

# Taichi: An Introduction

## Languages Seminar

Russell Bentley

Stony Brook

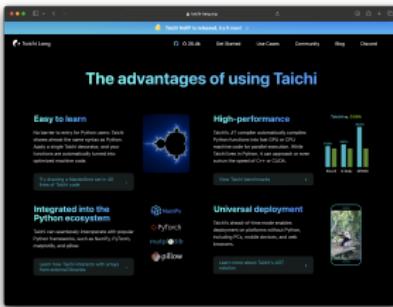
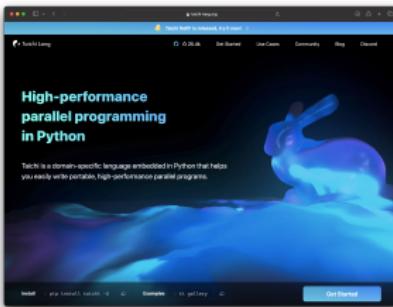
Thursday 3<sup>rd</sup> October, 2024



Stony Brook University

# Useful Links

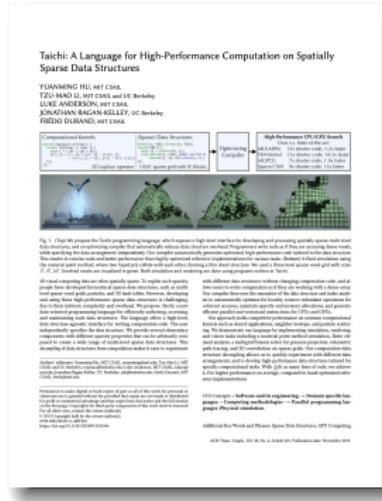
- ▶ Taichi homepage:  
<https://www.taichi-lang.org>
- ▶ Official docs / guides:  
<https://docs.taichi-lang.org>
- ▶ API reference:  
<https://docs.taichi-lang.org/api/>
- ▶ Github:  
<https://github.com/taichi-dev/taichi>



Stony Brook University

# Taichi?

- ▶ Large and Growing Project
  - ▶ Language for compute kernels
  - ▶ Data structure DSL
  - ▶ Python frontend
  - ▶ Compiler Infrastructure
  - ▶ Focus on Graphics



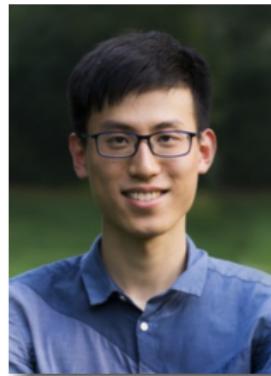
2019 [4]



Stony Brook University

# Taichi Background

- ▶ Taichi is Yuanming Hu's 2021 Thesis
- ▶ Growing collection of related papers



The Taichi High-Performance and Univerifiable Programming Language for Sparse and Quantized Visual Computing  
by  
Yuanming Hu

Submitted to the Department of Electrical Engineering and Computer Science  
on May 18, 2021, in partial fulfillment of the  
Degree of Philosophy in Computer Science

#### Abstract

Using traditional programming languages such as C++ and CUDA, writing high-performance sparse and quantized visual computing programs is error-prone and performance-critical. This implies an inextricable trade-off between performance and correctness. In this thesis, we propose Taichi, a high-level language for sparse and quantized visual computing. It supports sparse and quantized tensors, data structure operators, differentiable programming, and quantized computation. It also provides a high-level interface for memory management and parallelism. To address these issues, we propose Taichi, an imperative and parallel programming model, which is based on the concept of “univerifiable” and “verifiable” computation. Taichi leverages domain-specific features of visual computing tasks, provides a high-level interface for memory management, and supports both imperative and parallel programming, and quantization. With Taichi’s optimization, the final generated code is as efficient as hand-tuned CUDA programs, yet it is much easier to write. We demonstrate its expressiveness by implementing a state-of-the-art sparse convolutional neural network. We also demonstrate its efficiency by implementing a sparse matrix multiplication using 1/10 times of code as sparse computation, efficiently develop 3D deformable registration, and efficiently simulate 20 million material point method (MPM) particles on a single GPU.

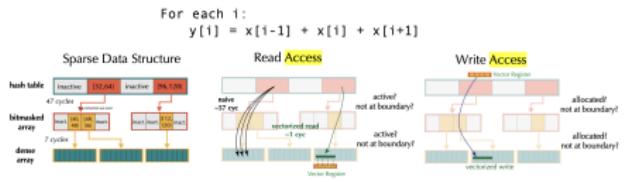
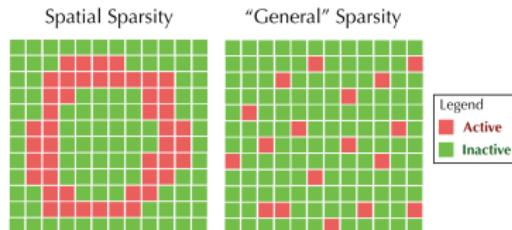
Thesis Supervisor: Frédo Durand  
Title: Professor of Electrical Engineering and Computer Science

Thesis Supervisor: William T. Freeman  
Title: Professor of Electrical Engineering and Computer Science

3

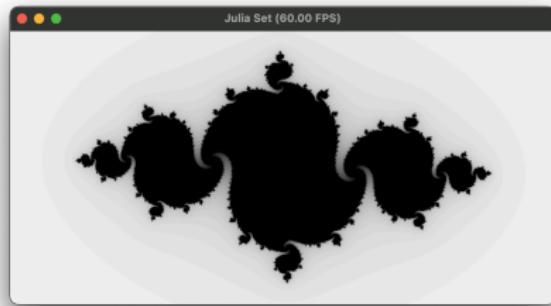
# Design Goals

- ▶ Explore data representation
  - ▶ Sparse data structures
  - ▶ Quantized data types
- ▶ Decouple data and algorithms
- ▶ Differentiable kernels
- ▶ Optimize access patterns
- ▶ Optimize parallelism



```
is = ti.quant.int(bits=5)
u19 = ti.quant.int(bits=19, signed=False)
```

# Interactive Example

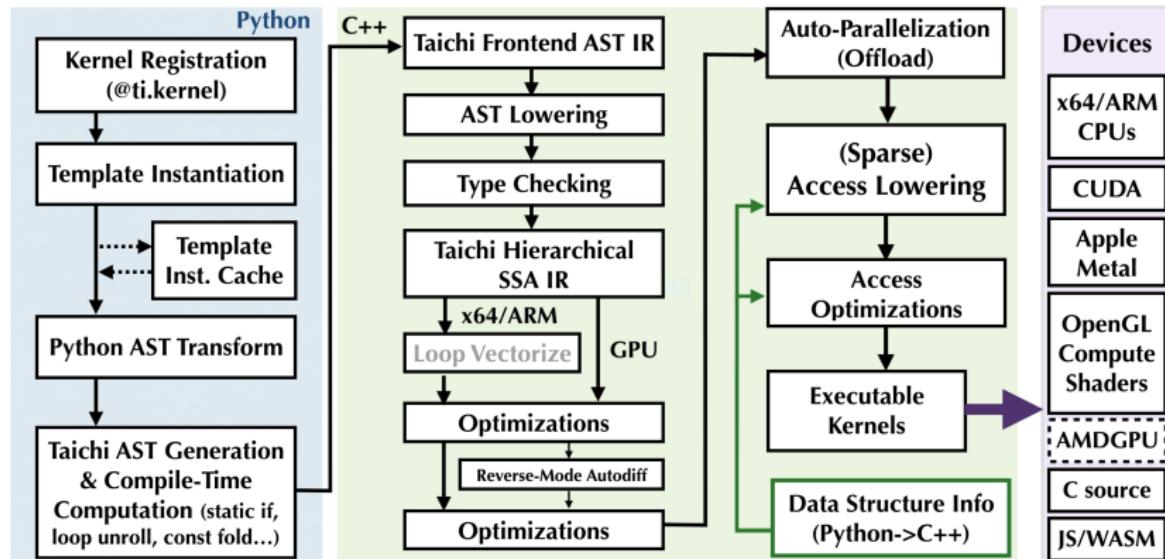


Based on [Hello, World!](#) guide



Stony Brook University

# Life of a Taichi Kernel



Stony Brook University

# Why use a Python Frontend?

## Taichi's old C++ frontend

Taichi used to be embedded in C++14. However, that solution is mostly abandoned, because

- ① C++ itself is too complex for most users to learn, not to say a DSL embedded in C++. E.g., [►math/svd.h](#) 😞
- ② Getting C++ AST is almost impossible. We had to heavily use templates/macros tricks, which harm readability. 😞

## Taichi's new Python frontend

Now the whole Taichi system is deeply embedded in Python.

- ① Python is easy to learn and widely adopted. 😊
- ② Python allows flexible AST inspection and manipulation. 😊

Both of these allow us to invent a new high-performance language out of Python.



Stony Brook University



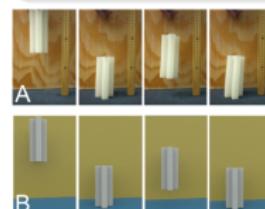
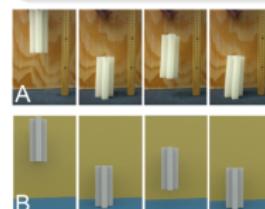
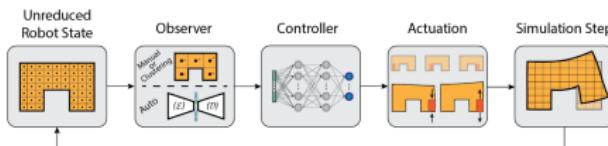
# Applications: Differentiable Physics

## ► Soft Robotics

- Difficult to design
- Difficult to control

## ► Simulation

- Sparse representation
- Differentiable for learning / optimizing



2019 [2]

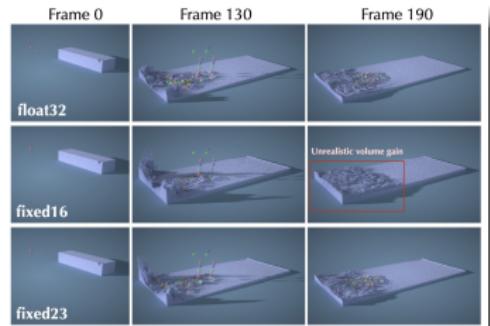
2019 [8]

# Applications: Quantized Graphics

- ▶ Decouple algorithms from data representation
- ▶ Lower memory requirements
  - ▶ 8x Game of Life
  - ▶ 1.5x Eulerian Fluid
  - ▶ 1.7x Material point method
- ▶ Leads to faster runtimes

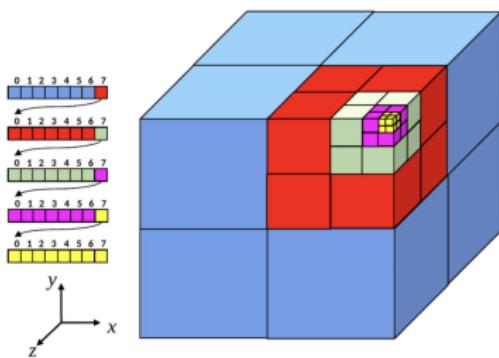


2021 [3]



# Sparse Data Structures: Octree

- ▶ Easy to implement
- ▶ Hard to write access patterns
- ▶ Not cache friendly
- ▶ Can only represent data in a bounded space



# OpenVDB

- ▶ Originally from Dreamworks animation
- ▶ Published in ACM Trans. Graphics
  - ▶ Describes VDB Datastructure
  - ▶ Introduced OpenVDB implementation
- ▶ Widely used for special effects and animation
  - ▶ Supported by industry standard software like Houdini and Blender
- ▶ Often used with implicit geometry

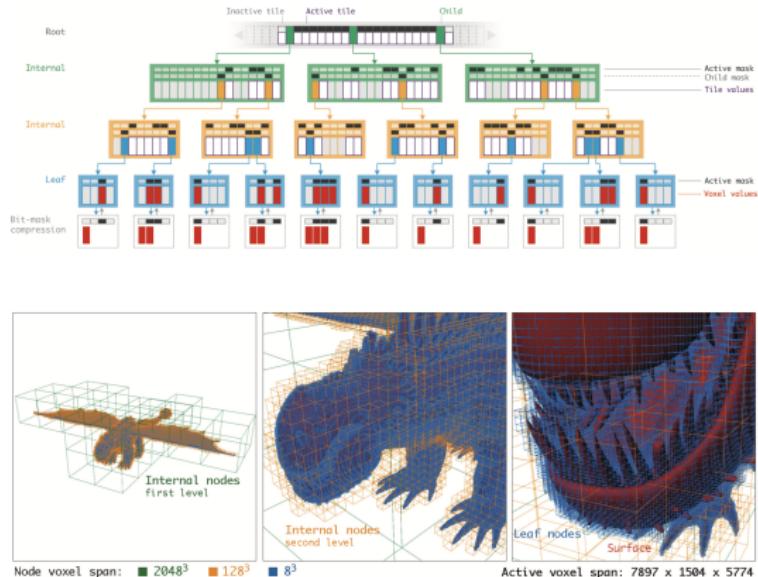


2013 [6]



# VDB Datastructure

- ▶ Similar to B+Trees
- ▶ Map based root
  - ▶ Dynamic branching factor
  - ▶ Virtually infinite spatial bounds
- ▶ Fixed height
- ▶ Large (fixed) branching factors for children



- ▶ ACM Trans. Graphics 2014 [7]
- ▶ Key ideas
  - ▶ Large finite grid replaces infinite map
  - ▶ Morton key as virtual memory address
  - ▶ Translation Lookahead Buffer (TLB) handles address mapping

**SPGrid: A Sparse Paged Grid structure applied to adaptive smoke simulation**

Reinhard Fidler, Michael Attwells, Ross Lasseter, Edouard Mélotin  
University of Wisconsin - Madison



Figure 3: Dense flow-particle field (10M particles), adaptive and sparse on the right.

**Abstract:**  
We present a new method for fluid simulation on high-resolution adaptive grids which yields the densification and partitioning properties of a sparse grid. Our approach is based on a hierarchical tree structure, called SPGrid, where a tree of nodes stores the current range and offset information for each leaf node. This allows for efficient construction and removal of active handles (handles are device specific pointers to memory locations). We also propose a novel way to represent particles using Morton keys, which are integer values that encode the position of a particle within a regular grid. This allows for efficient insertion and removal of particles from a sparse paged grid, and the resulting data structure is highly cache friendly. Finally, we demonstrate an adaptive multigrid solver for the Poisson equation that uses the SPGrid structure to handle phenomena with a modest memory footprint. Our method is competitive with state-of-the-art adaptive grid solvers in terms of both efficiency and memory usage. We demonstrate the efficacy of our method on two simulations of smoke flow.

CC Categories: I.3.3 [Computer Graphics]: Computational Geometry and Object Representations; I.3.7 [Computer Graphics]: Computational Geometry and Object Representations; I.3.8 [Computer Graphics]: Applications

Keywords: fluid simulation, sparse grids, adaptive multigrid

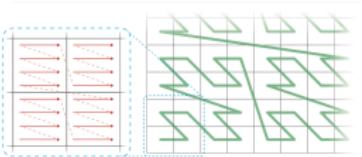
Link: <http://www.cs.wisc.edu/~rlas/SPGrid.pdf>

**1 Introduction**

Computer graphics research has explored opportunities for spatial indexing and partitioning of data structures to support efficient rendering of complex scenes, including adaptive octree grids (Lisicki et al. 2004; Liu et al. 2005; Saito et al. 2005; Saito et al. 2006; Saito et al. 2007; Saito et al. 2008; Saito et al. 2009; Saito et al. 2010; Saito et al. 2011; Saito et al. 2012). For the same scene, octrees can be more memory efficient than a uniform grid, but they are usually irregular in volume grids. However, the computational cost of rendering a scene with an octree is higher than with a uniform grid, and space-time storage is significantly more efficient with a uniform grid. Octrees are also less cache friendly than uniform grids, increasing rendering efficiency and capacity utilization.

This has prompted many efforts to pursue balanced compromise solutions. One approach is to use a hybrid representation to achieve the best of both worlds in terms of both detail and performance (Chen et al. 2003; Chen et al. 2004; Chen et al. 2005; Chen et al. 2006; Chen et al. 2007; Chen et al. 2008; Chen et al. 2009; Chen et al. 2010; Chen et al. 2011; Chen et al. 2012). Another approach is to use a hierarchical tree structure to store data in a more cache friendly manner (Fidler et al. 2005; Fidler et al. 2006; Fidler et al. 2007; Fidler et al. 2008; Fidler et al. 2009; Fidler et al. 2010; Fidler et al. 2011; Fidler et al. 2012; Fidler et al. 2013; Fidler et al. 2014; Fidler et al. 2015; Fidler et al. 2016; Fidler et al. 2017; Fidler et al. 2018; Fidler et al. 2019; Fidler et al. 2020; Fidler et al. 2021). These methods have been used for various applications, such as ray tracing (Fidler et al. 2005; Fidler et al. 2006; Fidler et al. 2007; Fidler et al. 2008; Fidler et al. 2009; Fidler et al. 2010; Fidler et al. 2011; Fidler et al. 2012; Fidler et al. 2013; Fidler et al. 2014; Fidler et al. 2015; Fidler et al. 2016; Fidler et al. 2017; Fidler et al. 2018; Fidler et al. 2019; Fidler et al. 2020; Fidler et al. 2021), volume rendering (Fidler et al. 2005; Fidler et al. 2006; Fidler et al. 2007; Fidler et al. 2008; Fidler et al. 2009; Fidler et al. 2010; Fidler et al. 2011; Fidler et al. 2012; Fidler et al. 2013; Fidler et al. 2014; Fidler et al. 2015; Fidler et al. 2016; Fidler et al. 2017; Fidler et al. 2018; Fidler et al. 2019; Fidler et al. 2020; Fidler et al. 2021), and adaptive simulation on areas of interest (Fidler et al. 2005; Fidler et al. 2006; Fidler et al. 2007; Fidler et al. 2008; Fidler et al. 2009; Fidler et al. 2010; Fidler et al. 2011; Fidler et al. 2012; Fidler et al. 2013; Fidler et al. 2014; Fidler et al. 2015; Fidler et al. 2016; Fidler et al. 2017; Fidler et al. 2018; Fidler et al. 2019; Fidler et al. 2020; Fidler et al. 2021).

© 2014 IEEE. Reprinted with permission. All rights reserved.



# Composable Data Structures

```
root.hash(ijk, 32).dense(ijk, 16).pointer()  
    .dense(ijk, 8).place(u, v, w);
```

(a) 3D VDB-style [Museth 2013] structure with configuration [5, 4, 3]. The root-level hash table allows negative coordinates to be accessed, providing the user with an unbounded domain.

```
root.dense(ijk, 512).morton().bitmasked()  
    .dense(ijk, {8, 4, 4}).place(flags, u, v, w);
```

(b) 3D SPGrids [Setaluri et al. 2014] occupying voxels in the bounding box  $[0, 4096] \times [0, 2048] \times [0, 2048]$ . The data structure is relatively shallow (only two levels), so root-to-leaf accesses have relatively low cost.

```

// "Hierarchical Particle Buckets": each leaf block contains all indices of particles within its range
root.dynamic(l, 2048).place(particle_x, particle_y, particle_z, particle_mass);
root.hash(ijk, 512).dense(ijk, 32).pointer().dense(ijk, 8).pointer().dynamic(l, 2048).place(particle_index);

// "SPVDB": Unbounded data structures with bitmasks and Morton coding. (VDB and SPGrid combined.)
root.hash(ijk, 512).dense(ijk, 512).morton().bitmasked().dense(ijk, {8, 4, 4}).place(flags, u, v, w);

```

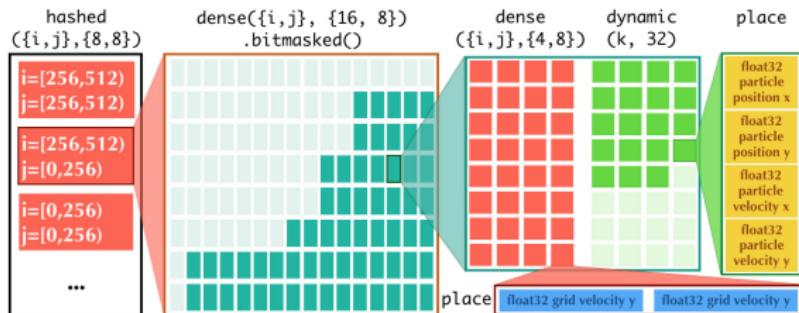
```

// "Hybrid Eulerian-Lagrangian Grid"
auto &block = root
    .hash(ij, 8)
    .dense({ij, {16, 8}})
    .bitmasked();

// Child 1: grid nodes
block.dense(ij, {4, 8})
    .place(grid_vx)
    .place(grid_vy);

// Child 2: particles
block.dynamic(i, 32)
    .place(part_x)
    .place(part_y)
    .place(part_vx)
    .place(part_vy);

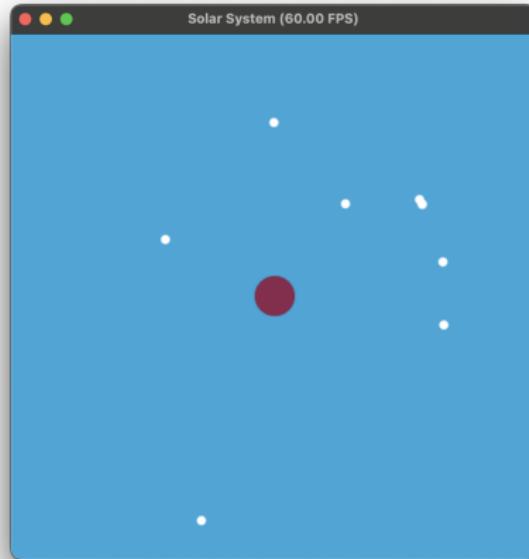
```



(c) "HPB", "SPVDB", "HLEG": We can easily design new data structures with customized features.

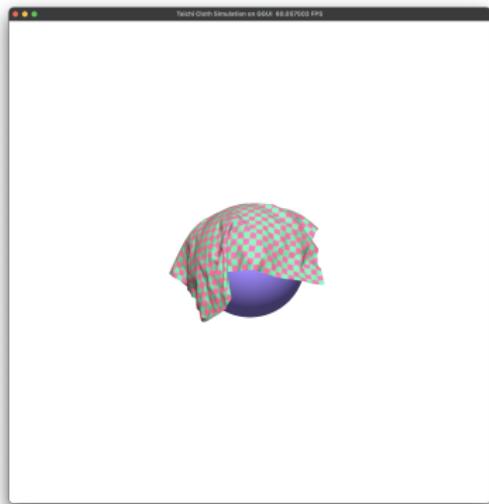
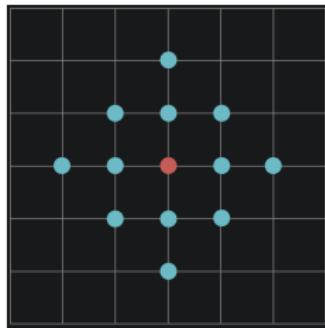
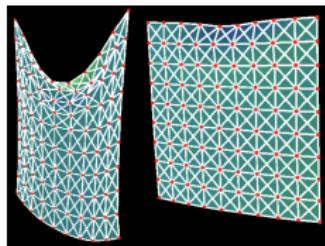
Fig. 6. The layout language allows users to define data structures using our building blocks. We can reproduce two popular multi-level sparse grid used in simulation (a) (b). Furthermore, we can use our language to design new data structures (c) by chaining and forking elementary components. Hybrid Eulerian-Lagrangian simulations (e.g. FLIP [Zhu and Bridson 2005] and MPM [Stomakhin et al. 2013]) often need to maintain both particles and grids, and the required data structures are usually complicated. Using these building blocks, we easily found a data structure with a hierarchical pointer list of particles, which we call *Hierarchical Particle Buckets*, that is especially useful for the material point method simulation (Sec.6.1).

# Objective Data-oriented Programming Example



Stony Brook University

# Physical Simulation Example



Stony Brook University

# Conclusion

- ▶ Today we've introduced Taichi
  - ▶ Python based implementation
  - ▶ Design Goals
  - ▶ Applications and Research Projects
  - ▶ Sparse Data Structures
- ▶ Questions?



Stony Brook University

# References I

- [1] Niels Aage et al. "Giga-voxel computational morphogenesis for structural design". In: *Nature* 550.7674 (2017), pp. 84–86.
- [2] Yuanming Hu et al. "ChainQueen: A Real-Time Differentiable Physical Simulator for Soft Robotics". In: *Proceedings of IEEE International Conference on Robotics and Automation (ICRA)* (2019).
- [3] Yuanming Hu et al. "QuanTaichi: A Compiler for Quantized Simulations". In: *ACM Transactions on Graphics (TOG)* 40.4 (2021).
- [4] Yuanming Hu et al. "Taichi: a language for high-performance computation on spatially sparse data structures". In: *ACM Transactions on Graphics (TOG)* 38.6 (2019), p. 201.
- [5] Haixiang Liu et al. "Narrow-band topology optimization on a sparsely populated grid". In: *ACM Transactions on Graphics (TOG)* 37.6 (2018), pp. 1–14.

## References II

- [6] Ken Museth. "VDB: High-resolution sparse volumes with dynamic topology". In: *ACM transactions on graphics (TOG)* 32.3 (2013), pp. 1–22.
- [7] Rajsekhar Setaluri et al. "SPGrid: A sparse paged grid structure applied to adaptive smoke simulation". In: *ACM Transactions on Graphics (TOG)* 33.6 (2014), pp. 1–12.
- [8] Andrew Spielberg et al. "Learning-in-the-loop optimization: End-to-end control and co-design of soft robots through learned deep latent representations". In: *Advances in Neural Information Processing Systems* 32 (2019), pp. 8284–8294.

# Source Files

The source files live on github:

[https://github.com/SallySoul/stencil\\_latex/tree/main/slides/taichi](https://github.com/SallySoul/stencil_latex/tree/main/slides/taichi)

≈ 570 lines of latex

35 images



Stony Brook University