



Computer Science and Artificial Intelligence Laboratory  
Technical Report

MIT-CSAIL-TR-2015-017

May 26, 2015

---

## Simit: A Language for Physical Simulation

Fredrik Kjolstad, Shoaib Kamil, Jonathan Ragan-Kelley, David I.W. Levin, Shinjiro Sueda, Desai Chen, Etienne Vouga, Danny M. Kaufman, Gurtej Kanwar, Wojciech Matusik, and Saman Amarasinghe



# Simit: A Language for Physical Simulation

FREDRIK KJOLSTAD and SHOAB KAMIL

Massachusetts Institute of Technology

JONATHAN RAGAN-KELLEY

Stanford University

DAVID I.W. LEVIN

Disney Research

SHINJIRO SUEDA

California Polytechnic State University

DESAI CHEN

Massachusetts Institute of Technology

ETIENNE VOUGA

University of Texas at Austin

DANNY M. KAUFMAN

Adobe Systems

and

GURTEJ KANWAR, WOJCIECH MATUSIK, and SAMAN AMARASINGHE

Massachusetts Institute of Technology

Using existing programming tools, writing high-performance simulation code is labor intensive and requires sacrificing readability and portability. The alternative is to prototype simulations in a high-level language like Matlab, thereby sacrificing performance. The Matlab programming model naturally describes the behavior of an entire physical system using the language of linear algebra. However, simulations also manipulate individual geometric elements, which are best represented using linked data structures like meshes. Translating between the linked data structures and linear algebra comes at significant cost, both to the programmer and the machine. High-performance implementations avoid the cost by rephrasing the computation in terms of linked or index data structures, leaving the code complicated and monolithic, often increasing its size by an order of magnitude.

In this paper, we present Simit, a new language for physical simulations that lets the programmer view the system both as a linked data structure in the form of a hypergraph, and as a set of global vectors, matrices and tensors depending on what is convenient at any given time. Simit provides a novel assembly construct that makes it conceptually easy and computationally efficient to move between the two abstractions. Using the information provided by the assembly construct, the compiler generates efficient in-place computation on the graph. We demonstrate that Simit is easy to use: a Simit program is typically shorter than a Matlab program; that it is high-performance: a Simit program running sequentially on a CPU performs comparably to hand-optimized simulations; and that it is portable: Simit programs can be compiled for GPUs with no change to the program, delivering 5-25x speedups over our optimized CPU code.

Categories and Subject Descriptors: I.3.7 [Computer Graphics]: Three-Dimensional Graphics and Realism—*Animation*

General Terms: Languages, Performance

Additional Key Words and Phrases: Graph, Matrix, Tensor, Simulation

## 1. INTRODUCTION

Efficient large-scale computer simulations of physical phenomena are notoriously difficult to engineer, requiring careful optimization to achieve good performance. This stands in stark contrast to the elegance of the underlying physical laws; for example, the behavior of an elastic object, discretized (for ease of exposition) as a network of masses connected by springs, is determined by a single quadratic equation, Hooke’s law, applied homogeneously to every spring in the network. While Hooke’s law describes the local behavior of the mass-spring network, it tells us relatively little about its global, emergent behavior. This global behavior, such as how an entire object will deform, is also described by simple but coupled systems of equations.

Each of these two aspects of the physical system—its local interactions and global evolution laws—admit different useful abstractions. The local behavior of the system can be naturally encoded in a graph, with the degrees of freedom stored on vertices, and interactions between degrees of freedom represented as edges. These interactions are described by local physical laws (like Hooke’s law from above), applied uniformly, like a stencil, over all of the edges. However, this stencil interpretation is ill-suited for representing the coupled equations which describe global behaviors. Once discretized and linearized, these global operations are most naturally expressed in the language of linear algebra, where all of the system data is aggregated into huge but sparse matrices and vectors.

The easiest way for a programmer to reason about a physical simulation, and hence a common idiom when implementing one, is to swap back and forth between the global and local abstractions. First, a graph or mesh library might be used to store a mass spring system. Local forces and force Jacobians are computed with uniform, local stencils on the graph and copied into large sparse matrices and vectors, which are then handed off to optimized sparse linear algebra

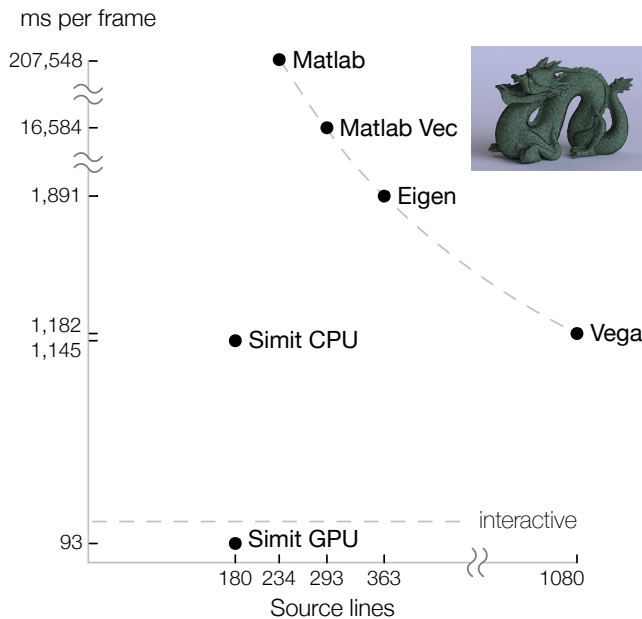


Fig. 1: Scatter plot that shows the relationship between the code size and runtime of a Neo-Hookean FEM simulation implemented using (Vectorized) Matlab, the optimized Eigen Linear Algebra library, the hand-optimized Vega FEM framework, and Simit. The runtimes are for a dragon with 160,743 tetrahedral elements. The trend is that you get more performance by writing more code, however, with Simit you get both performance and productivity. Simit requires fewer lines of code than the Matlab implementation and runs faster than the hand-optimized Vega library on a single-threaded CPU. On a GPU, the Simit implementation runs 10x faster with no code changes.

libraries to calculate the updated global state of the simulation. Finally this updated state is copied back onto the graph.

While straightforward to conceptualize, the strategy of copying data back and forth between the graph and matrix representations incurs high performance costs as a result of data translation, and the inability to optimize globally across linear algebra operations. To overcome this inefficiency, highly-optimized simulations, like those used for games and other real-time applications, are built as monolithic codes that perform assembly and linear algebra on a single set of data structures, often by computing in-place on the graph. Building such a monolithic code requires enormous programmer effort and expertise. Doing so while keeping the system maintainable and extensible, or allowing retargeting of the same code to multiple architectures such as GPUs and CPUs, is nearly impossible.

## The Simit Language

To allow users to take advantage of implicit local-global structure without the performance pitfalls described above, we propose a new programming language called Simit that *natively supports switching between the graph and matrix views* of the simulation. Because Simit is aware of the local-global duality at the language level, it allows simulation code to be concise, fast (see Figure 1), and portable (compiling to both CPU and GPU with no source code change). Simit makes use of three key abstractions: first, the local view is defined using a hypergraph data structure, where nodes represent degrees of freedom and hyperedges relationships such as force stencils, finite elements, and joints. Hyperedges are used instead of regular edges to support relationships between more

than two vertices. Second, the local operations to be performed are encoded as functions acting on neighborhoods of the graph (such as Hooke’s law). Lastly, and most importantly, the user specifies how global vectors and matrices are related to the hypergraph and local functions. For instance, the user might specify that the global force vector is to be built by applying Hooke’s law to each spring and summing the forces acting on each mass. The key point is that defining a global matrix in this way is *not* an imperative instruction for Simit to materialize a matrix in memory; rather, it is an abstract definition of the matrix (much as one would define the matrix in a mathematical paper). The programmer can then operate on that abstract matrix using linear algebra operations; Simit analyzes these operations and translates them into operations on the hypergraph. Because Simit understands the map between the matrix and the graph, it can globally optimize the code it generates while still allowing the programmer to reason about the simulation in the most natural way: as both local graph operations and linear algebra on sparse matrices.

Simit’s performance comes from its design and is made possible by Simit’s careful choice of abstractions. Three features (Section 9) come together to yield the surprising performance shown in Figure 1:

*In-place Computation* is made possible by the tensor assembly construct that lets the compiler understand the relationship between global operations and the graph and turn global linear algebra into in-place local operations on the graph structure. This means Simit does not need to generate sparse matrix index structures or allocate any matrix or vector memory at runtime;

*Index Expression Fusion* is used to fuse linear algebra operations, yielding loops that perform multiple operations at once. Further, due to in-place computation even sparse operations can be fused; and

*Simit’s Type System*, with natively blocked vectors, matrices and tensors, lets it perform efficient dense block computation by emitting dense loops as sub-computations of sparse operations.

Simit’s performance could be enhanced even further by emitting vector instructions or providing multi-threaded CPU execution, optimizations that are planned for a future version of the language.

## Scope

Simit is designed for algorithms where local stencils are applied to a graph of fixed topology to form large, global matrices and vectors, to which numerical algorithms are applied and the results written back onto the graph. This abstraction perfectly fits many physical simulation problems such as mass-spring networks (where hyperedges are the springs), cloth (the bending and stretching force stencils), viscoelastic deformable bodies (the finite elements), etc. At this time Simit does not natively support graphs that change topology over time (such as occurs with fracture or penalty-based impact forces) or simulation elements that do not fit the graph abstraction (collision detection spatial data structures, semi-Lagrangian advection of fluids). As discussed in Section 5, Simit is interoperable with C++ code and libraries, which can circumvent some of these limitations.

The target audience for Simit is researchers, practitioners, and educators who want to develop physical simulation code that is more readable, maintainable, and retargetable than MATLAB or C++, while also being significantly more efficient (comparable to optimized physics libraries like SOFA). Of course Simit programs will not outperform hand-tuned (and complex, unmaintainable) CUDA code, or be simpler to use than problem-specific tools like FreeFem++, but Simit occupies a sweet spot balancing these goals (see Figure 1 and benchmarks in Section 8) that is ideal for a general-purpose physical simulation language.

## Contributions

Simit is the first platform that allows the development of physics code that is simultaneously:

*Concise* The Simit language has MATLAB-like syntax that lets algorithms be implemented in a compact, readable form that closely mirrors their mathematical expression. In addition, Simit matrices specified from the hypergraph are indexed by hypergraph elements like vertices and edges rather than by raw integers, significantly simplifying indexing code and eliminating bugs.

*Expressive* The Simit language consists of linear algebra operations augmented with control flow that let developers implement a wide range of algorithms ranging from finite elements for deformable bodies, to cloth simulations and more. Moreover, the hypergraph abstraction is powerful enough to allow easy specification of complex geometric data structures.

*Fast* The Simit compiler produces high-performance executable code comparable to that of hand-optimized end-to-end libraries and tools, as validated against the state-of-the-art SOFA and VEGA real-time simulation frameworks. Simulations can now be written as easily as a traditional prototype and yet run as fast as a high performance implementation without manual optimization.

*Performance Portable* A single Simit program can be compiled to CPUs and GPUs with no additional programmer effort, while generating efficient code for each architecture. Where Simit delivers performance comparable to hand-optimized CPU code on the same processor, the same simple Simit program delivers roughly an order of magnitude higher performance on a modern GPU in our benchmarks, with no changes to the program.

*Interoperable* Simit hypergraphs and program execution are exposed as C++ APIs, so developers can seamlessly integrate with existing C++ programs, algorithms and libraries.

## 2. RELATED WORK

The Simit programming model draws on ideas from programming systems, numerical and simulation libraries, and physical and mathematical frameworks.

### Libraries for Physical Simulation

A wide range of libraries for the physical simulation of deformable bodies with varying degrees of generality are available [Pommier and Renard 2005; Faure et al. 2007; Dubey et al. 2011; Sin et al. 2013; Comsol 2005; Hibbett et al. 1998; Kohnke 1999], while still others specifically target rigid and multi-body systems with domain specific custom optimizations [Coumans et al. 2006; Smith et al. 2005; Liu 2014]. These simulation codes are broad and many serve double duty as both production codes and algorithmic testbeds. As such they often provide collections of algorithms rather than customizations suited to a particular timestepping and/or spatial discretization model. With broad scope comes convenience but even so interlibrary communication is often hampered by data conversion while generality often limits the degree of optimization.

For very specific problems, previous work developed highly-optimized code, but these instances are limited in scope. For example, multigrid solvers on CPU and GPU [McAdams et al. 2011; Dick et al. 2011] are very fast, but limited to corotated linear material models for tri-linear hexahedral finite elements.

### Mesh data structures

Simulation codes often use third-party libraries that support higher-level manipulation of the simulation mesh. A *half-edge data structure* [Eastman and Weiss 1982] (from, e.g., the OpenMesh library [Botsch et al. 2002]) is one popular method for describing a mesh while allowing efficient connectivity queries and neighborhood circulation. Alternatives targeting different application requirements (manifold vs. nonmanifold, oriented vs unoriented, etc.) abound, such as winged-edge [Baumgart 1972] or quad-edge [Guibas and Stolfi 1985] data structures, and modern software packages like CGAL [cga ] have built sophisticated tools on top of many of these data structures for performing common geometry operations. Simit differs from these approaches in that its hierarchical hyper-edges provide sufficient expressiveness to let users build semantically rich data structures like many of these meshes, while not limiting the user to any specific mesh data structure.

### DSLs for computer graphics

Graphics has a long history of using domain-specific languages and abstractions to provide high performance, and performance portability, from relatively simple code. Most visible are shading languages and the graphics pipeline [Hanrahan and Lawson 1990; Segal and Akeley 1994; Mark et al. 2003; Blythe 2006]. Image processing languages also have a long history [Holzmann 1988; Elliott 2001; Ragan-Kelley et al. 2012], and more recently domain-specific languages have been proposed for new domains like 3D printing [Vidimčič et al. 2013]. In physical simulation, Guenter et al. built the  $D^*$  system for symbolic differentiation, and demonstrated its application to modeling and simulation [Guenter and Lee 2009].  $D^*$  is an elegant abstraction, but its implementation focuses less on optimized simulation performance, and its model cannot express features important to many of our motivating applications.

### Graph programming models

A number of programming systems address computation over graphs or graph-like data structures, including GraphLab [Low et al. 2010], Galois [Pingali et al. 2011], Liszt [DeVito et al. 2011], Socialite [Jiwon Seo 2013], and GreenMarl [Sungpack Hong and Olukotun 2012]. In these systems, programs are generally written as explicit in-place computations using stencils on the graph, providing a much lower level of abstraction than linear algebra over whole systems. Of these, GraphLab and Socialite focus on distributed systems, where we currently focus on single-node/shared memory execution. Socialite and GreenMarl focus on scaling traditional graph algorithms (e.g., breadth-first search and betweenness centrality) to large graphs. Liszt exposes a programming model over meshes. Computations are written in an imperative fashion, but must look like stencils, so it only allows element-wise operations and reductions. This is similar to the programming model used for assembly in Simit, but it has no corollary to Simit's linear algebra for easy operation on whole systems. Galois exposes explicit in-place programming via a similarly low-level but extremely dynamic programming model, which inhibits compiler analysis and optimization.

### Programming systems for linear algebra

Our linear algebra syntax is explicitly inspired by MATLAB [2014], the most successful high-productivity tool in this domain, though we believe our syntax is improved in key ways for our applications. In particular, the combination of coordinate-free indexing and the assembly map operator, with hierarchically blocked tensors, dramat-

ically reduces indexing complexity, while also exposing structure critical to our compiler optimizations. Eigen is a C++ library for linear algebra which uses aggressive template metaprogramming to specialize and optimize linear algebra computations at compile time, including fusion of multiple operations and vectorization [Guennebaud et al. 2010]. It does an impressive job exposing linear algebra operations to C++, and aggressive vectorization delivers impressive inner-loop performance, but assembly is still both challenging for programmers and computationally expensive during execution.

### 3. FINITE ELEMENT METHOD EXAMPLE

To make things concrete, we start by discussing an example of a paradigmatic Simit program: a Finite Element Method (FEM) statics simulation that uses Newton’s method to compute the final configuration of a deforming object. Figure 2 shows the source code for this example. The implementation of `compute_tet_stiffness` and `compute_tet_force` depends on the material model chosen by the user and are omitted. In this section we introduce Simit concepts with respect to the example, but we will come back to them in Section 4 with rigorous definitions.

As is typical, this Simit application consists of five parts: (1) graph definitions, (2) functions that are applied to each graph vertex or edge to compute new values based on neighbors, (3) functions that compute local contributions of vertices and edges to global vectors and matrices, (4) assemblies that aggregate the local contributions into global vectors and matrices, and (5) code that computes with vectors and matrices.

Step 1 is to define a graph data structure. Graphs consist of elements (objects) that are organized in vertex sets and edge sets. Lines 1–16 define a Simit graph where edges are tetrahedra and vertices their degrees of freedom. Lines 1–5 define an element of type `Vertex` that represents a tetrahedron’s degrees of freedom. It has two fields: a coordinate `x` and a velocity `v`. Next, lines 7–12 define a `Tet` element that represents an FEM Tetrahedron with four fields: shear modulus `u`, Lamé’s first parameter `l`, volume `W`, and the strain-displacement  $3 \times 3$  matrix `B`. Finally, lines 15–16 define a vertex set `verts` with `Vertex` elements, and an edge set `tets` with `Tet` elements. Since `tets` is an edge set, its definition lists the sets containing the edges’ endpoints; a tetrahedron connects four vertices (see Figure 3). That is, Simit graphs are hypergraphs, which means that edges can connect any fixed number of vertices.

Step 2 is to define and apply a function `precompute_vol` to precompute the the volume of every tetrahedron. In Simit this can be done by defining the *stencil function* `precompute_vol` shown on lines 19–22. Simit stencil functions are similar to the update functions of data-graph libraries such as GraphLab [Low et al. 2010] and can be applied to every element of a set (`verts`) and its endpoints (`verts`). Stencil functions define one or more `inout` parameters that have pass-by-reference semantics and that can be modified. Lines 24–26 shows the Simit procedure `init` that can be called to precompute volumes. The procedure contains a single statement that applies the stencil function to every tetrahedron in `tets`.

Step 3 defines functions that compute the local contributions of a vertex or an edge to global vectors and matrices. Lines 29–34 define `tet_force` that computes the forces exerted by a tetrahedron on its vertices. The function takes two arguments, the tetrahedron and a tuple containing its vertices. It returns a global vector `f` that contains the local force contributions of the tetrahedron. Line 32 computes the tetrahedron forces and assigns them to `f`. Since `tet_force` only has access to the vertices of one tetrahedron, it can only write to four locations in the global vector. This is sufficient, however, since a tetrahedron only *directly* influences its own vertices.

```

1 element Vertex
2   x : vector[3](float); % position
3   v : vector[3](float); % velocity
4   fe : vector[3](float); % external force
5 end
6
7 element Tet
8   u : float; % shear modulus
9   l : float; % Lamé's first parameter
10  W : float; % volume
11  B : matrix[3,3](float); % strain-displacement
12 end
13
14 % graph vertices and (tetrahedron) hyperedges
15 extern verts : set{Vertex};
16 extern tets : set{Tet}(verts, verts, verts, verts);
17
18 % precompute tetrahedron volume
19 func precompute_vol(inout t : Tet, v : (Vert*4))
20   t.B = compute_B(v);
21   t.W = -det(B)/6.0;
22 end
23
24 proc init
25   apply precompute_vol to tets;
26 end
27
28 % computes the force of a tetrahedron on its vertices
29 func tet_force(t : Tet, v : (Vertex*4))
30   -> f : vector[verts](vector[3](float))
31   for i in 0:4
32     f(v(i)) = compute_tet_force(t,v,i);
33   end
34 end
35
36 % computes the stiffness of a tetrahedron
37 func tet_stiffness(t : Tet, v : (Vertex*4))
38   -> K : matrix[verts,verts](matrix[3,3](float))
39   for i in 0:4
40     for j in 0:4
41       K(v(i),v(j)) = compute_tet_stiffness(t,v,i,j);
42     end
43   end
44 end
45
46 % newton's method timestepper
47 proc newton_method
48   tol = 1e-6;
49   while abs(f - verts.fe) > tol
50     f = map tet_force to tets reduce +;
51     K = map tet_stiffness to tets reduce +;
52
53     verts.x = vert.x + K\((verts.fe - f));
54   end
55 end

```

Fig. 2: Simit code for a Finite Element Method (FEM) simulation of statics. The code contains: (1) graph definitions, (2) a function that is applied to each tetrahedron to precompute its volume, (3) functions that compute local contributions of each tetrahedron to global vectors and matrices, (4) assemblies that aggregate the local contributions into global vectors and matrices, and (5) code that computes with those vectors and matrices. The `compute_*` functions have been omitted for brevity. For more Simit code, see Appendix A for a full Finite Element Dynamics Simulation.

Step 4 uses a Simit assembly map to aggregate the local contributions computed from vertices and edges in the previous stage into global vectors and matrices. Line 50 assembles a global force vector by summing the local force contributions computed by applying `tet_force` to every tetrahedron. The result is a dense force vector `f` that contains the force of every tetrahedron on its vertices.

Finally, Step 5 is to compute with the assembled global vectors and matrices. The results of these computations are typically vectors

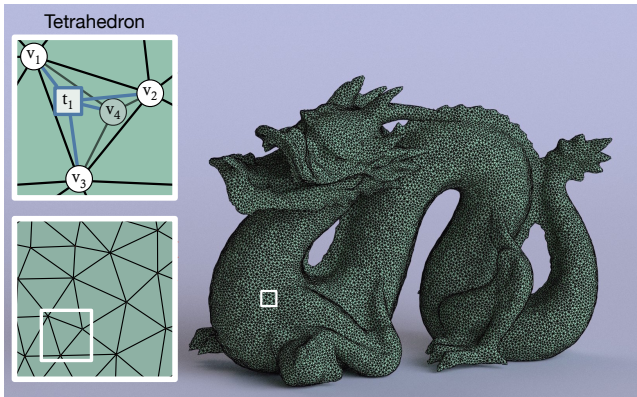


Fig. 3: Tetrahedral dragon mesh consisting of 160,743 tetrahedra that connect 46,779 vertices. Windows zoom in on a region on the mesh and a single tetrahedron. The tetrahedron is modeled as a Tet edge  $t_1 = \{v_1, v_2, v_3, v_4\}$ .

that are stored to fields of graph vertices or edges. Line 53 reads the  $x$  position field from the `verts` set, computes a new position and writes that position back to `verts.x`. Reading a field from a set results in a global vector whose blocks are the fields of the set's elements. Writing a global vector to a set field works the same way; the vector's blocks are written to the set elements. In this example the computation uses the linear solve  $\backslash$  operator to perform a linearly implicit time-step, but many other approaches are possible.

#### 4. PROGRAMMING MODEL

Simit's programming model is designed around the observation that a physical system is typically graph structured, while computation on the system is best expressed as global linear and multi-linear algebra. Thus, the Simit **data model**<sup>1</sup> consists of two abstract data structures: **hypergraphs** and **tensors**. Hypergraphs generalize graphs by letting edges connect any  $n$ -element subset of vertices instead of just pairs. Tensors generalize scalars, vectors and matrices, that respectively are indexed by 0, 1 and 2 indices, to an arbitrary number of indices.

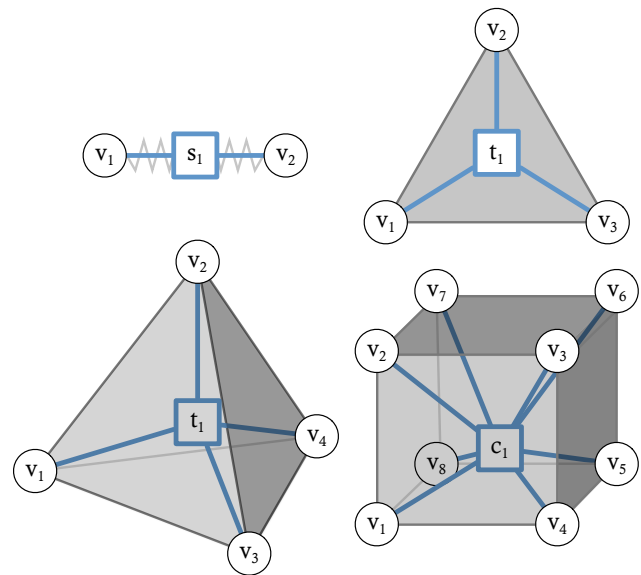
We also describe two new operations on hypergraphs and tensors: **tensor assemblies**, and **index expressions**. Tensor assemblies map tensors to graphs, while index expressions compute with tensors.

##### 4.1 Hypergraphs with Hierarchical Edges

Hypergraphs are ordered pairs  $H = (V, E)$ , comprising a set  $V$  of vertices and a set  $E$  of hyperedges that are  $n$ -element subsets of  $V$ . The number of elements a hyperedge connects is its **cardinality**, and we call a hyperedge of cardinality  $n$  an  **$n$ -edge**. Thus, hypergraphs generalize graphs where edges must have a cardinality of two. Hypergraphs are useful for describing relationships between vertices that are more complex than binary relationships. Figure 4 shows four examples of hyperedges: a 2-edge, a 3-edge, a 4-edge and an 8-edge (in black) are used to represent a spring, a geometric triangle, a tetrahedron and a hexahedron (in grey). Although, these hyperedges are used to model geometric mesh relationships, hyperedges can be used to model any relationship. For example, a 2-edge can also be used to represent a joint between two rigid bodies or the relationship between two neurons, and a 3-edge can represent a clause in a 3-SAT instance.

Simit hypergraphs generalize normal hypergraphs and are ordered tuples  $H_{Simit} = (S_1, \dots, S_m)$ , where  $S_i$  is the  $i$ th set whose elements connect 0 or  $n$  elements from other sets, called its **endpoints**.

<sup>1</sup>In this section we will use bold when we first mention a new concept.



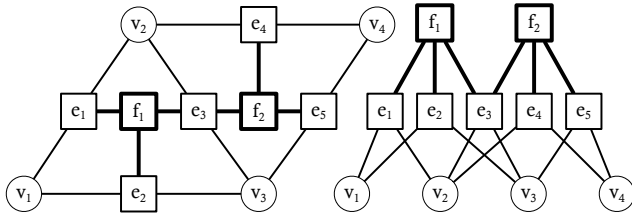
```
springs : set{Spring}(verts,verts);
trigs   : set{Triangle}(verts,verts,verts);
tets    : set{Tetrahedron}(verts,verts,verts,verts);
hexes   : set{Hexahedron}(verts,verts,verts,verts,
                          verts,verts,verts,verts);
```

Fig. 4: Four geometric elements are shown in gray: a spring, a triangle, a tetrahedron and a cube. Degrees of freedom are shown as black circles. Simit graph nodes match the degrees of freedom, while Simit edges of cardinality three, four and eight are shown as blue squares. Note that these Simit edges represent area/volume and not mesh edges.

We refer to a set with cardinality 0 as a **vertex set** and a set of cardinality  $n$  greater than 0 as an **edge set**. So Simit hypergraphs can have any number of vertex and edge sets, and edge sets have one or more endpoints. The endpoints of an  $n$ -cardinality edge set are a set relation over  $n$  other vertex or edge sets. That is, each endpoint of an edge is an element from the corresponding endpoint set. This means that edge sets can connect multiple distinct sets, and we call such edge sets **heterogeneous edge sets**.

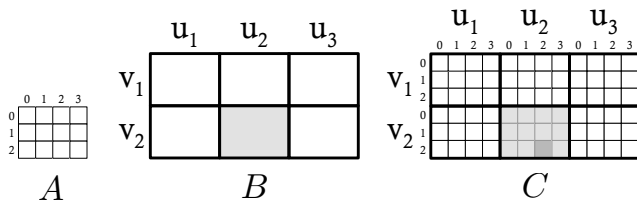
Elements (vertices and edges) of hypergraph sets can contain data. An element's data is a tuple whose entries are called **fields**. This is equivalent to record or struct types in other languages. Fields can be scalars, vectors, matrices or tensors. For example, the `Vertex` element on lines 2–5 in Figure 2 has two vector fields  $x$  and  $v$ . We say that a hypergraph set has the same fields as its elements; however, the field of a set is a vector whose blocks are the fields of the set's elements. Blocked vector types are described in Section 4.2.

Edge sets can connect other edge sets, so we say that Simit supports **hierarchical edges**. Hierarchical edges have important applications in physical simulation because they let us represent the topology of a mesh. Figure 5 demonstrates how hierarchical edges can be used to capture the topology in a triangle mesh with faces, triangle edges and vertices. The left hand side shows two triangles that share the edge  $e_3$ . Each triangle has three vertices that are connected in pairs by graph edges  $\{e_1, e_2, e_3, e_4, e_5\}$  that represent triangle edges. However, these edges are themselves connected by face edges  $\{f_1, f_2\}$ , thus forming the hierarchy shown on the right hand side. By representing both triangle edges and faces we can store different quantities on them. Moreover, it becomes possible to accelerate typical mesh queries such as finding the adjacent faces



```
verts : set{Vertex};
edges : set{Edge}(verts,verts);
faces : set{Face}(edges,edges,edges);
```

Fig. 5: Hierarchical hyperedges model triangles with faces, edges, and vertices. On the left two triangles with faces  $f_1$  and  $f_2$  are laid flat. On the right the same triangles are arranged to show the topological hierarchy.



```
A : matrix[3,4](float)
B : matrix[V,U](float)
C : matrix[V,U][3,4](float)
```

Fig. 6: Three Simit matrices. On left is a basic  $3 \times 4$  matrix  $A$ . In the middle is a matrix  $B$  whose dimensions are the sets  $V$  and  $U$ . The matrix is indexed by pairs of elements from  $U, V$ , e.g.  $B(v_2, u_2)$ . Finally, on the right is a blocked matrix  $C$  with  $V \times U$  blocks of size  $3 \times 4$ . A block of this matrix is located by a pair of elements from  $V, U$ , e.g.  $C(v_2, u_2)$ , and an element can be indexed using a pair of indices per matrix hierarchy, e.g.  $C(v_2, u_2)(2, 2)$ .

of a face by inserting topological indices, such as a half-edge index, into the graph structure.

## 4.2 Tensors with Blocks

We use zero indices to index a scalar, one to index a vector and two to index a matrix. Tensors generalize scalars, vectors and matrices to an arbitrary number of indices. We call the number of indices required to index into a tensor its **order**. Thus, scalars are 0th-order tensors, vectors are 1st-order tensors, and matrices are 2nd-order tensors. Further, we refer to the  $n$ th tensor index as its  $n$ th dimension. Thus, the first dimension of an  $m \times n$  matrix is the rows  $m$ , while the second dimension is the columns  $n$ .

The dimensions of a Simit tensor are sets: either integer ranges or hypergraph sets. Thus, a Simit vector is more like a dictionary than an array, and an  $n$ -order tensor is an  $n$ -dimensional dictionary, where an  $n$ -tuple of hypergraph set elements map to a tensor component. For example, Figure 6 (center) depicts a matrix  $B$  whose dimensions are  $(V, U)$ , where  $V = \{v_1, v_2\}$  and  $U = \{u_1, u_2, u_3\}$ . We can index into the matrix using an element from each set. For example,  $B(v_2, u_2)$  locates the grey component.

Simit tensors can also be blocked. In a blocked tensor each dimension consists of a hierarchy of sets. For example, a hypergraph set that maps to tensor blocks, where each block is described by an integer range. Blocked tensors are indexed using **hierarchical indexing**. This means that if we index into a blocked tensor using an element from the top set of a dimension, the result is a tensor block. For example, Figure 6 (right) shows a blocked matrix  $C$  whose

dimensions are  $(V \times 3, U \times 4)$ , which means there are  $|V| \times |U|$  blocks of size  $3 \times 4$ . As before, we can index into the matrix using an element from each set,  $C(v_2, u_2)$ , but now the index operation results in the  $3 \times 4$  grey block matrix. If we index into the matrix block,  $C(v_2, u_2)(2, 2)$  we locate the dark grey component. In addition to being convenient for the programmer, blocked tensors let Simit produce efficient code. By knowing that a sparse matrix consists of dense inner blocks, Simit can emit dense inner loops for sparse matrix-vector multiplies with that matrix.

## 4.3 Tensor Assembly using Maps

A **tensor assembly** is a map from the triple  $(S, f, r)$  to one or more tensors, where  $S$  is a hypergraph set,  $f$  an **assembly function**, and  $r$  an associative and commutative reduction operator. The tensor assembly applies the assembly function to every element in the hypergraph set, producing per-element tensor contributions. The tensor assembly then aggregates these tensor contributions into a global tensor, using the reduction operator to combine values. The result of the tensor assembly is one or more global tensors, whose dimensions can be the set  $S$  or any of its endpoints. The diagram in Figure 7 shows this process. On the left is a graph where the edges  $E = \{e_1, e_2\}$  connect the vertices  $V = \{v_1, v_2, v_3\}$ . The function  $f$  is applied to every edge to compute contributions to the global  $V \times V$  matrix. The contributions are shown in grey and the tensor assembly aggregates them by adding the per-edge contribution matrices.

Assembly functions are pure functions whose arguments are an element and its endpoints, and that return one or more tensors that contain the element's global tensor contributions. The arguments of an assembly function are supplied by a tensor assembly as it applies the function to every element of a hypergraph set, and the same tensor assembly combines the contributions of every assembly function application. The center of Figure 7 shows code for  $f$ : a typical assembly function that computes the global matrix contributions of a 2-edge and its vertex endpoints. The function takes as arguments an edge  $e$  of type `Edge`, and a tuple  $v$  that contains  $e$ 's two `Vertex` endpoints. The result is a  $V \times V$  matrix with  $3 \times 3$  blocks as shown in the figure. Notice that  $f$  can only write to four locations in the resulting  $V \times V$  matrix, since it has access to only two vertices. In general, an assembly function that maps a  $c$ -edge to an  $n$ -dimensional tensor, can write to exactly  $c^n$  locations in the tensor. We call this property **coordinate-free indexing**, since each assembly function locally computes and writes its matrix contributions to the global matrix using opaque indices (the vertices) without regards to where in the matrix those those contributions end up. Further, since the global matrix is blocked, the  $3 \times 3$  matrix  $k$  can be stored into it with one assignment, by only specifying block coordinates and not intra-block coordinates. As described in Section 4.2 we call this property **hierarchical indexing**, and the resulting *coordinate-free hierarchical indexing* removes a large class of indexing bugs, and make assembly functions easy to write and read.

We have so far discussed the functional semantics of tensor assemblies, but it is also important to consider their performance semantics. The way they are defined above, if executed literally, would result in very inefficient code where ultra-sparse tensors are created for every edge, followed by a series of tensor additions. However, as we will see in Section 6, the tensor assembly abstraction lets Simit's compiler produce code that stores tensor blocks on the graph elements corresponding to one of the tensor dimensions. Thus, memory can be pre-allocated and indexing structures pre-built, and assembly becomes as cheap as computing and storing blocks in a contiguous segmented array. Generally, as discussed in Section 9, the tensor assembly construct lets the Simit compiler know where

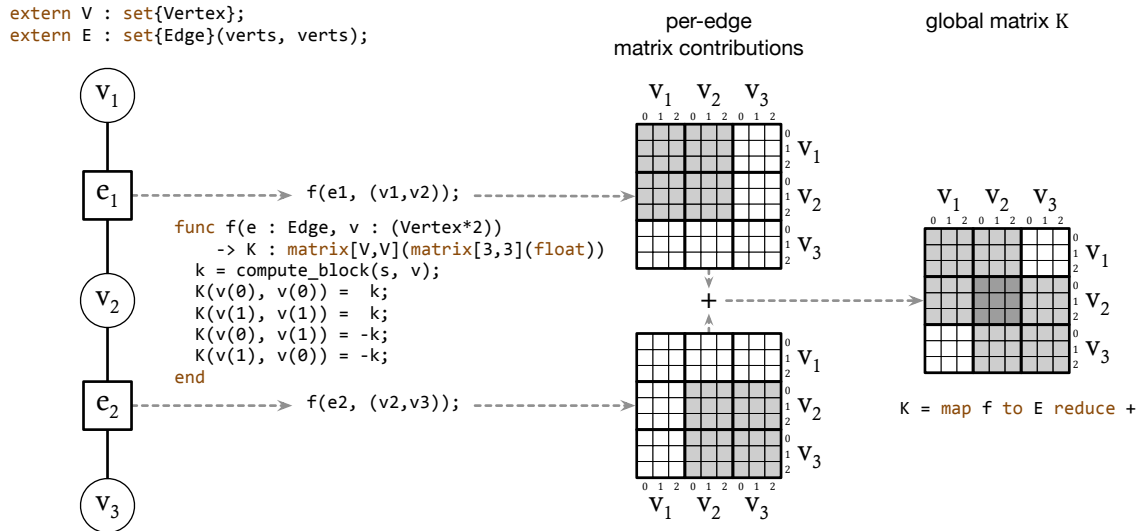


Fig. 7: Global matrix  $K$  assembly. The assembly map (on the right) applies  $f$  to every edge in  $E$ , and sums the resulting matrix contributions.  $f$  computes a block  $k$  that is stored into the positions of  $K$  that correspond to the current edge’s endpoints. Since each edge only has two endpoints,  $f$  can only store into four locations of the matrix. This is sufficient, however, since if  $v_1$  and  $v_3$  do not directly interact; if they did there should have been an edge between them. As a result of this restriction, the top right and lower left entries of the final matrix  $K$  are empty. Finally, note the collision at  $(v_2, v_2)$  due to  $v_2$  being connected by both  $e_1$  and  $e_2$ .

global vectors and matrices come from, which lets it emit in-place code that executes very fast.

#### 4.4 Tensor Computation using Index Expressions

So far, we have used linear algebra to compute with scalars, vectors and matrices. Linear algebra is familiar and intuitive for programmers, so we provide it in the Simit language, but it has two important drawbacks. First, it does not extend to higher-order tensors. Second, it is riddled with operators that have different meanings depending on the operands, and does not cleanly let us express computations that perform multiple operations simultaneously. This makes linear algebra ill-suited as compute operators in the Simit programming model. Instead we have designed index expressions, which are a generalization of tensor index notation [Ricci-Curbastro and Levi-Civita 1901] in expression form. Index expressions have all of the properties we seek, and as an added benefit we can build all the basic linear algebra operations on top of them. Thus, programmers can program with familiar linear algebra when that is convenient, and the linear algebra can be lowered to index expressions that are easier to optimize and generate efficient code from (see Section 7).

An index expression computes a tensor, and consists of a scalar expression and one or more **index variables**. Index variables come in two variants—**free variables** and **reduction variables**—and are used to index tensor operands. In addition, free variables determine the dimensions of the tensor resulting from the index expression, and reduction variables combine values. Thus, an index expression takes the form:

(free-variable\*) reduction-variable\* scalar-expression

where scalar-expression is a normal scalar expression with operands such as +, -, \* and /, and scalar operands that are typically indexed tensors.

Free index variables are variables that can take the values of a tensor dimension: an integer range, a hypergraph set, or a hierarchical set as shown in Figure 6. The values an index variable can

take are called its range. The range of the free index variables of an index expression determine the dimensions of the resulting tensor. To compute the value of one component of this tensor the index variables are bound to the component’s coordinate, and the index expression is evaluated. To compute every component of the resulting tensor, the index expression is evaluated for every value of the set product of the free variables’ ranges. For example, consider an index expression that computes a vector addition:

(i)  $a(i) + b(i)$

In this expression  $i$  is a free index variable whose range is implicitly determined by the dimensions of the vectors  $a$  and  $b$ . Note that an index variable can only be used to index into a tensor dimension that is the same as its range. Thus, the vector addition requires that the dimensions of  $a$  and  $b$  are the same. Further, the result of this index expression is also a vector whose dimension is the range of  $i$ . Next, consider a matrix transpose:

(i, j)  $A(j, i)$

Here  $i$  and  $j$  are free index variables whose ranges are determined by the second and first dimensions of  $A$  respectively. As expected, the dimensions of the resulting matrix are the reverse of  $A$ ’s, since the order of the free index variables in the list determines the order of the result dimensions. Finally, consider an index expression that adds  $A$  and  $B^T$ :

(i, j)  $A(i, j) + B(j, i)$

Since index variables ranges take on the values of the dimensions they index, this expression requires that the first and second dimensions of  $A$  are the same as the second and first dimensions of  $B$  respectively. This example shows how one index expression can simultaneously evaluate multiple linear algebra operations.

Reduction variables, like free variables, range over the values of tensor dimensions. However, unlike free variables, reduction variables do not contribute to the dimensions of the resulting tensor. Instead, they describe how a range of computed values must be



combined to produce a result component. How these values are combined is determined by the reduction variables's reduction operator, which must be an associative and commutative operation. For example, a vector dot product ( $\sum_i a(i) * b(i)$ ) can be expressed using an addition reduction variable:

```
+i a(i) * b(i)
```

As with free index variables, the range of  $i$  is implicitly determined by the dimension of  $a$  and  $b$ . However, instead of resulting in a vector with one component per value of  $i$ , the values computed for each  $i$  are added together, resulting in a scalar. Free index variables and reduction index variables can also be combined in an index expression, such as in the following matrix-matrix multiplication:

```
(i, k) +j A(i,j) * B(j,k)
```

The two free index variables  $i$  and  $k$  determine the dimensions of the resulting matrix.

In this section we showed how index expressions can be used to express linear algebra in a simpler and cleaner framework. We showed a simple example where two linear algebra operations (addition and transpose) were folded into a single index expression. As we will see in Section 7, index expressions make it easy for a compiler to combine basic linear algebra expressions into arbitrarily complex expressions that share index variables. After code generation, this results in fewer loop nests and less memory traffic.

## 5. INTERFACING WITH SIMIT

To use Simit in an application there are four steps:

- (1) Specify the structure of a system as a hypergraph with vertex sets and edge sets, using the C++ Set API described in Section 5.1,
- (2) Write a program in the Simit language to compute on the hypergraph, as described in Sections 3 and 4,
- (3) Load the program, bind hypergraph sets to it, and compile it to one or more Function object, as described in Section 5.2,
- (4) Call the Function object's run method for each solve step (e.g., time step, static solve, etc.), as described in Section 5.2.

Collision detection, fracturing, and topology changes are not currently expressed in Simit's language, but can be invoked in C++ using Simit's Set API to dynamically add or remove vertices and edges in the hypergraph between each Simit program execution. For example external calls to collision detection code between time steps would then express detected contacts as edge sets between colliding elements.

### 5.1 Set API

Simit's Set API is a set of C++ classes and functions that create hypergraph sets with tensor fields. The central class is the Set class, which creates sets with any number of endpoints, that is, both vertex sets and edge sets. When a Set is constructed, the Set's endpoints are passed to the constructor. Next, fields can be added using the Set's addField method and elements using its add method.

The following code shows how to use the Simit Set API to construct a pyramid from two tetrahedra that share a face. The vertices and tetrahedra are given fields that match those in the running example from Section 3:

```
Set verts;
Set tets(verts, verts, verts, verts);

// create fields (see the FEM example in Figure 2)
```

```
FieldRef<double,3> x = verts.addField<double,3>("x");
FieldRef<double,3> v = verts.addField<double,3>("v");
FieldRef<double,3> fe = verts.addField<double,3>("fe");

FieldRef<double> u = tets.addField<double>("u");
FieldRef<double> l = tets.addField<double>("l");
FieldRef<double> W = tets.addField<double>("W");
FieldRef<double,3,3> B = tets.addField<double,3,3>("B");

// create a pyramid from two tetrahedra
Array<ElementRef> v = verts.add(5);
ElementRef t0 = tets.add(v(0), v(1), v(2), v(4));
ElementRef t1 = tets.add(v(1), v(2), v(3), v(4));

// initialize fields
x(v0) = {0.0, 1.0, 0.0};
// ...
```

First, we create the `verts` vertex set and `tets` edge set, whose tetrahedron edges each connects four `verts` vertices. We then add to the `verts` and `tets` sets the fields from the running example from Section 3. The `addField` method is a variadic template method whose template parameters describe the tensors stored at each set element. The first template parameter is the tensor field's component type (`double`, `int`, or `boolean`), followed by one integer literal per tensor dimension. The integers describe the size of each tensor dimension; since the `x` field above is a position vector there is only one integer. Thus, to add a  $3 \times 4$  matrix field we would write: `addField<double,3,4>`. Finally, we create five vertices and the two tetrahedra that connects them together, and initialize the fields.

### 5.2 Program API

Once a hypergraph has been built (Section 5.1) and a Simit program written (Sections 3 and 4), the Program API can be used to compile and run the program on the hypergraph. To do this, the programmer creates a Program object, loads source code into it, and compiles a procedure in the source code to a Function object. Next, the programmer binds hypergraph sets to externs in the Function object's Simit program, and the Function::run method is called to execute the program on the bound sets.

The following code shows how to load the FEM code in Figure 2, and run it on tetrahedra we created in Section 5.1:

```
Program program;
program.loadFile("fem_statics.sim");

Function func = program.compile("main");
func.bind("verts", &verts);
func.bind("tets", &tets);

func.runSafe();
```

In this example we use the Function::runSafe method, which lazily initializes the function. For more performance the initialization and running of a function can be split into a call to Function::init followed by repeated calls to Function::run.

## 6. RUNTIME DATA LAYOUT AND EXECUTION

In Sections 3 and 4 we described the language and abstract data structures (hypergraphs and tensors) that a Simit programmer works with. Since the abstract data structures are only manipulated through global operations (tensor assemblies and index expressions) the Simit system is freed from implementing them literally, an important property called physical data separation [Codd 1970]. Simit exploits this separation to compile global operations to efficient local operations on compact physical data structures that look very different from the graphs and tensors the programmer works with. In the rest of this section we go into detail on how the current Simit

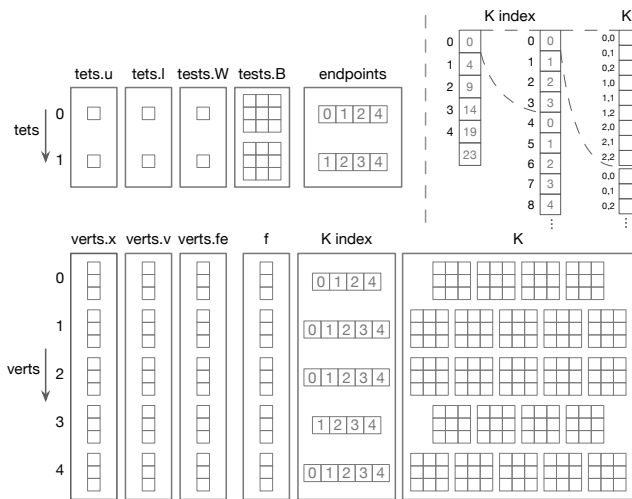


Fig. 8: Table that shows all the per-element data stored for each `tets` tetrahedron and `verts` vertex in the two tetrahedra we constructed in Section 5.1. The global vector `f` and global matrix `K` from the code in Figure 2 are stored on the `verts` set, since `verts` is their first dimension. The table is stored by columns (struct of arrays) and matrices are stored row major. The matrix `K` is sparse so it is stored as a segmented array (top right), consisting of two index arrays and one value array. Since `K` is assembled from the `tets` set its sparsity is known and the index can be precomputed, and it can also be shared with other matrices assembled from the same edge set.

implementation lays out data in memory and what kind of code it emits. This is intended to give the reader a sense of how the physical data separation lets Simit pin global vectors and matrices to graph vertices and edges for in-place computation. It also demonstrates how the tensor assembly construct lets Simit use the graph as the index of matrices, and how the blocked matrix types lets Simit emit dense inner loops when computing with sparse matrices.

As described in Section 5, Simit graphs consist of vertex sets and edge sets, which are the same except that edges have endpoints. A Simit set stores its size (one integer) and field data. In addition, edge sets store a pointer to each of its  $n$  endpoint sets and for each edge,  $n$  integer indices to the edge's endpoints within those endpoint sets. Figure 8 (top) shows all the data stored on each Tet in the `tets` set we built in Section 5.1 for the FEM example in Figure 2. Each Tet stores the fields `u`, `l`, `w`, and `B`, as well as an `endpoints` array of integer indexes into the `verts` set. The set elements and their fields form a table that can be stored by rows (arrays of structs) or by columns (structs of arrays). The current Simit implementation stores this table by columns, so each field is stored separately as a contiguous array with one scalar, vector, dense matrix or dense tensor per element. Furthermore, dense matrices and tensors are stored in row-major order within the field arrays.

Global vectors and matrices are also stored as fields of sets. Specifically, a global vector is stored as a field of its dimension set, while a global matrix is stored as a field of one of its dimension sets. That is, either matrix rows are stored as a field of the first matrix dimension or matrix columns are stored as a field of the second matrix dimension. This shows the equivalence in Simit of a set field and global vector whose dimension is a set; a key organizing property. Figure 8 (bottom) shows all the data stored on each Vertex in the `verts` set, including the global vector `f` and global matrix `K` from the `main` procedure in Figure 2. Since `K` is sparse, its rows can have different sizes and each row is therefore stored as

an array of column indices (`K index`) and a corresponding array of data (`K`). The column indices of `K` are the `verts` vertices that each row vertex can reach through non-empty matrix components. Since the `K` matrix was constructed from a tensor assembly over the `tets` edge sets, the neighbors through the matrix is the same as the neighbors through `tets` and can be precomputed. Since global matrix indices and values are set fields with a different number of entries per set element, they are stored in a segmented array, as shown for `K` in Figure 8 (upper right). Thus, Simit matrix storage is equivalent to Blelloch's segmented vectors [Blelloch 1990] and the BCSR (Blocked Compressed Sparse Row) matrix storage format.

A Simit vector or matrix map statement is compiled into a loop that computes the tensor values and stores them in the global vector or matrix data structures. The loop iterates over the map's target set and each loop iteration computes the local contributions of one target set element using the map function, which is inlined for efficiency. Equivalent sequential C code to the machine code generated for the map statement on line 51 of Figure 2 is:

```
for (int t=0; t<tets.len; t++) {
  for (int i=0; i<4; i++) {
    double[3] tmp = // inlined compute_tet_force(t,v,i)
    for (int j=0; j<3; j++) {
      f[springs.endpoints[s*4 + i]*3 + j] = tmp[j];
    }
  }
}
```

The outer loop comes from the map statement itself and iterates over the tetrahedra. Its loop body is the inlined `tet_force` function, which iterates over the four endpoints of the tetrahedra and for each endpoint computes a tet force that is stored in the `f` vector. A global matrix is assembled similarly with the exception that the location of a matrix component must be computed from the matrix's index array as follows (taken from the Simit runtime library):

```
int loc(int v0, int v1, int *elems, int *elem_nbrs) {
  int l = elems[v0];
  while (elem_nbrs[l] != v1) l++;
  return l;
}
```

The `loc` function turns a two-dimensional coordinate into a one-dimensional array location, given a matrix index consisting of the arrays `elems` and `elem_nbrs`. It does this by looking up the location in `elem_nbrs` where the row (or column) `v0` starts. That is, it finds the correct segment of `v0` in the segmented array `elem_nbrs`. It then scans down this segment to find the location of the element neighbor `v1`, which is then returned.

Figure 9 shows an end-to-end example where a matrix is assembled from a 2-uniform graph and multiplied by a vector field of the same graph. The top part shows the abstract data structure views that the programmer works with, which were described in Section 4. The arrows show how data from the blue edge is put into the matrix on matrix assembly, how data from the  $p_4$  vertex becomes the  $p_4$  block of the `b` vector, and how the block in the  $(p_2, p_4)$  matrix component is multiplied with the block in the  $p_4$  `b` vector component to form the  $p_2$  `c` vector component when the `A` is multiplied with `b`. The bottom part shows the physical data structures; the vertex set has a field `points.b` and the edge set has a field `edges.m`. The stippled arrows show how the `loc` function is used to find the correct location in the array of `A` values when storing matrix contributions of the blue edge. The full arrows show how the `edges.m` field is used to fill in values in the  $(p_2, p_4)$  matrix component (the sixth block of the `A` array), and how this block is multiplied directly with the  $p_4$  `points.b` vector component to form the  $p_2$  `points.c` vector component when the `A` is multiplied with `b`.



*Index Expression Lowering.* In this phase, all index expressions are transformed into loops. For every index expression, the compiler replaces each index variable with a corresponding dense or sparse loop. The index expression is then inserted at the appropriate places in the loop nest, with index variables replaced by loop variables. This process is demonstrated in the middle and right panes of Figure 10.

*Lowering Tensor Accesses.* In the next phase, tensor accesses are lowered. The compiler takes statements that refer to multidimensional matrix locations and turns them into concrete array loads and stores. A major optimization in this phase is to use context information about the surrounding loops to make sparse matrix indexing more efficient. For example, if the surrounding loop is over the same sets as the sparse system matrix, we can use existing generated variables to index into the sparse matrix instead of needing to iterate through the column index of the neighbor data structure for each read or write.

*Code Generation.* After the transformation phases, a Simit program consists of imperative code, with explicit loops, allocations and function calls that are easy to turn into low-level code. The code generation phase turns each Simit construct into the corresponding LLVM operations, using information about sets and indices to assist in generating efficient code. Currently, the backend calls LLVM optimization passes to perform inter-procedural optimization on the Simit program, as well as other standard compiler optimizations. Only scalar code is generated; we have not yet implemented vectorization or parallelism, and LLVM’s auto-vectorization passes cannot automatically transform our scalar code into vector code. Future work will implement these optimizations during code generation, prior to passing the generated code to LLVM’s optimization passes. Our index expression representation is a natural form in which to perform these transformations.

*GPU Code Generation.* Code generation for GPU targets is performed as an alternative code generation step specified by the user. Making use of Nvidia’s NVVM framework allows us to code generate from a very similar LLVM structure as the CPU-targeted code generation. Because CUDA kernels are inherently parallel, a GPU-specific lowering pass is performed to translate loops over global sets into parallel kernel structures. Broadly, to convert global for-loops into parallel structures, reduction operations are turned into atomic operations and the loop variable is replaced with the CUDA thread ID. Following this, we perform a GPU-specific analysis to fuse these parallel loops wherever possible to reduce kernel launch overhead and increase parallelism. Using a very similar pipeline for CPU and GPU code generation helps us ensure that behavior on the GPU and CPU are identical for the code structures we generate.

## 8. RESULTS

To evaluate Simit, we implemented three realistic simulation applications, Implicit Springs, Neo-Hookean FEM and Elastic Shells (Section 8.1), using Simit, Matlab and Eigen. In addition, we compared to SOFA and Vega, two hand-optimized state-of-the-art real-time physics engines (Section 8.2). Note that Vega did not support Implicit Springs, and neither Vega nor SOFA supported Elastic Shells. We then conducted three experiments that show that:

With the traditional approaches we evaluated you get better performance by writing more code. With Simit you can get both performance and productivity (Section 8.3)

You can compile a Simit program to GPUs with no change to the source code to get about 10× more performance. (Section 8.4)



Fig. 11: Still from an Implicit Springs simulation of the Stanford bunny. The bunny consists of 36,976 vertices and 220,147 springs, which can be simulated by Simit on a GPU at 12 frames per second. (Only surface vertices and springs are shown.) The Simit code is only 93 lines, which includes a conjugate gradient solver implementation.

Simit scales well when the data set size increases (Section 8.5)

All CPU timings are taken on an Intel Xeon E5-2695 v2 at 2.40GHz with 128 GB of memory running Linux. All CPU measurements are single-threaded—none of the libraries we compare to support multi-threaded CPU execution, nor does the current Simit compiler, though parallelization will be added in the future. Simit and cuSPARSE GPU timings are taken on an Nvidia Titan GK110.

### 8.1 Applications

We implemented three simulation applications with different edge topologies and computational structure, and paired each with a suitable data set (bunny, dragon and cloth).

For Implicit Springs and Neo-Hookean FEM we chose to implement the CG solver in Simit instead of using an external solver library. The reason for this is that Simit offers automatic portability to GPUs, natively compiles SPMV operations to be blocked with a dense inner loop (see Section 9), and because this avoids data translation going from Simit to the external solver library. To evaluate the performance benefit of performing CG in Simit, we ran an experiment where Simit instead used Eigen to perform the CG solve. This resulted in a 30% slowdown, due to data translation and because Eigen does not support the Blocked CSR format.

**8.1.1 Implicit Springs.** Our first example is a volumetric elasticity simulation using implicit springs. We tetrahedralized the Stanford bunny to produce 37K vertices and 220K springs, and passed this to Simit as a vertex set connected by an edge set. Our implementation uses two assembly maps—one on the vertex set to compute the mass and damping matrices, and one on the edge set to compute the stiffness matrix. To solve for new vertex velocities, we implement a linearly-implicit time stepper and use the method of conjugate gradients (CG) as our linear solver. One benefit of CG, as an indirect solver, is that it can be implemented without materializing the system matrix  $A$ : only a function to compute  $Ax$  and  $A^T x$  given an arbitrary vector  $x$  is needed. For this reason, highly-optimized CG libraries require the user to implement a pair of callbacks rather than passing in the matrix. Simit offers the same benefit without the hassle: the user constructs and passes the  $A$  matrix to a straightforward implementation of CG that can be implemented in Simit. By analyzing data flow Simit can automatically avoid materializing it if possible. Per common practice when implementing a performance-sensitive

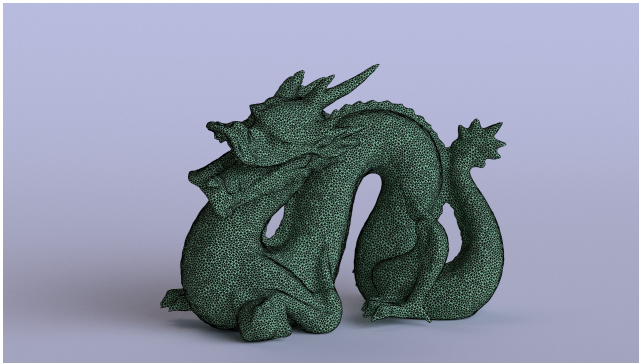


Fig. 12: Still from a Tetrahedral FEM (Finite Element Method) simulation of a dragon with 46,779 vertices and 160,743 elements, using the Neo-Hookean material model. Simit performs the simulation at 11 frames per second with only 154 non-comment lines of code shown in Appendix A. This includes 15 lines for global definitions, 14 lines for utility functions, 69 lines for local operations, 23 lines to implement CG and 13 lines for the global linearly-implicit timestepper procedure to implement the simulation, as well as 20 lines to precompute tet shape functions.

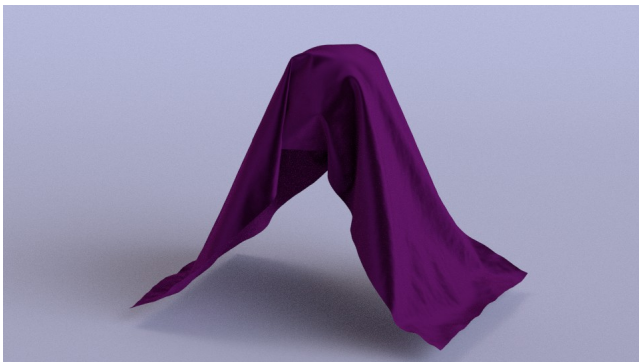


Fig. 13: Still from an Elastic Shells simulation of a cloth with 499,864 vertices, 997,012 triangle faces and 1,495,518 hinges. Elastic shells require two hyperedge sets: one for the triangle faces and one for the hinges. The Simit implementation ran at 15 frames per second on a GPU.

implicit solver, we limit (in all three implementations) the maximum number of conjugate gradients iterations (we choose 50).

**8.1.2 Neo-Hookean FEM.** Our second example is one of the most common methods for animating deformable objects, tetrahedral finite elements with linear shape functions, in Simit. We use the non-linear Neo-Hookean material model [Mooney 1940], as it is one of the standard models for animation and engineering, and we set the stiffness and density of the model to realistic physical values. The Simit implementation uses three maps, one to compute forces on each element, one to build the stiffness matrix and one to assemble the mass matrix, and then solves for velocities. We then use the conjugate gradient method to solve for the equations of motion, again relying on an implementation in pure Simit.

**8.1.3 Elastic Shells.** As a final example, we implemented an elastic shell code and used it to simulate the classic example of a rectangular sheet of cloth draping over a rigid, immobile sphere. The input geometry is a triangle mesh with 500K vertices, and is encoded as a hypergraph using one vertex set (the mesh vertices) and two hyperedge sets: one for the triangle faces, which are the stencil for a constant-strain Saint Venant-Kirchhoff stretching force; and one

for pairs of triangles meeting at a *hinge*, the stencil for the Discrete Shells bending force of Grinspun et al [Grinspun et al. 2003]. The bending force is a good example of the benefit of specifying force computation as a function mapped over an edge set: finding the two neighboring triangles and their vertices given their common edge, in the correct orientation, typically involves an intricate and bug-prone series of mesh traversals. Simit provides the bending kernel with the correct local stencil automatically. The Simit implementation uses a total of five map operations over the three sets to calculate forces and the mass matrix before updating velocities and positions for each vertex using explicit Velocity Verlet integration.

## 8.2 Languages and Libraries

We implemented each application in Matlab and in C++ using the Eigen high-performance linear algebra library. In addition, we used the SOFA simulation framework to implement Implicit Springs and Neo-Hookean FEM, and the Vega FEM simulator library to implement Neo-Hookean FEM.

**8.2.1 Matlab.** Matlab is a high-level language that was developed to make it easy to program with vectors and matrices [MATLAB 2014]. Matlab can be seen as a scripting language on top of high-performance linear algebra operations implemented in C and Fortran. Even though Matlab’s linear algebra operations are individually very fast, they don’t typically compose into fast simulations. The main reasons for this is Matlab’s high interpretation overhead, and the fact that individually optimized linear algebra foregoes opportunities for fusing operations (see Sections 7 and 9).

**8.2.2 Eigen.** Eigen is an optimized and vectorized linear algebra library written in C++ [Guennebaud et al. 2010]. To get high performance it uses template meta-programming to produce specialized and vectorized code for common operations, such as  $3 \times 3$  matrix-vector multiply. Furthermore, Eigen defers execution through its object system, so that it can fuse certain linear algebra operations such as element-wise addition of dense vectors and matrices.

**8.2.3 SOFA.** SOFA is an open source framework, originally designed for interactive, biomechanical simulation of soft tissue [Faure et al. 2007]. SOFA’s design is optimized for use with iterative solvers. It uses a scene graph to model the interacting objects in a simulation and, during each solver iteration, visits each one in turn, aggregating information such as applied forces. Using this traversal SOFA avoids forming global sparse matrices which is the key to its performance.

**8.2.4 Vega.** Vega is a graphics-centric, high-performance finite element simulation package [Sin et al. 2013]. Vega eschews special scene management structures in favor of a general architecture that can be used with iterative or direct solvers. It achieves high performance using optimized data-structures and clever rearrangement of material model computations to reduce operation count.

## 8.3 Simit Productivity and Performance

The general trend in the traditional systems we evaluate is that you get more performance by writing more code. Table I shows this for the three applications from Section 8.1. For each application and each language/library we report the performance (seconds per frame), source lines and memory consumption. For two applications we vectorized the Matlab code to remove all the loops (Matlab Vec). This made the Matlab code about one order of magnitude faster, but took 9 and 16 hours of additional development time for domain experts who are proficient with Matlab, and made the code very hard to read.

Table I. : Comparison of three applications implemented with Matlab, Vectorized Matlab, Eigen, hand-optimized C++ (SOFA and Vega) and Simit, showing the productivity and performance of Simit. For example, the Simit Neo-Hookean FEM is just 180 lines of code (includes CG solver), but simulates 160,743 tetrahedral elements in just 1.15 seconds using a single non-vectorized CPU thread. For each implementation we report the non-comment source lines of code, seconds per frame and peak memory in megabytes ( $1024^2$  bytes), as well as the size of each number relative to Simit. The trend is that you can get better performance by writing more code, however, with Simit you get both performance and productivity. For example, the Implicit Springs Simit implementation is shorter than any other implementation at 93 lines of code, yet runs faster at 0.6 seconds per frame. Matlab ran out of memory running the cloth simulation. The applications were run on the bunny, dragon and cloth data sets respectively with **double precision floating point** on an Intel Xeon E5-2695 v2 running at 2.40GHz with 128 GB of memory.

		ms per frame		Source lines		Memory (MB)	
Implicit Springs	Matlab	13,280	23.7×	142	1.5×	1,059	6.1×
	Matlab Vec	2,115	3.8×	230	2.5×	1,297	7.5×
	Eigen	883	1.6×	314	3.4×	339	2.0×
	SOFA	588	1.1×	1,404	15.1×	94	0.5×
	Simit	559	1.0×	<b>93</b>	1.0×	173	1.0×
Neo-Hookean FEM	Matlab	207,538	181.3×	234	1.3×	1,564	12.8×
	Matlab Vec	16,584	14.5×	293	1.6×	53,949	442.2×
	Eigen	1,891	1.7×	363	2.0×	626	5.1×
	SOFA	1,512	1.3×	1,541	8.6×	324	1.8×
	Vega	1,182	1.0×	1,080	6.0×	614	5.0×
	Simit	1,145	1.0×	<b>180</b>	1.0×	173	1.0×
Elastic Shells	Matlab			203	1.1×	OOM	
	Simit	600	1.0×	<b>190</b>	1.0×	785	1.0×
	Eigen	489	0.8×	453	2.4×	354	0.5×

For example, the Implicit Springs Eigen implementation requires 60% more code than the Matlab implementation, but runs almost 15 times faster. In general, higher performance meant more code had to be written in a lower-level language. Simit, however, breaks this tradeoff and produces high performance with few lines of code. In every case the Simit implementation is within 20% of the performance of the fastest alternative we found, while being fewer lines of code than Matlab. For example, Simit obtains equivalent performance to Vega, an optimized FEM library (Section 8.2.4), but requires 6× fewer lines of code, is more general, and compiles to GPUs for 10× more performance (Section 8.4). Furthermore, we plan to vectorize and multi-thread Simit in the future, which will further increase its performance.

For the Elastic Shells, Simit is 18% slower than Eigen. This is due to the application’s use of an explicit method to simulate the cloth, which is amenable to Eigen’s vectorized dense vector routines. Simit, on the other hand, does not yet implement vectorization. However, by running on the GPU, Simit can drastically increase performance on this code due to its large number of small dense operations (see Section 8.4) while still using code that is less than half the size of the Eigen implementation.

**8.3.1 Memory Usage.** While Matlab enables productive experimentation, its memory usage relative to Eigen and Simit is quite large; in some cases, Matlab uses an order of magnitude more memory. Interestingly, vectorizing the Matlab code can drastically increase memory usage, as in the case of the Neo-Hookean FEM example, where memory usage increased by 35×. In two applications, Simit is much more memory-efficient than Eigen due to its single representation for both the graph and matrices. Table I shows that SOFA’s strategy of not building system matrices does indeed reduce memory usage. In Discrete Elastic Shells, Simit consumes about double the memory that Eigen uses.

**8.3.2 Compilation and Initialization.** Simit applications generate and compile code during two operations: Simit code is just-in-time compiled from source prior to the first time it is executed, and graph-specific code is generated to bind datasets to the appro-

Table II. : Compilation and initialization times for Simit applications, in milliseconds. Simit code is just-in-time compiled prior to starting the timestepper for the first time. Initialization time measures how long it takes to construct and compile code passing in the graph prior to a timestep. The applications were run on the bunny, dragon and cloth data sets respectively with **double precision floating point** on an Intel Xeon E5-2695 v2 running at 2.40GHz with 128 GB of memory.

	Compilation (ms)	Initialization (ms)
Implicit Springs	37	263
Neo-Hookean FEM	52	48
Elastic Shells	38	31

Table III. : Comparison of Simit applications running on a CPU and a GPU, and a cuSPARSE hybrid CPU-GPU implementation. The cuSPARSE FEM implementation was hand written. The applications were run on the bunny, dragon and cloth data sets respectively with **single precision floating point**. GPU measurements were taken on an Nvidia Titan GK110. With a GPU Simit achieves interactive rates on these data sets: 13 fps, 9 fps and 59 fps.

		ms per frame		Source lines	
Implicit Springs	Simit CPU	431	1.0×	93	1.0×
	Simit GPU	90	0.2×	93	1.0×
Neo-Hookean FEM	Simit CPU	944	1.0×	180	1.0×
	cuSPARSE	801	0.8×	464	2.6×
	Simit GPU	93	0.1×	180	1.0×
Elastic Shells	Simit CPU	427	1.0×	190	1.0×
	Simit GPU	17	0.04×	190	1.0×

appropriate pointers prior to calling the timestepper. Table II shows the compilation and initialization times for the three applications. In all three applications, Simit program compilation and initialization code compilation takes a fraction of a second.

## 8.4 Simit GPU Execution

Simit programs compile to GPUs with no change to the source code. We compiled all three applications to run on an Nvidia GPU, and the resulting code executed 5-25× faster than the CPU implementation. Shown in Table III is a comparison of execution times and lines of

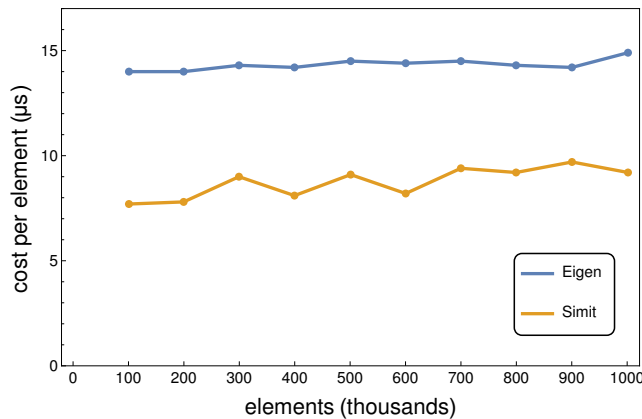


Fig. 14: Scaling performance of Eigen and Simit implementations of Neo-Hookean Tetrahedral FEM, as the number of elements in the dragon mesh are increased from approximately 100,000 elements to 1 million elements.

code for the single-precision implementations of each application. We saw the largest gains in the most computationally intense application, the Elastic Shells application, because GPUs are designed to excel at floating point operations over small dense vectors. For Elastic Shells, the runtime was improved by  $25\times$  over GPU Simit.

GPU Simit also performs on par with an efficient hand-written CUDA implementation. To compare against hand-written CUDA code, the Eigen code for the Neo-Hookean FEM example was rewritten into a hybrid CPU-GPU implementation. The matrix assembly was performed using the Eigen library, using code identical to the full Eigen implementation. Using the CUDA cuSPARSE library to perform fast sparse matrix operations, the resulting matrix was then passed to the GPU to be quickly solved via the conjugate gradient method. This hybrid code approach is common for simulation developers to gain performance from a GPU by taking advantage of available libraries. At the cost of roughly  $2.5\times$  the lines of code, the CG solve step, previously half the computation time of the Eigen solution, was improved to 48ms. The GPU Simit code achieves the same speed for the CG solve step, with fewer lines of code. Due to Amdahl’s Law, the hybrid implementation is limited in the speedup it can achieve, since the assembly step (which is roughly half of the runtime) is not sped up by cuSPARSE. As a result, the Simit GPU implementation still significantly outperforms this hybrid version.

## 8.5 Simit Scaling to Large Datasets

Figure 14 shows the performance of Eigen and Simit as the number of elements are increased in the Neo-Hookean Tetrahedral FEM application. In both cases, the implementations scale linearly in the number of elements, as we expect, given that both of the major computational steps (assembly and CG with a capped maximum number of steps) will scale linearly. Across the various sizes, Eigen averages 14us per element per timestep, while Simit averages 8.7us per element per timestep.

## 9. REASONS FOR PERFORMANCE

Simit’s performance comes from its design and is made possible by the choice of abstract data structures (hypergraphs and blocked tensors), and the choice of collection-oriented operations on these (stencils, tensor assemblies and index expressions). Little effort has so far gone into low-lever performance engineering. For example, Simit does not currently emit vector instructions, nor does it optimize memory layout order or run in parallel on multi-core machines,

but we plan to add this in the future. Specifically, there are three keys to the performance of Simit presented in this article:

*In-place Computation.* The assembly map construct unifies the hypergraph with computations expressed as index expressions on system tensors. This lets the compiler reason about where tensors come from, and thus how their non-empty values are distributed relative to the graph. Simit uses this information to pin every value in a system tensor to a graph node or edge, which lets it schedule computation in-place on the graph. In-place computation allows efficient matrix assembly, since the graph itself becomes the sparse index structure of matrices. Matrix values can be placed directly in their storage location without the need to construct and compress a dedicated sparse index structure. Traditional sparse matrix libraries require separate insertion and compression steps; in Simit, these are replaced by the map construct, and compression disappears. In some cases it is even possible to completely remove matrix assembly by fusing it with the first computation that uses it.

Furthermore, knowledge about matrix assembly even lets the compiler fuse sequences of sparse matrix operations, as it knows when sparse matrices have the same or overlapping values distributions. For example, adding two sparse matrices that come from the same graph becomes a local operation where each graph node adds the matrix values it owns for both matrices without any communication. These kinds of static optimizations on sparse matrices have not been possible before, as explained by Vuduc et al.: “while sparse compilers could be used to provide the underlying implementations of sparse primitives, they do not explicitly make use of matrix structural information available, in general, only at run-time.” [Vuduc et al. 2005]. The knowledge of matrix assembly provided by Simit’s assembly map breaks the assumptions underlying this assertion and opens up for powerful static sparse optimizations.

Finally and importantly, in-place computation makes efficient parallel computation on massively parallel GPUs straightforward by assigning graph nodes and edges to parallel threads and computing using the owner-computes rule. This works well because given an assignment, Simit also knows how to parallelize system vector and matrix operations. Furthermore, these are parallelized in a way that matches the parallelization of previous stages thereby reducing synchronization and communication. Without the in-place representation, parallel code generation would not be as effective.

*Index Expression Fusion.* Tensor index expressions are a powerful internal representation of computation that simplify make program transformation. Index expressions are at once simpler and more general than linear algebra. Their power and uniformity makes it easy for the compiler to perform transformations like tensor operation fusion to remove tensor temporaries and the resulting memory bandwidth costs. Moreover, it can do these optimizations without the need to build in the many rules of linear algebra.

*Dense Block Computation.* Natively blocked tensors and hierarchical indexing not only simplify the programmer’s code, they also make it trivial for the compiler to use blocked representations. Blocked matrix representations result in efficient dense or unrolled inner loops within sparse computations such as SpMV. This greatly reduces the need to look up values through an index.

## 10. CONCLUSIONS

A key insight in Simit is that the best abstraction for describing a system’s structure is different from the best abstraction for describing its behavior. Sparse systems are naturally graph-structured, while behavior is often best described using linear algebra over the whole

system. Simit is a new programming model that takes advantage of this duality, using assembly maps to bridge between the abstractions. Using information about system sparsity, combined with a new representation of operations as index expressions, Simit compiles to fast code while retaining the expressibility of high-level languages like MATLAB. With the ability to run programs on CPUs and GPUs, Simit attains an unprecedented level of performance portability.

We believe Simit has the potential to obtain higher performance while retaining its expressibility. So far, our implementation has only scratched the surface of what kinds of optimizations are possible with assemblies and index expressions. Future work will extend our optimization strategies for manipulating index expressions resulting from linear algebra operations. Furthermore, we have not yet implemented parallelization or vectorization of CPU code, which can provide further factors of speedup. Finally, distributed and hybrid code generation is possible given the Simit abstractions and will further improve performance.

Simit lets programmers write code at a high level and get the performance of optimized low-level code. Simit enables MATLAB-like productivity with the performance of manually optimized C++ code.

## APPENDIX

### A. NEO-HOOKEAN FINITE ELEMENT METHOD

We show a complete implementation of a finite element method with linear tetrahedral elements. Our implementation includes the constitutive model, the assembly stage of forces and stiffness matrices, and a linearly-implicit dynamics integrator. In the example, we implemented the Neo-Hookean material model. Different material models can be plugged in by changing the stress and stress differential functions. Next, our assembly code only defines how to locally compute stiffness and forces for a single element. The global assembly is handled by Simit. Lastly, we show an implementation of a linearly-implicit time integrator with a conjugate gradient linear solver. The time stepper is written purely in terms of linear algebra. It is agnostic of the underlying finite element structures and therefore can be applied to different element types.

```
const grav = [0.0, -10.0, 0.0];

element Tet
  u : float;
  l : float;
  W : float;
  B : tensor[3,3](float);
end

element Vert
  x : tensor[3](float);
  v : tensor[3](float);
  c : int;
  m : float;
end

extern verts : set{Vert};
extern tets : set{Tet}(verts, verts, verts, verts);

% precompute volume and shape function gradient for tets
func precomputeTetMat(inout t : Tet, v : (Vert*4))
  -> (m:tensor[verts](float))
  var M:tensor[3,3](float);
  for ii in 0:3
    for jj in 0:3
      M(jj,ii) = v(ii).x(jj) - v(3).x(jj);
    end
  end
  t.B = inv(M);
```

```
vol = -(1.0/6.0) * det(M);
t.W = vol;

rho = 1000.0;
for ii in 0:4
  m(v(ii))=0.25*rho*vol;
end
end

proc initializeTet
  m = map precomputeTetMat to tets;
  verts.m = m;
end

% first Piola-Kirchoff stress
func PK1(u : float, l : float, F : tensor[3,3](float))
  -> (P : tensor[3,3](float))
  J = log(det(F));
  Finv = inv(F)';
  P = u*(F-Finv) + l*J*Finv;
end

% gradient of first Piola-Kirchoff stress
func dPdF(u : float, l : float, F : tensor[3,3](float),
  dF : tensor[3,3](float))
  -> (dP : tensor[3,3](float))
  J = log(det(F));
  Fi = inv(F);
  FidF = Fi*dF;
  dP = u*dF + (u-l*J) * Fi' * FidF' + l*trace(FidF)*Fi';
end

% assemble lumped mass matrix and gravitational force
func compute_mass(v : Vert)
  -> (M : tensor[verts,verts](tensor[3,3](float)),
  fg : tensor[verts](tensor[3](float)))
  M(v,v) = v.m * I;
  if(v.c <= 0)
    fg(v) = v.m * grav;
  end
end

% assemble internal forces, fixed vertices contribute
% no force
func compute_force(e : Tet, v : (Vert*4))
  -> (f : tensor[verts](tensor[3](float)))
  var Ds : tensor[3,3](float);
  for ii in 0:3
    for jj in 0:3
      Ds(jj,ii) = v(ii).x(jj) - v(3).x(jj);
    end
  end
  F = Ds*e.B;

  P = PK1(e.u, e.l, F);
  H = -e.W * P * e.B';

  for ii in 0:3
    fi = H(:,ii);

    if (v(ii).c <= 0)
      f(v(ii)) = fi ;
    end

    if (v(3).c <= 0)
      f(v(3)) = -fi;
    end
  end
end

% assemble stiffness matrix, fixed vertices contribute
% nothing to stiffness matrix
func compute_stiffness(e : Tet, v : (Vert*4))
  -> (K : tensor[verts,verts](tensor[3,3](float)))
  var Ds : tensor[3,3](float);
  var dFRow : tensor[4,3](float);
  m = 0.01;
```



```

for ii in 0:3
  for jj in 0:3
    Ds(jj,ii) = v(ii).x(jj)-v(3).x(jj);
  end
end

F = Ds*e.B;
for ii in 0:3
  for ll in 0:3
    dFRow(ii,ll) = e.B(ii,ll);
  end
  dFRow(3, ii) = -(e.B(0, ii)+e.B(1, ii)+e.B(2, ii));
end

for row in 0:4
  var Kb : tensor[4,3,3](float) = 0.0;
  for kk in 0:3
    var dF : tensor[3,3](float) = 0.0;
    for ll in 0:3
      dF(kk, ll) = dFRow(row, ll);
    end
    dP = dPdF(e.u, e.l, F, dF);
    dH = -e.W * dP * e.B';

    for ii in 0:3
      for ll in 0:3
        Kb(ii,ll, kk) = dH(ll, ii);
      end
      Kb(3, ii, kk) = -(dH(ii, 0)+dH(ii, 1)+dH(ii, 2));
    end
  end

  for jj in 0:4
    if(v(jj).c <= 0) and (v(row).c <= 0)
      K(v(jj), v(row)) = Kb(:, :, jj);
    end
  end
end
end

% conjugate gradient with no preconditioning.
func CG(A : tensor[verts,verts](tensor[3,3](float)),
      b : tensor[verts](tensor[3](float)),
      x0 : tensor[verts](tensor[3](float)),
      tol : float, maxIter : int)
-> (x : tensor[verts](tensor[3](float)))
r = b - (A*x0);
p = r;
iter = 0;
x = x0;
normr2 = dot(r, r);
while (normr2 > tol) and (iter < maxiters)
  Ap = A * p;
  denom = dot(p, Ap);
  alpha = normr2 / denom;
  x = x + alpha * p;
  normr2old = normr2;
  r = r - alpha * Ap;
  normr2 = dot(r, r);
  beta = normr2 / normr2old;
  p = r + beta * p;
  iter = iter + 1;
end
end

% linearly-implicit time stepper with CG solver
proc main
  h = 0.01;
  tol = 1e-12;
  maxiters = 50;

  M,fg = map compute_mass to verts reduce +;
  f = map compute_force to tets reduce +;
  K = map compute_stiffness to tets reduce +;

  A = M - (h*h) * K;
  b = M*verts.v + h*(f+fg);

```

```

x0 = verts.v;
verts.v = CG(A, b, x0, tol, maxIter);
verts.x = h * verts.v + verts.x;
end

```

## REFERENCES

- CGAL, Computational Geometry Algorithms Library. <http://www.cgal.org>.
- BAUMGART, B. G. 1972. Winged edge polyhedron representation. Tech. rep., Stanford University.
- BLELLOCH, G. E. 1990. *Vector models for data-parallel computing*. Vol. 356. MIT press Cambridge.
- BLYTHE, D. 2006. The direct3d 10 system. *ACM Transactions on Graphics* 25, 3 (July), 724–734.
- BOTSCH, M., STEINBERG, S., BISCHOFF, S., AND KOBBELT, L. 2002. Openmesh – a generic and efficient polygon mesh data structure.
- CODD, E. F. 1970. A relational model of data for large shared data banks. *Communications of the ACM* 13, 6, 377–387.
- COMSOL, A. 2005. Comsol multiphysics users guide. *Version: September*.
- COUMANS, E. ET AL. 2006. Bullet physics library. *Open source: bullet-physics.org* 4, 6.
- DEVITO, Z., JOUBERT, N., PALACIOS, F., OAKLEY, S., MEDINA, M., BARRIENTOS, M., ELSEN, E., HAM, F., AIKEN, A., DURAISAMY, K., DARVE, E., ALONSO, J., AND HANRAHAN, P. 2011. Liszt: A domain specific language for building portable mesh-based pde solvers. In *Proceedings of 2011 International Conference for High Performance Computing, Networking, Storage and Analysis*. SC '11. ACM, New York, NY, USA, 9:1–9:12.
- DICK, C., GEORGII, J., AND WESTERMANN, R. 2011. A real-time multi-grid finite hexahedra method for elasticity simulation using cuda. *Simulation Modelling Practice and Theory* 19, 2, 801–816.
- DUBEY, P., HANRAHAN, P., FEDKIW, R., LENTINE, M., AND SCHROEDER, C. 2011. Physbam: Physically based simulation. In *ACM SIGGRAPH 2011 Courses*. SIGGRAPH '11. ACM, New York, NY, USA, 10:1–10:22.
- EASTMAN, C. AND WEISS, S. 1982. Tree structures for high dimensionality nearest neighbor searching. *Information Systems* 7, 2, 115–122.
- EINSTEIN, A. 1916. The Foundation of the General Theory of Relativity. *Annalen der Physik* 354, 769–822.
- ELLIOTT, C. 2001. Functional image synthesis. In *Proceedings of Bridges*.
- FAURE, F., ALLARD, J., COTIN, S., NEUMANN, P., BENSOUSSAN, P.-J., DURIEZ, C., DELINGETTE, H., AND GRISONI, L. 2007. SOFA: A modular yet efficient simulation framework. In *Surgetica 2007 - Computer-Aided Medical Interventions: tools and applications*, P. Merloz and J. Troccaz, Eds. Surgetica 2007, Gestes médicaux chirurgicaux assistés par ordinateur. Chambéry, France, 101–108.
- GRINSPUN, E., HIRANI, A. N., DESBRUN, M., AND SCHRÖDER, P. 2003. Discrete shells. In *Proceedings of the 2003 ACM SIGGRAPH/Eurographics Symposium on Computer Animation*. SCA '03. Eurographics Association, Aire-la-Ville, Switzerland, Switzerland, 62–67.
- GUENNEBAUD, G., JACOB, B., ET AL. 2010. Eigen v3. <http://eigen.tuxfamily.org>.
- GUENTER, B. AND LEE, S.-H. 2009. *Symbolic Dynamics and Geometry*. AK Peters.
- GUIBAS, L. AND STOLFI, J. 1985. Primitives for the manipulation of general subdivisions and the computation of voronoi diagrams. *ACM Trans. Graph.* 4, 2, 74–123.
- HAMMING, R. 2003. History of computer—software. In *Art of Doing Science and Engineering: Learning to Learn*. CRC Press.
- HANRAHAN, P. AND LAWSON, J. 1990. A language for shading and lighting calculations. In *Computer Graphics (Proceedings of SIGGRAPH 90)*. 289–298.

- HECHT, F. 2015. private communication.
- HIBBETT, KARLSSON, AND SORENSEN. 1998. *ABAQUS/standard: User's Manual*. Vol. 1. Hibbitt, Karlsson & Sorensen.
- HOLZMANN, G. 1988. *Beyond Photography*. Prentice Hall.
- JIWON SEO, STEPHEN GUO, M. S. L. 2013. Socialite: Datalog extensions for efficient social network analysis. In *IEEE 29th International Conference on Data Engineering*.
- KOHNKE, P. 1999. *ANSYS theory reference*. Ansys.
- LIU, K. 2014. Dynamic animation and robotics toolkit.
- LOW, Y., GONZALEZ, J., KYROLA, A., BICKSON, D., GUESTRIN, C., AND HELLERSTEIN, J. M. 2010. Graphlab: A new parallel framework for machine learning. In *Conference on Uncertainty in Artificial Intelligence*.
- MARK, W. R., GLANVILLE, R. S., AKELEY, K., AND KILGARD, M. J. 2003. Cg: A system for programming graphics hardware in a c-like language. *ACM Transactions on Graphics* 22, 3 (July), 896–907.
- MATLAB. 2014. *version 8.3.0 (R2014a)*. The MathWorks Inc., Natick, Massachusetts.
- MCADAMS, A., ZHU, Y., SELLE, A., EMPEY, M., TAMSTORF, R., TERAN, J., AND SIFAKIS, E. 2011. Efficient elasticity for character skinning with contact and collisions. In *ACM Transactions on Graphics (TOG)*. Vol. 30. ACM, 37.
- MOONEY, M. 1940. A theory of large elastic deformation. *Journal of applied physics* 11, 9, 582–592.
- PINGALI, K., NGUYEN, D., KULKARNI, M., BURTSCHER, M., HASSAAN, M. A., KALEEM, R., LEE, T.-H., LENHARTH, A., MANEVICH, R., MÉNDEZ-LOJO, M., PROUNTZOS, D., AND SUI, X. 2011. The tao of parallelism in algorithms. In *Proceedings of the 32Nd ACM SIGPLAN Conference on Programming Language Design and Implementation. PLDI '11*. ACM, New York, NY, USA, 12–25.
- POMMIER, J. AND RENARD, Y. 2005. Getfem++, an open source generic c++ library for finite element methods.
- RAGAN-KELLEY, J., ADAMS, A., PARIS, S., LEVOY, M., AMARASINGHE, S., AND DURAND, F. 2012. Decoupling algorithms from schedules for easy optimization of image processing pipelines. *ACM Trans. Graph.* 31, 4 (July), 32:1–32:12.
- RICCI-CURBASTRO, G. AND LEVI-CIVITA, T. 1901. Mthodes de calcul diffrentiel absolu et leurs applications. *Mathematische Annalen* 54, 125–201.
- SEGAL, M. AND AKELEY, K. 1994. *The OpenGL graphics system: a specification*.
- SIN, F. S., SCHROEDER, D., AND BARBIČ, J. 2013. Vega: Non-linear fem deformable object simulator. In *Computer Graphics Forum*. Vol. 32. 36–48.
- SMITH, R. ET AL. 2005. Open dynamics engine.
- SUNGPACK HONG, HASSAN CHAFI, E. S. AND OLUKOTUN, K. 2012. Green-marl: A dsl for easy and efficient graph analysis. In *17th International Conference on Architectural Support for Programming Languages and Operating Systems*.
- VIDIMČE, K., WANG, S.-P., RAGAN-KELLEY, J., AND MATUSIK, W. 2013. Openfab: A programmable pipeline for multi-material fabrication. *ACM Transactions on Graphics (TOG)* 32, 4, 136.
- VUDUC, R., DEMMEL, J. W., AND YELICK, K. A. 2005. OSKI: A library of automatically tuned sparse matrix kernels. In *Proc. SciDAC, J. Physics: Conf. Ser.* Vol. 16. 521–530.

