formulated as dense matrix-matrix products, making the method
429 suitable for the current generation of SIMD computing architectures.

430 Salvus [1] contains a flexible implementation of SEM, separating the wave prop-
431 agation physics, the spatial discretization and finite-element shape mappings into
432 distinct and functionally orthogonal components. It uses modern C++ features to
433 ensure that and this does not affect runtime performance. It is parallelized using
434 MPI and GPU-accelerated with CUDA. PETSc DMPlex ([section 2](#)) is used for mesh
435 management so the developments of this paper can be directly applied.

436 **5. Performance results.** This section presents scalability tests of the new par-
437 allel simulation startup phase ([section 3](#)) used within seismic wave propagation mod-
438 eling ([section 4](#)).

439 **5.1. Hardware.** All benchmarks were run at Piz Daint, the flagship system of
440 the Swiss National Supercomputing Centre (CSCS). Piz Daint consists of 5704 12-core
441 Cray XC50 nodes with GPU accelerators, and 1813 36-core Cray XC40 nodes without
442 accelerators. All benchmarks presented in this paper ran on the XC50 nodes. Each
443 of them is equipped with one 12-core Intel Xeon E5-2690 v3 (Haswell) processor,
444 one NVIDIA Tesla P100 16GB GPGPU accelerator and 64 GB RAM. Piz Daint
445 has 8.8 PB shared scratch storage with the peak performance of 112 GiB/s. It is
446 implemented using Cray Sonexion 3000 [16] scale-out storage system equipped with
447 the Lustre parallel file system [42], 40 object storage targets (OSTs) and 2 metadata
448 servers (MDSs).

449 **5.2. Middleware (Lustre, MPI-IO, HDF5) settings.** We always used the
450 single shared file approach, i.e., every process reads its disjoint chunk from the common
451 file. There are many good reasons for such choice, such as reduction of metadata
452 accesses, but the main reason is the flexibility in number of processes using the same
453 file, unlike the file-per-process approach. As for Lustre file system settings, we used
454 stripe count of 40 (maximum on Piz Daint) and stripe size of 16 MB. Regarding
455 HDF5/MPI-IO, we always non-collective reading. We tested also collective reading
456 with various numbers of aggregators (MPI-IO hint `cb_nodes`) but never saw any
457 significant benefit. The raw file reading always took less than 2 seconds; we cannot
458 exclude that some scenarios with much bigger files and/or node counts could require
459 more deliberate settings but such scenarios are irrelevant within the context of this
460 paper.

461 **5.3. Cube benchmark.** Our performance benchmark consists in elastic wave
462 propagation in homogeneous isotropic media from a point source in a cubic geometry.
463 The cube is discretized into equally sized hexahedral cells, handled as an unstructured
464 mesh. Each cell hosts a 4-th order spectral element with 125 spatial DOFs. Since

465 a 3-D vector equation was solved, this resulted in 1125 field variables per element,
 466 together representing acceleration, velocity, and displacement. [Figure 7](#) illustrates the
 467 solution of the benchmark problem at two different timesteps. [Table 1](#) summarizes
 468 dimensions of the stored topology and geometry datasets and resulting file sizes for
 469 different numbers of elements in x-direction (NEX). Our sequence of NEX was chosen
 470 so that the total number of elements (NE) of each successive mesh is approximately
 471 doubled, starting at 8 million.

472 We present performance of the new parallel startup phase in several graphs.
 473 Graphs in [Figure 8](#) present strong scalability for the mesh size of 16 million elements,
 474 and serve mainly for comparison of the new parallel startup with the original serial
 475 startup. 16 million is an upper bound for the mesh size for the serial startup imposed
 476 by the memory of a single Piz Daint Cray XC50 node. Overcoming this limit is for
 477 us the most important achievement of the startup phase parallelization. However,
 478 obviously the performance improvement is very significant as well. At 1024 nodes,
 479 the serial startup takes an amount of wall time equivalent to 202'886 timesteps in the
 480 subsequent time loop, whereas the parallel startup takes only 1827 timesteps, which
 481 means 111x speedup. The number of timesteps in production simulations varies, but
 482 generally 30'000 or more are required.

483 Graphs in [Figure 9](#) are similar to [Figure 8](#) but show the only the parallel startup
 484 scalability for various mesh sizes, so that the displayed time can be limited to 70
 485 seconds. These graph gather all stages in a single graph and different mesh sizes
 486 are presented separately. By contrast, graphs in [Figure 10](#) show strong and weak
 487 scalability for each stage separately, gathering all mesh sizes in a single graph. We can
 488 see that the topological interpolation scales almost perfectly and becomes insignificant
 489 for high number of nodes even for very large meshes. The other stages do not scale
 490 that well; however, their absolute wall times are rather small for the mesh sizes and
 491 node counts of interest. The scalability of the significant redistribution stage breaks
 492 at about 256 nodes and the mesh size no more dictates the wall time. ParMetis
 493 was used with default settings; there might be some space for slight improvement by
 494 tuning its parameters but in general it is well known that current graph partitioners
 495 do not scale beyond 10'000 cores. We can hardly do anything about it apart from
 496 perhaps testing alternative approaches such as space-filling curves. There might be
 497 some space for improvement left for the distributed plex construction; nevertheless,
 498 from the stagnation between 256 and 1024 nodes for the largest mesh we conclude
 499 that such optimization is probably not worth the effort, at least for now.

500 **6. Application: seismicity on Mars.** In late 2018, the NASA InSight mission
 501 [\[6\]](#) placed a highly sensitive seismometer [\[39\]](#) on Mars' surface and recorded the first
 502 seismic signals ever observed on Mars [\[21\]](#). The observation of seismic waves is a
 503 crucial source of information to investigate the interior structure and composition
 504 of Mars. However, as the data shows significant differences to seismic data from
 505 both Moon and Earth, numerical simulations of seismic wave propagation on Mars
 506 that account for topography as well as 3D scattering due to lateral variations of the
 507 material parameters are key to assist the interpretation of the observational data.

508 Full waveform simulations are essential to constrain the planet's structure using
 509 data in the frequency band recorded by the probe. The seismic response to marsquakes
 510 or asteroid impacts is governed by a coupled system of the elastic/acoustic wave
 511 equation, which models seismic waves propagating through Mars' mantle and the
 512 liquid core, respectively. This can be simulated efficiently using the Spectral Element
 513 Method (SEM).

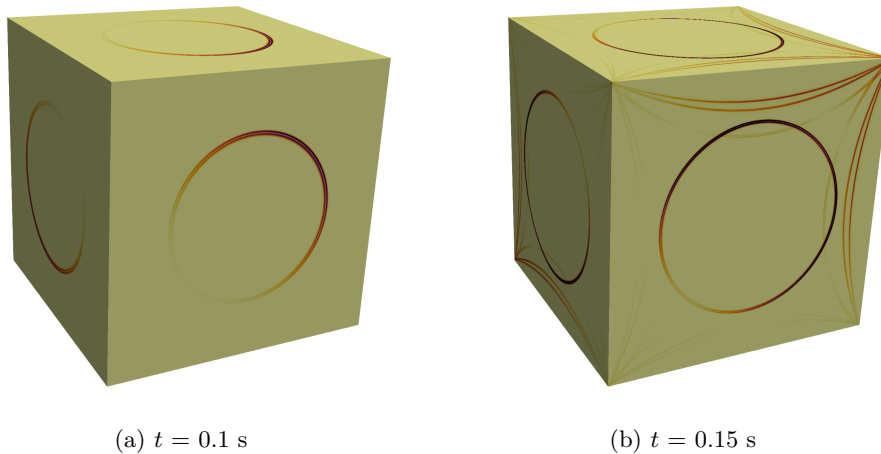


Fig. 7: Cube benchmark with a moment-tensor type source located at the center of the domain. Isotropic elastic material model without attenuation was applied. The mesh size was 256M elements, and 512 Piz Daint nodes (6144 CPUs, 512 GPUs) were used for the computation. The normalized magnitude of the velocity vector at time t is visualized. We can see the P-wave propagating from the source in (a), while (b) shows the S-wave and reflections of the P-wave from the cube boundary.

NEX	topology (int64)		geometry (double)		file size (GB)
	rows = $NE = (NEX)^3$	columns	rows	columns	
200	8'000'000	8	8'120'601	3	0.66
252	16'003'008	8	16'194'277	3	1.32
318	32'157'432	8	32'461'759	3	2.64
400	64'000'000	8	64'481'201	3	5.26
504	128'024'064	8	128'787'625	3	10.51
635	256'047'875	8	257'259'456	3	21.01

Table 1: Cube benchmark datafiles: number of elements in x-direction (NEX); total number of elements (NE); topology (connectivity) and geometry (vertices) dataset sizes; file sizes. Both topology (integer numbers) and geometry (real numbers) are stored with 64-bit precision.

514 For the computation of the seismic response of Mars, we rely on Salvus' imple-
 515 mentation of the SEM (section 4). Salvus' internal mesher uses custom algorithms
 516 to generate fully unstructured conforming 3D hexahedral meshes [49], efficiently rep-
 517 resenting topography and the extreme crustal thickness variations of Mars (~ 5 –120
 518 km), see Figure 11. The solver then represents these meshes in memory using DMPlex
 519 (section 2).

520 The CFL condition for Salvus' explicit second-order Newmark time-stepping
 521 scheme, coupled with a required minimum number of points-per-wavelength in each

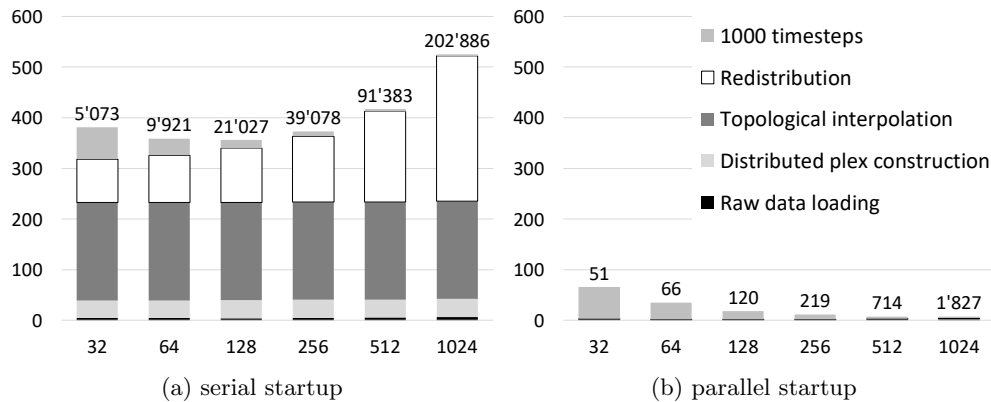


Fig. 8: **Cube benchmark. Strong scalability, serial/parallel startup, 16 million mesh elements.** The serial and parallel simulation startup phase are compared with each other and 1000 steps of the Salvus time loop in terms of wall time. Approximate wall time for any different number of timesteps can be obtained using simple proportionality. The number of timesteps in production simulations varies, but generally 30'000 or more are required. The particular startup stages are described in [section 3](#). **X-axis:** number of Piz Daint nodes, each equipped with 12 cores and 1 GPU per node. **Y-axis:** wall time in seconds. **Labels above bars:** the number of timesteps that take the same wall time as the startup phase. **Order of colors** in the bars is the same as in the legend.

522 dimension, results in the computational complexity of a simulation scaling with fre-
 523 quency to the power of 4. When using 4-th order spectral elements, which is common
 524 for planetary-scale wave propagation, more than 6 grid points per shortest wave-
 525 length are needed to accurately resolve seismic waves [19]. As the quakes are small in
 526 magnitude and the noise level increases at low frequencies, large-scale simulations are
 527 required to reach the parameter regime of the observations, and the required number
 528 of spectral elements can easily reach hundreds of millions.

529 Such mesh sizes necessitate the parallel simulation startup presented in [section 3](#).
 530 Prior to these developments, the largest possible mesh size was limited by the available
 531 memory of a single Piz Daint node to approximately 16 million elements. Moreover,
 532 loading a mesh of such size took more than four minutes. With the parallel startup
 533 in hand, these limitations vanish.

534 [Figure 11\(d\)](#) shows a snapshot of the surface displacement resulting from a simu-
 535 lation of a hypothetical quake on Mars. Here, the discretized coupled elastic wave
 536 equation has approximately 124 million 4-th order spectral elements, and we compute
 537 100'455 timesteps representing a simulated time of 30 minutes. Each element has 125
 538 spatial DOFs, each hosting 9 dynamic field components (vector displacement, veloc-
 539 ity, and acceleration). These parameters lead to an unprecedented resolved period
 540 of 3.2 s. Using again all 12 cores per Piz Daint node as well as the attached Tesla
 541 P100 GPU, this simulation took approximately 2.4 hours of wall time on 256 Piz
 542 Daint nodes. From this total wall time, raw data loading took 0.9 s, distributed plex
 543 construction 4.5 s, topological interpolation 3.1 s and redistribution 2.3 s, i.e., the
 544 whole startup phase took less than 11 s.

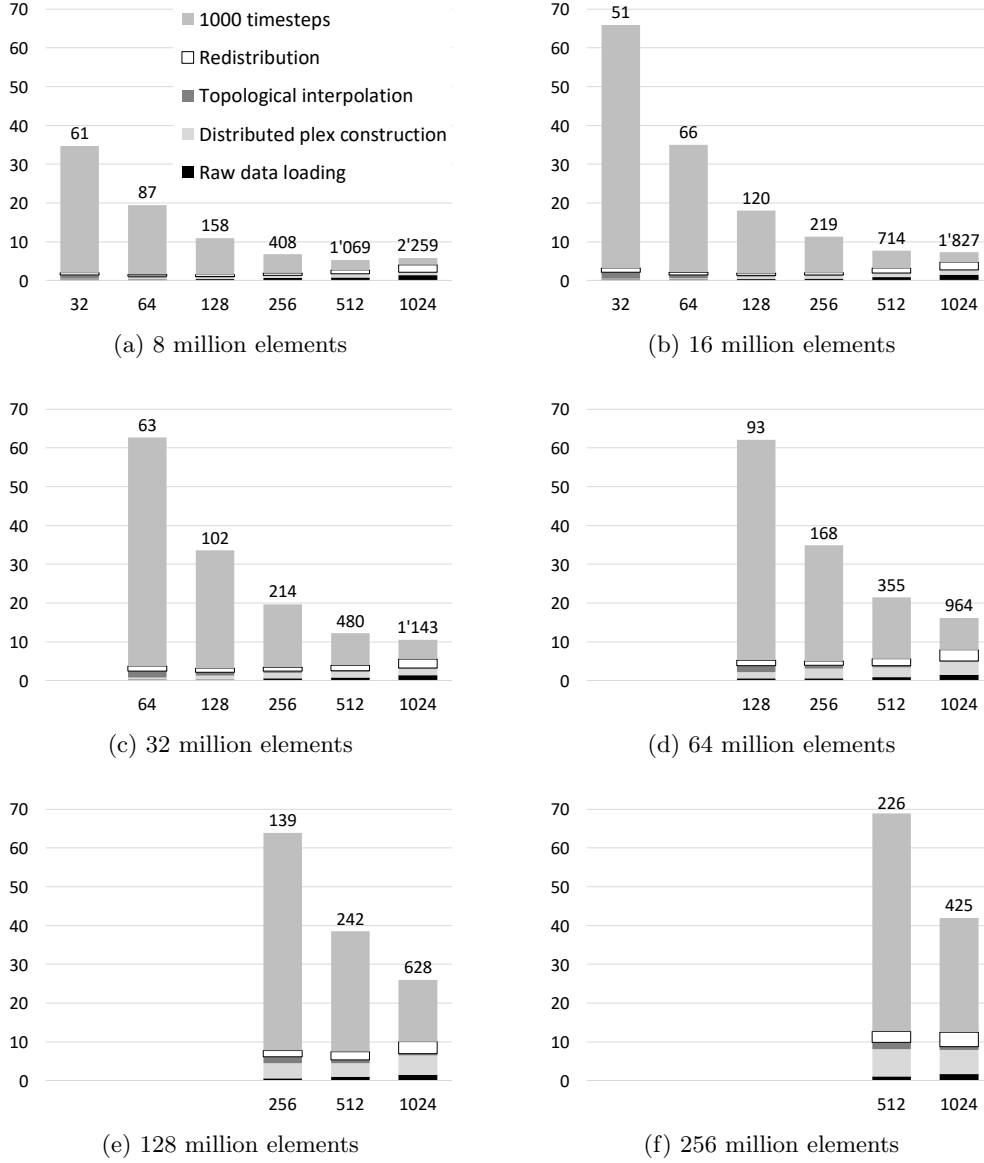


Fig. 9: **Cube benchmark. Strong scalability, parallel startup, various mesh sizes.** The parallel simulation startup phase is compared with 1000 steps of the Salvus time loop in terms of wall time. Approximate wall time for any number of timesteps can be obtained using simple proportionality. The number of timesteps in production simulations varies, but generally 30'000 or more are required. The particular stages are described in [section 3](#). **X-axis:** number of Piz Daint nodes, each with 12 cores and 1 GPU. **Y-axis:** wall time in seconds. **Labels above bars:** the number of timesteps that take the same wall time as the startup phase. **Missing bars:** out of memory failure during the time loop caused by the memory limit of the GPUs. **Order of colors** in the bars is the same as in the legend. **Note:** (a) is the same as [Figure 8\(b\)](#) but with a re-scaled time axis.

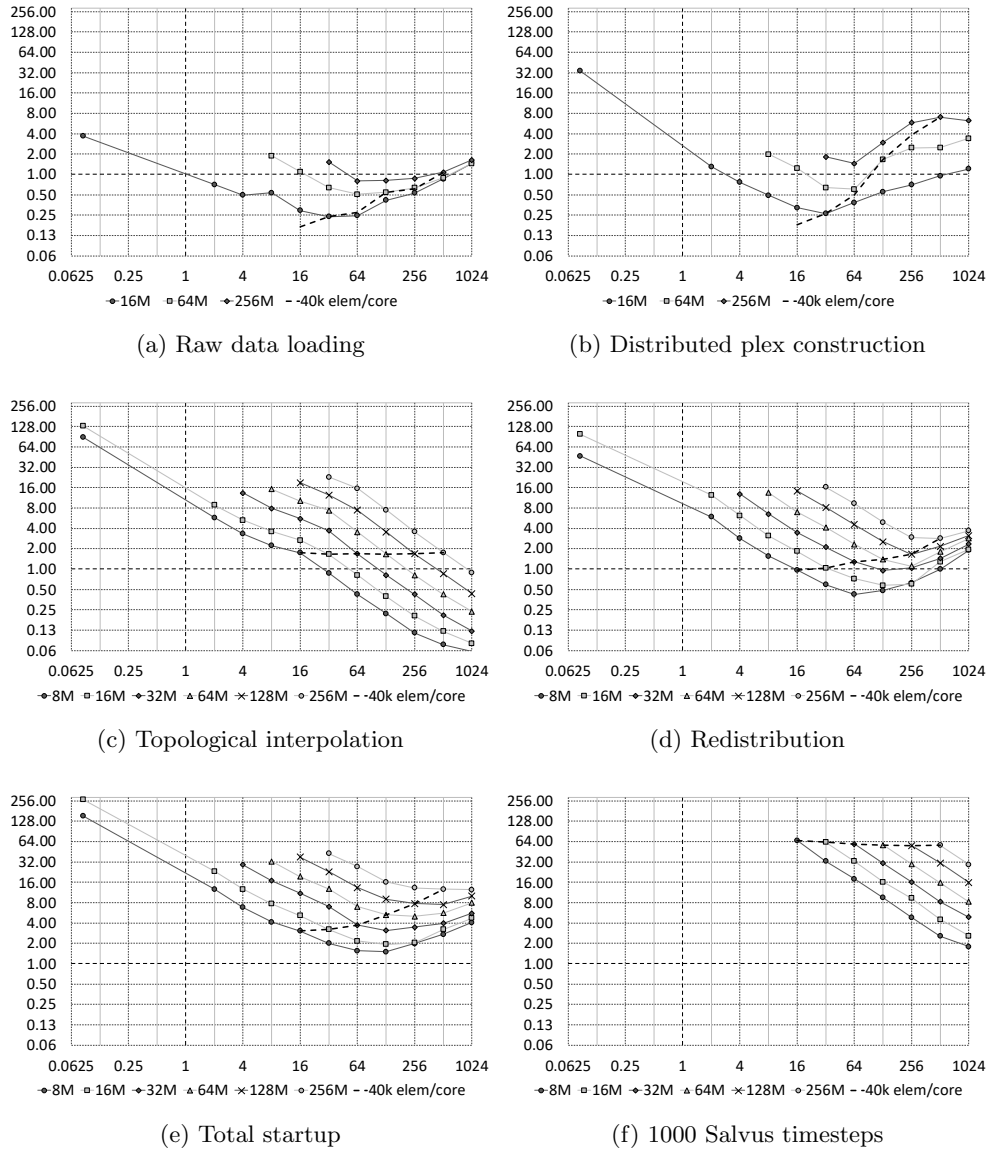
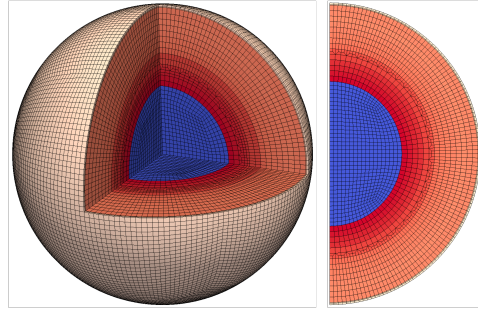
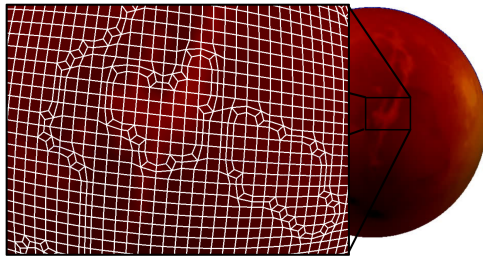


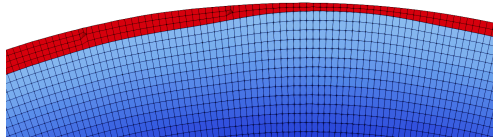
Fig. 10: **Cube benchmark. Strong and weak scalability per stage.** **X-axis:** number of Piz Daint nodes (log2 scale), each with 12 cores and 1 GPU; 0.0625 = 1/12 means a single core (serial) run. Note that a single node run was not possible due to the out-of-memory failure during the redistribution phase. **Y-axis:** wall time in seconds (log2 scale). **Solid lines** show the strong scalability for different mesh sizes. **Dashed line** shows the weak scalability for the mesh size of approximately 40'000 elements per core which is the upper bound for the Salvus timeloop imposed by limited GPU memory. **Order of line styles** is the same in the plots and in the legends.



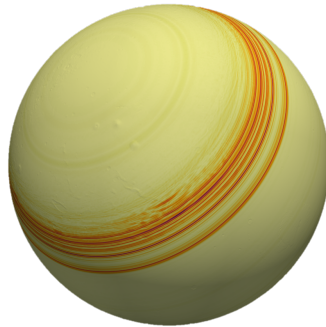
(a) cubed sphere with fluid core (blue) and solid mantle (red)



(b) modeling 3D topography on the surface



(c) crustal thickness variations (red)



(d) snapshot of the seismic wavefield on the surface (result of the simulation)

Fig. 11: Anisotropic mesh refinements to accurately model the structure of Mars on a conforming hexahedral mesh. The Mars texture map is based on NASA elevation and imagery data.

545 **7. Conclusions.** We presented algorithmic strategies for handling unstructured
 546 meshes for high-order finite-element simulations on complex domains. In particular,
 547 we demonstrated new capabilities for parallel mesh reading and topological interpo-
 548 lation in PETSc DMplex, which enables a fully parallel workflow starting from the
 549 initial data file reading. This is beneficial not only for direct users of DMplex but
 550 also users of software libraries and packages employing DMplex, such as Firedrake
 551 [45], Salvus [1], or HPDDM [27].

552 This work in a sense follows up [35] and addresses the main task stated in their
 553 Future Work: “*Most crucially perhaps is the development of a fully parallel mesh input*
 554 *reader in PETSc in order to overcome the remaining sequential bottleneck during*
 555 *model initialisation.*” Moreover, that paper mentions the “HDF5-based XDMF
 556 output format” but it has become an *input* format as well within this work. Hence,
 557 HDF5/XDMF has become the first widely used mesh format supported by PETSc
 558 which is both readable and writable in parallel.

559 The implementation is agnostic to the type of finite elements in the mesh and
 560 completely decoupled from the governing equations, and is thus applicable in many
 561 scientific disciplines. In particular, our solution overcomes bottlenecks in numerical
 562 modeling of seismic wave propagation on Mars and shows excellent parallel scalability
 563 in large-scale simulations on more than 12'000 cores.

564 **Acknowledgments.** All presented PETSc developments have been made pub-
 565 licly available in PETSc since its release 3.13.

566 We gratefully acknowledge support from the Swiss National Supercomputing Cen-
 567 tre (CSCS) under projects s922 and s961; the Platform for Advanced Scientific Com-
 568 puting (PASC) under the project “Salvus”; the Swiss National Science Foundation
 569 (SNF) BRIDGE fellowship program under the grant agreement No. 175322; the Eu-
 570 ropean Research Council (ERC) from the EU’s Horizon 2020 programme under grant
 571 agreement No. 714069; and the ETH Zurich Postdoctoral Fellowship Program which
 572 in turn received funding from the EU’s Seventh Framework Programme under the
 573 grant agreement No. 608881.

574 REFERENCES

- 575 [1] M. AFANASIEV, C. BOEHM, M. VAN DRIEL, L. KRISCHER, M. RIETMANN, D. A. MAY, M. G.
 576 KNEPLEY, AND A. FICHTNER, *Modular and flexible spectral-element waveform modelling in*
 577 *two and three dimensions*, Geophysical Journal International, 216 (2019), pp. 1675–1692,
 578 <https://doi.org/10.1093/gji/ggy469>.
 579 [2] R. ANDERSON, J. ANDREJ, A. BARKER, J. BRAMWELL, J.-S. CAMIER, J. CERVENY, V. DOBREV,
 580 Y. DUDOUT, A. FISHER, T. KOLEV, W. PAZNER, M. STOWELL, V. TOMOV, J. DAHM,
 581 D. MEDINA, AND S. ZAMPINI, *Mfem: a modular finite element methods library*, 2019,
 582 <https://arxiv.org/abs/1911.09220>.
 583 [3] S. BALAY, S. ABHYANKAR, M. F. ADAMS, J. BROWN, P. BRUNE, K. BUSCHELMAN, L. DALCIN,
 584 A. DENER, V. ELJKHOUT, W. D. GROPP, D. KARPEYEV, D. KAUSHIK, M. G. KNEPLEY,
 585 D. A. MAY, L. C. MCINNES, R. T. MILLS, T. MUNSON, K. RUPP, P. SANAN, B. F. SMITH,
 586 S. ZAMPINI, H. ZHANG, AND H. ZHANG, *PETSc web page*, <https://www.mcs.anl.gov/petsc>
 587 (accessed 2020-03-30).
 588 [4] S. BALAY, S. ABHYANKAR, M. F. ADAMS, J. BROWN, P. BRUNE, K. BUSCHELMAN, L. DALCIN,
 589 A. DENER, V. ELJKHOUT, W. D. GROPP, D. KARPEYEV, D. KAUSHIK, M. G. KNEPLEY,
 590 D. A. MAY, L. C. MCINNES, R. T. MILLS, T. MUNSON, K. RUPP, P. SANAN,
 591 B. F. SMITH, S. ZAMPINI, H. ZHANG, AND H. ZHANG, *PETSc users manual*, Tech. Re-
 592 port ANL-95/11 - Revision 3.13, Argonne National Laboratory, 2020, [https://www.mcs.](https://www.mcs.anl.gov/petsc/petsc-current/docs/manual.pdf)
 593 [anl.gov/petsc/petsc-current/docs/manual.pdf](https://www.mcs.anl.gov/petsc/petsc-current/docs/manual.pdf) (accessed 2020-03-30).
 594 [5] S. BALAY, W. D. GROPP, L. C. MCINNES, AND B. F. SMITH, *Efficient management of parallel-*
 595 *ism in object oriented numerical software libraries*, in Modern Software Tools in Scientific

- 596 Computing, E. Arge, A. M. Bruaset, and H. P. Langtangen, eds., Birkhäuser Press, 1997,
597 pp. 163–202, https://doi.org/10.1007/978-1-4612-1986-6_8.
- [6] W. B. BANERDT, S. E. SMREKAR, D. BANFIELD, D. GIARDINI, M. GOLOMBEK, C. L. JOHN-
598 SON, P. LOGNONNÉ, A. SPIGA, T. SPOHN, C. PERRIN, S. C. STÄHLER, D. ANTONANGELI,
599 S. ASMAR, C. BEGHEIN, N. BOWLES, E. BOZDAG, P. CHI, U. CHRISTENSEN, J. CLINTON,
600 G. S. COLLINS, I. DAUBAR, V. DEHANT, M. DRILLEAU, M. FILLINGIM, W. FOLKNER, R. F.
601 GARCIA, J. GARVIN, J. GRANT, M. GROTT, J. GRYGORCZUK, T. HUDSON, J. C. E. IRVING,
602 G. KARGL, T. KAWAMURA, S. KEDAR, S. KING, B. KNAPMEYER-ENDRUN, M. KNAPMEYER,
603 M. LEMMON, R. LORENZ, J. N. MAKI, L. MARGERIN, S. M. MCLENNAN, C. MICHAUD,
604 D. MIMOUN, A. MITTELHOLZ, A. MOCQUET, P. MORGAN, N. T. MUELLER, N. MURDOCH,
605 S. NAGIHARA, C. NEWMAN, F. NIMMO, M. PANNING, W. T. PIKE, A.-C. PLESA, S. RO-
606 DRIGUEZ, J. A. RODRIGUEZ-MANFREDI, C. T. RUSSELL, N. SCHMERR, M. SIEGLER, S. STAN-
607 LEY, E. STUTZMANN, N. TEANBY, J. TROMP, M. VAN DRIEL, N. WARNER, R. WEBER, AND
608 M. WIECZOREK, *Initial results from the InSight mission on Mars*, Nat. Geosci., 13 (2020),
609 pp. 183–189, <https://doi.org/10.1038/s41561-020-0544-y>.
- [7] N. BARRAL, M. G. KNEPLEY, M. LANGE, M. D. PIGGOTT, AND G. J. GORMAN, *Anisotropic*
612 *mesh adaptation in firedrake with petsc dmplex*, 2016, <https://arxiv.org/abs/1610.09874>.
- [8] J. BROWN, *Star forests as a parallel communication model*, 2011, [https://jedbrown.org/files/](https://jedbrown.org/files/StarForest.pdf)
613 [StarForest.pdf](https://jedbrown.org/files/StarForest.pdf) (accessed 2020-03-17).
- [9] J. BROWN, M. G. KNEPLEY, D. A. MAY, L. C. MCINNES, AND B. F. SMITH, *Composable linear*
614 *solvers for multiphysics*, in Proceedings of the 11th International Symposium on Parallel
615 and Distributed Computing (ISPDC 2012), IEEE Computer Society, 2012, pp. 55–62,
616 <https://doi.org/10.1109/ISPDC.2012.16>.
- [10] J. BROWN, M. G. KNEPLEY, AND B. F. SMITH, *Run-time extensibility and librarization of*
617 *simulation software*, Computing in Science Engineering, 17 (2015), pp. 38–45, [https://doi.](https://doi.org/10.1109/MCSE.2014.95)
618 [org/10.1109/MCSE.2014.95](https://doi.org/10.1109/MCSE.2014.95).
- [11] P. R. BRUNE, M. G. KNEPLEY, AND L. R. SCOTT, *Unstructured geometric multigrid in two and*
619 *three dimensions on complex and graded meshes*, SIAM Journal on Scientific Computing, 35
620 (2013), pp. A173–A191, <https://doi.org/10.1137/110827077>, [https://arxiv.org/abs/1104.](https://arxiv.org/abs/1104.0261)
621 [0261](https://arxiv.org/abs/1104.0261).
- [12] C. BURSTEDDE, O. GHATTAS, M. GURNIS, T. ISAAC, G. STADLER, T. WARBURTON, AND
622 L. WILCOX, *Extreme-scale AMR*, in SC '10: Proceedings of the 2010 ACM/IEEE Inter-
623 national Conference for High Performance Computing, Networking, Storage and Analysis,
624 2010, pp. 1–12, <https://doi.org/10.1109/SC.2010.25>.
- [13] C. BURSTEDDE, L. C. WILCOX, AND O. GHATTAS, *p4est: Scalable algorithms for parallel adap-*
625 *tive mesh refinement on forests of octrees*, SIAM Journal on Scientific Computing, 33
626 (2011), pp. 1103–1133, <https://doi.org/10.1137/100791634>.
- [14] H. BÉRIOT, A. PRINN, AND G. GABARD, *Efficient implementation of high-order finite elements*
627 *for Helmholtz problems*, International Journal for Numerical Methods in Engineering, 106
628 (2016), pp. 213–240, <https://doi.org/10.1002/nme.5172>.
- [15] C. CHEVALIER AND F. PELLEGRINI, *PT-Scotch: A tool for efficient parallel graph ordering*,
629 *Parallel Computing*, 34 (2008), pp. 318–331, <https://doi.org/10.1016/j.parco.2007.12.001>.
- [16] CRAY INC., *Cray Sonexion 3000 Storage System*, 2016, [https://www.cray.com/sites/default/](https://www.cray.com/sites/default/files/resources/WP-Cray-Sonexion-3000-Storage-Systems.pdf)
630 [files/resources/WP-Cray-Sonexion-3000-Storage-Systems.pdf](https://www.cray.com/sites/default/files/resources/WP-Cray-Sonexion-3000-Storage-Systems.pdf) (accessed 2020-03-20).
- [17] P. E. FARRELL, L. MITCHELL, AND F. WECHSUNG, *An augmented Lagrangian preconditioner*
631 *for the 3D stationary incompressible Navier–Stokes equations at high Reynolds number*,
632 *SIAM Journal on Scientific Computing*, 41 (2019), pp. A3073–A3096, [https://doi.org/10.](https://doi.org/10.1137/18M1219370)
633 [1137/18M1219370](https://doi.org/10.1137/18M1219370).
- [18] A. FERRONI, P. ANTONIETTI, I. MAZZIERI, AND A. QUARTERONI, *Dispersion-dissipation analy-*
634 *sis of 3-D continuous and discontinuous spectral element methods for the elastodynamics*
635 *equation*, Geophys. J. Int., 211 (2017), pp. 1554–1574, <https://doi.org/10.1093/gji/ggx384>.
- [19] A. FICHTNER, *Full Seismic Waveform Modelling and Inversion*, Advances in Geophysical and
636 *Environmental Mechanics and Mathematics*, Springer Berlin Heidelberg, 2011, [https://doi.](https://doi.org/10.1007/978-3-642-15807-0)
637 [org/10.1007/978-3-642-15807-0](https://doi.org/10.1007/978-3-642-15807-0).
- [20] J. E. FLAHERTY, R. M. LOY, M. S. SHEPHARD, B. K. SZYMANSKI, J. D. TERESCO, AND L. H.
638 ZIANTZ, *Adaptive local refinement with octree load balancing for the parallel solution of*
639 *three-dimensional conservation laws*, Journal of Parallel and Distributed Computing, 47
640 (1997), pp. 139–152, <https://doi.org/10.1006/jpdc.1997.1412>.
- [21] D. GIARDINI, P. LOGNONNÉ, W. B. BANERDT, W. T. PIKE, U. CHRISTENSEN, S. CEYLAN,
641 J. F. CLINTON, M. VAN DRIEL, S. C. STÄHLER, M. BÖSE, R. F. GARCIA, A. KHAN,
642 M. PANNING, C. PERRIN, D. BANFIELD, E. BEUCLER, C. CHARALAMBOUS, F. EUCH-
643 NER, A. HORLESTON, A. JACOB, T. KAWAMURA, S. KEDAR, G. MAINSANT, J.-R. SCHOLZ,

- 658 S. E. SMREKAR, A. SPIGA, C. AGARD, D. ANTONANGELI, S. BARKAOU, E. BARRETT,
 659 P. COMBES, V. CONEJERO, I. DAUBAR, M. DRILLEAU, C. FERRIER, T. GABSI, T. GUD-
 660 KOVA, K. HURST, F. KARAKOSTAS, S. KING, M. KNAPMEYER, B. KNAPMEYER-ENDRUN,
 661 R. LLORCA-CEJUDO, A. LUCAS, L. LUNO, L. MARGERIN, J. B. MCCLEAN, D. MIMOUN,
 662 N. MURDOCH, F. NIMMO, M. NONON, C. PARDO, A. RIVOLDINI, J. A. R. MANFREDI,
 663 H. SAMUEL, M. SCHIMMEL, A. E. STOTT, E. STUTZMANN, N. TEANBY, T. WARREN, R. C.
 664 WEBER, M. WIECZOREK, AND C. YANA, *The seismicity of Mars*, Nat. Geosci., 13 (2020),
 665 pp. 205–212, <https://doi.org/10.1038/s41561-020-0539-8>.
- [22] HDF5 GROUP, *HDF5*, <https://portal.hdfgroup.org/display/HDF5/HDF5> (accessed 2020-03-
 666 20).
 667
- [23] HDF5 GROUP, *Parallel HDF5*, <https://portal.hdfgroup.org/display/HDF5/Parallel+HDF5>
 668 (accessed 2020-03-20).
 669
- [24] T. HUGHES, *The Finite Element Method: Linear Static and Dynamic Finite Element Analysis*,
 670 vol. 78, Dover Publications, 2000.
 671
- [25] T. ISAAC, C. BURSTEDDE, L. C. WILCOX, AND O. GHATTAS, *Recursive algorithms for distributed*
 672 *forests of octrees*, SIAM Journal on Scientific Computing, 37 (2015), pp. C497–C531, <https://doi.org/10.1137/140970963>.
 673
 674
- [26] T. ISAAC AND M. G. KNEPLEY, *Support for non-conformal meshes in PETSc's DMplex inter-*
 675 *face*, 2015, <https://arxiv.org/abs/1508.02470>.
 676
- [27] P. JOLIVET, F. NATAF, ET AL., *HPDDM – high-performance unified framework for domain*
 677 *decomposition methods*, <https://github.com/hpddm/hpddm> (accessed 2020-04-18).
 678
- [28] G. KARYPIS ET AL., *METIS - serial graph partitioning and fill-reducing matrix ordering*, <http://glaros.dtc.umn.edu/gkhome/metis/metis/overview> (accessed 2020-03-17).
 679
 680
- [29] G. KARYPIS ET AL., *ParMETIS - parallel graph partitioning and fill-reducing matrix ordering*,
 681 <http://glaros.dtc.umn.edu/gkhome/metis/parmetis/overview> (accessed 2020-03-17).
 682
- [30] G. KARYPIS AND V. KUMAR, *A fast and high quality multilevel scheme for partitioning irregular*
 683 *graphs*, SIAM Journal on Scientific Computing, 20 (1998), pp. 359–392, [https://doi.org/](https://doi.org/10.1137/S1064827595287997)
 684 [10.1137/S1064827595287997](https://doi.org/10.1137/S1064827595287997).
 685
- [31] G. KARYPIS AND V. KUMAR, *A parallel algorithm for multilevel graph partitioning and sparse*
 686 *matrix ordering*, Journal of Parallel and Distributed Computing, 48 (1998), pp. 71–95,
 687 <https://doi.org/10.1006/jpdc.1997.1403>.
 688
- [32] M. G. KNEPLEY AND D. A. KARPEEV, *Mesh algorithms for PDE with Sieve I: Mesh distribution*,
 689 *Scientific Programming*, 17 (2009), pp. 215–230, <https://doi.org/10.3233/SPR-2009-0249>.
 690
- [33] M. G. KNEPLEY, M. LANGE, AND G. J. GORMAN, *Unstructured overlapping mesh distribution*
 691 *in parallel*, 2017, <https://arxiv.org/abs/1506.06194>.
 692
- [34] T. KÄRNÄ, S. C. KRAMER, L. MITCHELL, D. A. HAM, M. D. PIGGOTT, AND A. M. BAP-
 693 TISTA, *Thetis coastal ocean model: discontinuous Galerkin discretization for the three-*
 694 *dimensional hydrostatic equations*, Geoscientific Model Development, 11 (2018), pp. 4359–
 695 4382, <https://doi.org/10.5194/gmd-11-4359-2018>.
 696
- [35] M. LANGE, M. G. KNEPLEY, AND G. J. GORMAN, *Flexible, scalable mesh and data manage-*
 697 *ment using PETSc DMplex*, in Proceedings of the 3rd International Conference on Exascale
 698 Applications and Software, EASC '15, Edinburgh, Scotland, UK, 2015, University of Edin-
 699 burgh, pp. 71–76, <http://dl.acm.org/citation.cfm?id=2820083.2820097> (accessed 2020-03-
 700 30).
 701
- [36] M. LANGE, L. MITCHELL, M. G. KNEPLEY, AND G. J. GORMAN, *Efficient mesh management*
 702 *in Firedrake using PETSc DMplex*, SIAM Journal on Scientific Computing, 38 (2016),
 703 pp. S143–S155, <https://doi.org/10.1137/15M1026092>.
 704
- [37] M. G. LARSON AND F. BENGZON, *The finite element method: theory, implementation, and*
 705 *applications*, no. 10 in Texts in computational science and engineering, Springer, 2013,
 706 <https://doi.org/10.1007/978-3-642-33287-6>.
 707
- [38] A. LOGG, *Efficient representation of computational meshes*, International Journal of Compu-
 708 tational Science and Engineering, 4 (2009), pp. 283–295, [https://doi.org/10.1504/IJCSE.](https://doi.org/10.1504/IJCSE.2009.029164)
 709 [2009.029164](https://doi.org/10.1504/IJCSE.2009.029164).
 710
- [39] P. H. LOGNONNÉ, W. B. BANERDT, D. GIARDINI, W. T. PIKE, U. CHRISTENSEN, P. LAUDET,
 711 S. DE RAUCOURT, P. ZWEIFEL, S. CALCUTT, M. BIERWIRTH, K. J. HURST, F. IPELAAN,
 712 J. W. UMLAND, R. LLORCA-CEJUDO, S. A. LARSON, R. F. GARCIA, S. KEDAR,
 713 B. KNAPMEYER-ENDRUN, D. MIMOUN, A. MOCQUET, M. P. PANNING, R. C. WEBER,
 714 A. SYLVESTRE-BARON, G. PONT, N. VERDIER, L. KERJEAN, L. J. FACTO, V. GHARAKA-
 715 NIAN, J. E. FELDMAN, T. L. HOFFMAN, D. B. KLEIN, K. KLEIN, N. P. ONUFER, J. PAREDES-
 716 GARCIA, M. P. PETKOV, J. R. WILLIS, S. E. SMREKAR, M. DRILLEAU, T. GABSI, T. NEBUT,
 717 O. ROBERT, S. TILLIER, C. MOREAU, M. PARISE, G. AVENI, S. BEN CHAREF, Y. BENNOUR,
 718 T. CAMUS, P. A. DANDONNEAU, C. DESFOUX, B. LECOMTE, O. POT, P. REVUZ, D. MANCE,
 719

- 720 J. TENPIERICK, N. E. BOWLES, C. CHARALAMBOUS, A. K. DELAHUNTY, J. HURLEY, R. IR-
 721 SHAD, H. LIU, A. G. MUKHERJEE, I. M. STANDLEY, A. E. STOTT, J. TEMPLE, T. WARREN,
 722 M. EBERHARDT, A. KRAMER, W. KÜHNE, E.-P. MIETTINEN, M. MONECKE, C. AICARDI,
 723 M. ANDRÉ, J. BAROUKH, A. BORRIEN, A. BOUISSET, P. BOUTTE, K. BRETHOMÉ, C. BRY-
 724 BAERT, T. CARLIER, M. DELEUZE, J. M. DESMARRES, D. DILHAN, C. DOUCET, D. FAYE,
 725 N. FAYE-REFALO, R. GONZALEZ, C. IMBERT, C. LARIGAUDERIE, E. LOCATELLI, L. LUNO, J.-
 726 R. MEYER, F. MIALHE, J. M. MOURET, M. NONON, Y. PAHN, A. PAILLET, P. PASQUIER,
 727 G. PEREZ, R. PEREZ, L. PERRIN, B. POUILLOUX, A. ROSAK, I. SAVIN DE LARCLAUZE,
 728 J. SICRE, M. SODKI, N. TOULEMONT, B. VELLA, C. YANA, F. ALIBAY, O. M. AVALOS,
 729 M. A. BALZER, P. BHANDARI, E. BLANCO, B. D. BONE, J. C. BOUSMAN, P. BRUNEAU,
 730 F. J. CALEF, R. J. CALVET, S. A. D'AGOSTINO, G. DE LOS SANTOS, R. G. DEEN, R. W.
 731 DENISE, J. ERVIN, N. W. FERRARO, H. E. GENGL, F. GRINBLAT, D. HERNANDEZ, M. HET-
 732 ZEL, M. E. JOHNSON, L. KHACHIKYAN, J. Y. LIN, S. M. MADZUNKOV, S. L. MARSHALL,
 733 I. G. MIKELLIDES, E. A. MILLER, W. RAFF, J. E. SINGER, C. M. SUNDAY, J. F. VIL-
 734 LALVAZO, M. C. WALLACE, D. BANFIELD, J. A. RODRIGUEZ-MANFREDI, C. T. RUSSELL,
 735 A. TREBI-OLLENU, J. N. MAKI, É. BEUCLER, M. BÖSE, C. BONJOUR, J. L. BERENGUER,
 736 S. CEYLAN, J. F. CLINTON, V. CONEJERO, I. J. DAUBAR, V. DEHANT, P. DELAGE, F. EU-
 737 CHNER, I. ESTÈVE, L. FAYON, L. FERRAIOLI, C. L. JOHNSON, J. GAGNEPAIN-BEYNEIX,
 738 M. GOLOBEK, A. KHAN, T. KAWAMURA, B. KENDA, P. LABROT, N. MURDOCH, C. PARDO,
 739 C. PERRIN, L. POU, A. SAURON, D. SAVOIE, S. C. STÄHLER, É. STUTZMANN, N. A. TEANBY,
 740 J. TROMP, M. VAN DRIEL, M. A. WIECZOREK, R. WIDMER-SCHNIDRIG, AND J. WOOKEY,
 741 *SEIS: Insight's seismic experiment for internal structure of Mars*, Space Sci. Rev., 215
 742 (2019), p. 12, <https://doi.org/10.1007/s11214-018-0574-6>.
- 743 [40] MESSAGE PASSING INTERFACE FORUM, *MPI: A Message-Passing Interface Standard. Version*
 744 *3.1*, 2015, <https://www.mpi-forum.org/docs/mpi-3.1/mpi31-report.pdf> (accessed 2020-03-
 745 20).
- 746 [41] T. NISSEN-MEYER, A. FOURNIER, AND F. A. DAHLEN, *A two-dimensional spectral-element*
 747 *method for computing spherical-earth seismograms - II. Waves in solid-fluid media.*, Geo-
 748 *phys. J. Int.*, 174 (2008), pp. 873–888.
- 749 [42] ORACLE AND INTEL CORPORATION, *Lustre Software Release 2.x - Operations Manual*, 2017,
 750 http://doc.lustre.org/lustre_manual.pdf (accessed 2020-03-20).
- 751 [43] A. T. PATERA, *A spectral element method for fluid dynamics: Laminar flow in a channel*
 752 *expansion*, Journal of Computational Physics, 54 (1984), pp. 468–488.
- 753 [44] D. PETER, D. KOMATITSCH, Y. LUO, R. MARTIN, N. LE GOFF, E. CASAROTTI, P. LE LOHER,
 754 F. MAGNONI, Q. LIU, C. BLITZ, T. NISSEN-MEYER, P. BASINI, AND J. TROMP, *Forward and*
 755 *adjoint simulations of seismic wave propagation on fully unstructured hexahedral meshes*,
 756 *Geophys. J. Int.*, 186 (2011), pp. 721–739.
- 757 [45] F. RATHGEBER, D. A. HAM, L. MITCHELL, M. LANGE, F. LUPORINI, A. T. T. MCRAE, G.-
 758 T. BERCEA, G. R. MARKALL, AND P. H. J. KELLY, *Firedrake: Automating the finite*
 759 *element method by composing abstractions*, ACM Trans. Math. Softw., 43 (2016), <https://doi.org/10.1145/2998441>.
- 760 [46] J. RUDI, A. C. I. MALOSI, T. ISAAC, G. STADLER, M. GURNIS, P. W. J. STAAR, Y. INEICHEN,
 761 C. BEKAS, A. CURIONI, AND O. GHATTAS, *An extreme-scale implicit solver for complex*
 762 *PDEs: highly heterogeneous flow in earth's mantle*, in Proceedings of the International
 763 Conference for High Performance Computing, Networking, Storage and Analysis, SC '15,
 764 Association for Computing Machinery, 2015, pp. 1–12, [https://doi.org/10.1145/2807591.](https://doi.org/10.1145/2807591.2807675)
 765 [2807675](https://doi.org/10.1145/2807675).
- 766 [47] R. SCHNEIDERS, *Octree-based hexahedral mesh generation*, International Journal of Com-
 767 putational Geometry & Applications, 10 (2000), pp. 383–398, [https://doi.org/10.1142/](https://doi.org/10.1142/S021819590000022X)
 768 [S021819590000022X](https://doi.org/10.1142/S021819590000022X).
- 769 [48] T. J. TAUTGES, C. ERNST, C. STIMPSON, R. J. MEYERS, AND K. MERKLEY, *MOAB: a mesh-*
 770 *oriented database*, Tech. Report SAND2004-1592, Sandia National Laboratories, 2004,
 771 <https://doi.org/10.2172/970174>.
- 772 [49] M. VAN DRIEL, C. BOEHM, L. KRISCHER, AND M. AFANASIEV, *Accelerating numerical wave prop-*
 773 *agation using wavefield adapted meshes. Part I: forward and adjoint modelling*, Geophysical
 774 Journal International, 221 (2020), pp. 1580–1590, <https://doi.org/10.1093/gji/ggaa058>.
- 775 [50] *XDMF - eXtensible Data Model and Format*, <http://www.xdmf.org/> (accessed 2020-03-20).
- 776 [51] O. ZIENKIEWICZ, R. TAYLOR, AND J. ZHU, eds., *The Finite Element Method: its Basis and*
 777 *Fundamentals (Seventh Edition)*, Butterworth-Heinemann, Oxford, seventh edition ed.,
 778 2013, <https://doi.org/10.1016/C2009-0-24909-9>.