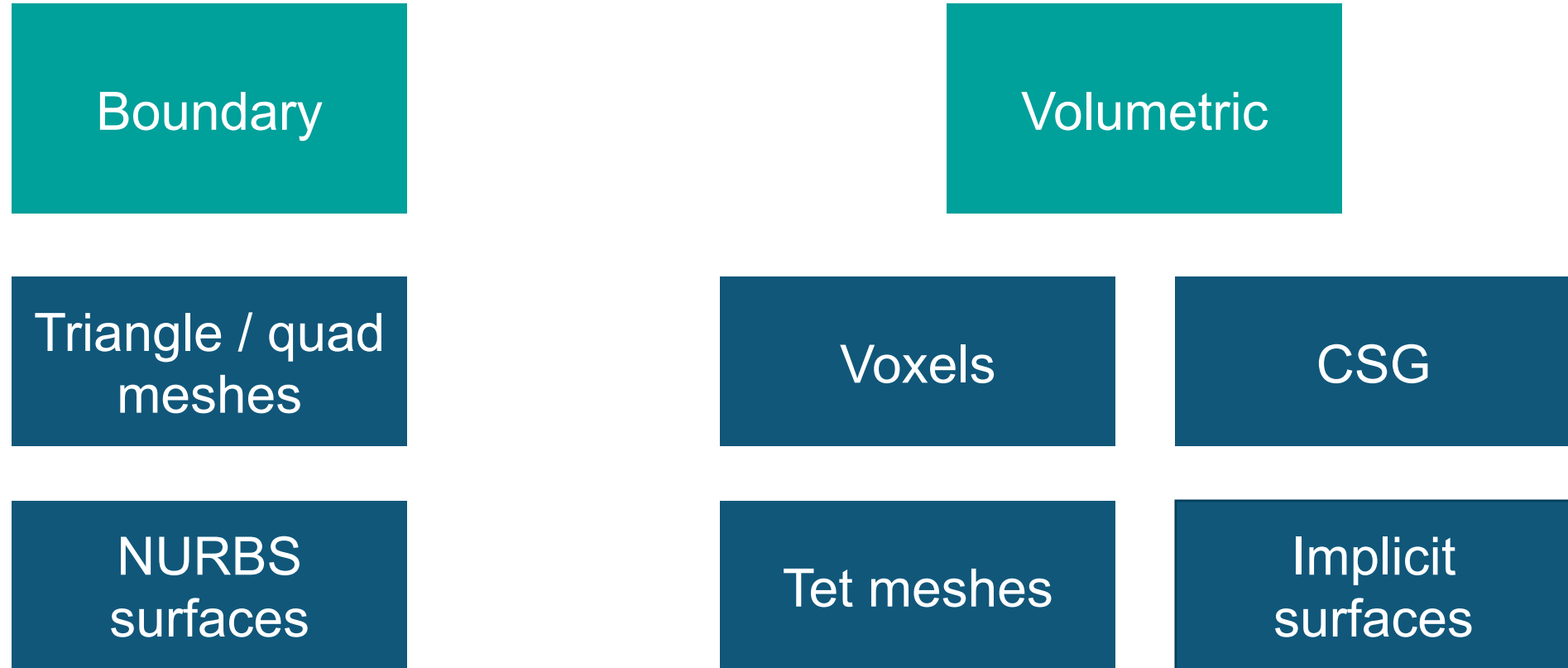


# **CLOSED-FORM IMPLICIT SURFACES ON THE GPU**

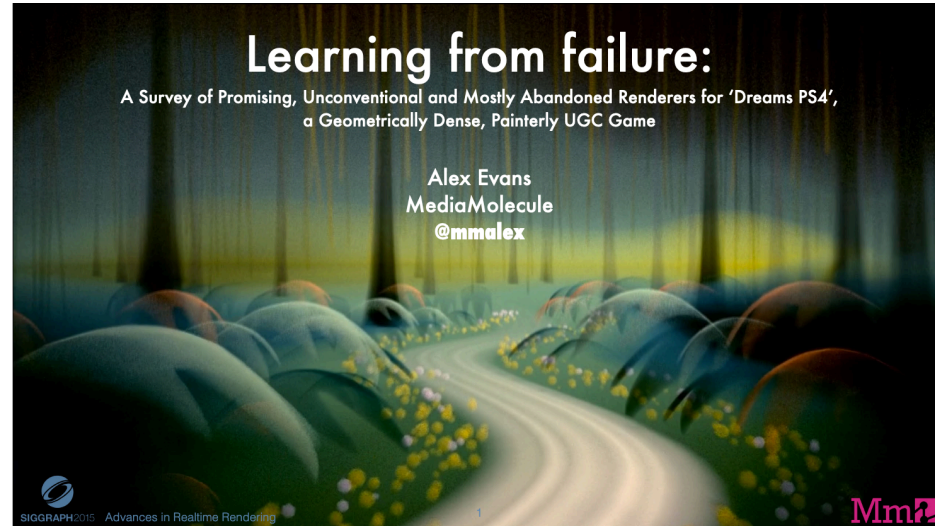
Matt Keeter

Independent researcher

# SHAPE REPRESENTATIONS



# VOLUMETRIC RENAISSANCE



## NeRF

Representing Scenes as Neural Radiance Fields for View Synthesis  
ECCV 2020 Oral - Best Paper Honorable Mention

Ben Mildenhall\*  
UC Berkeley

Pratul P. Srinivasan\*  
UC Berkeley

Matthew Tancik\*  
UC Berkeley

Jonathan T. Barron  
Google Research

Ravi Ramamoorthi  
UC San Diego

Ren Ng  
UC Berkeley

\* Denotes Equal Contribution

Paper

[Data](#)

[Code](#)

0:00 / 0:04

## MagicaCSG @ ephtracy

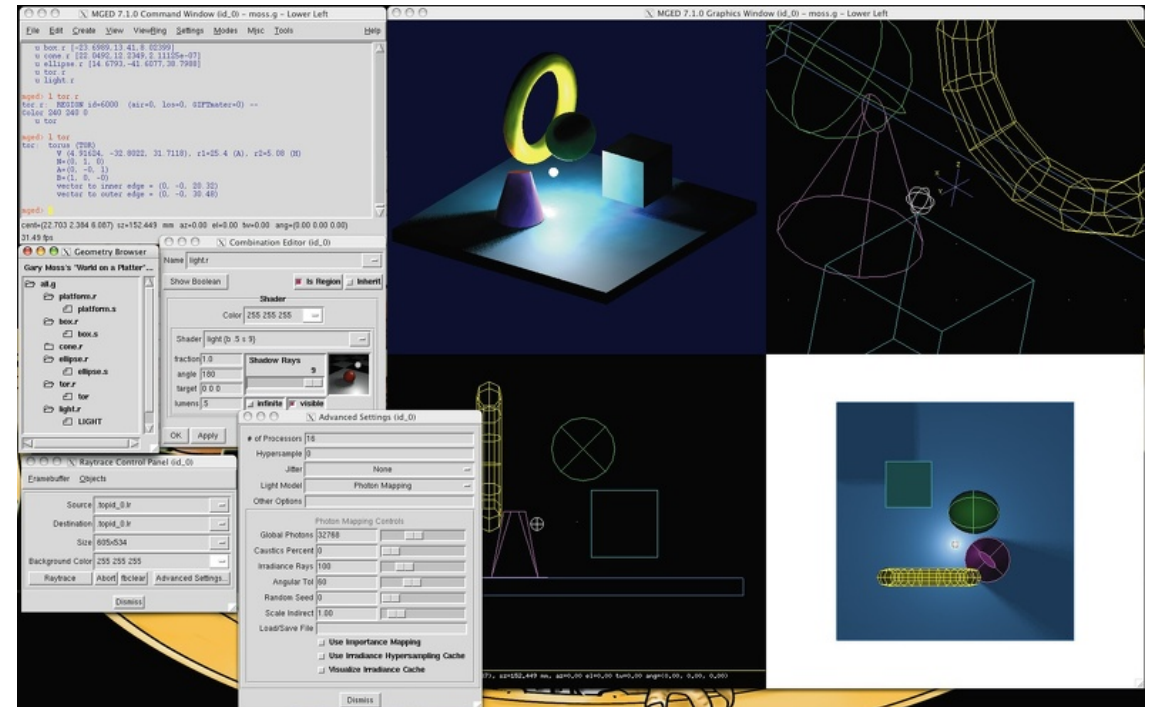
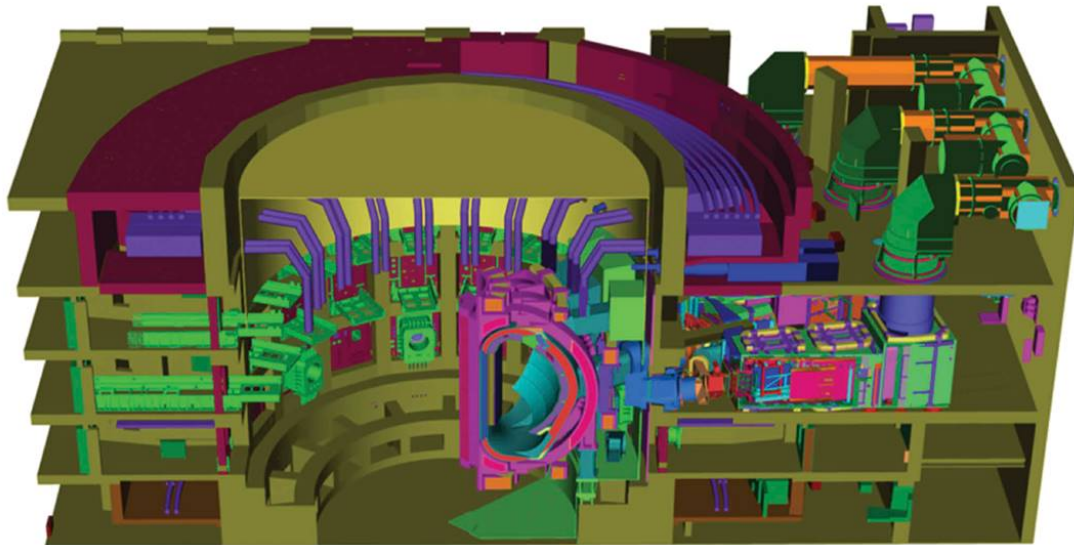
A lightweight signed distance field editor and path tracing renderer.

[win64 0.0.2 Download](#) [3rd Party Intro Video](#)

0:00 / 0:04

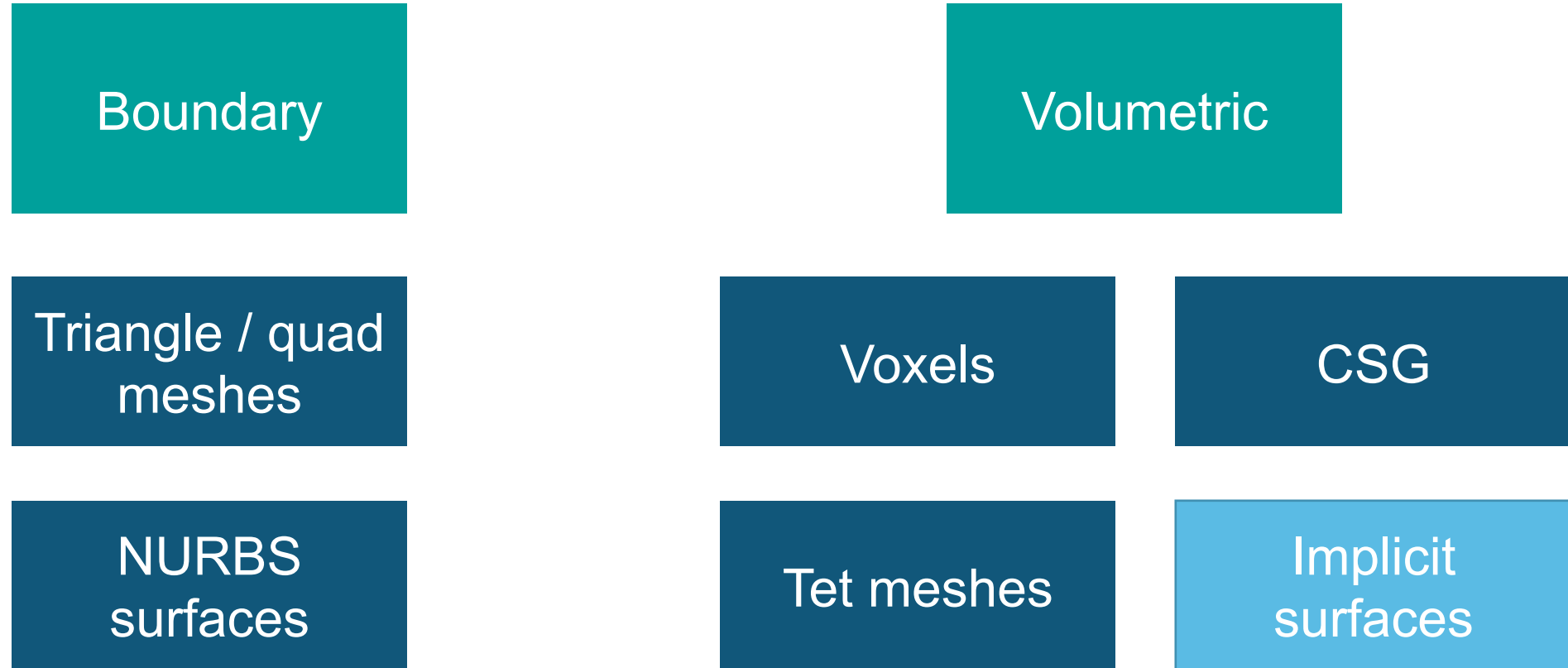
# HISTORICAL VOLUMETRIC MODELING

## Non-interactive Ray-traced CSG via Radiant

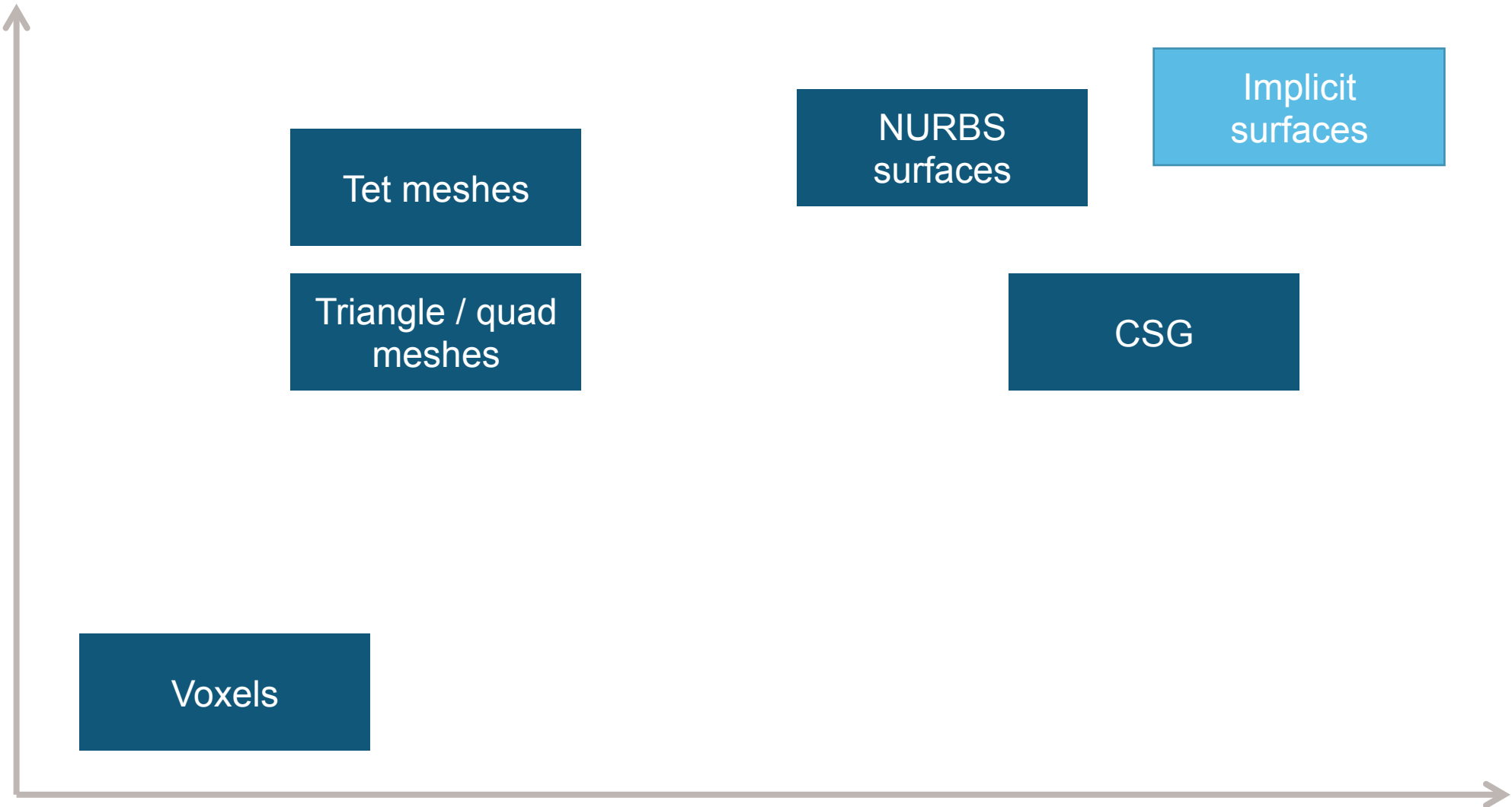




# SHAPE REPRESENTATIONS



Representative power



Encoding density

Ease of rendering

NURBS  
surfaces

Tet meshes

Triangle / quad  
meshes

Voxels

CSG

Implicit  
surfaces

Ease of writing kernel

# Massively Parallel Rendering of Complex Closed-Form Implicit Surfaces

MATTHEW J. KEETER, Independent researcher

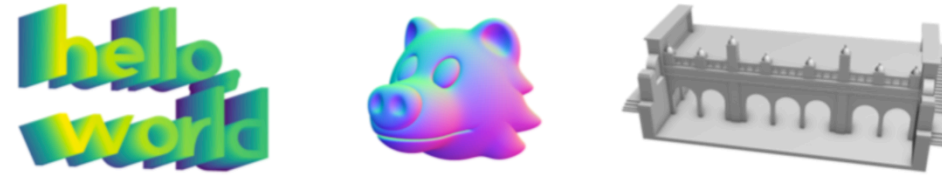


Fig. 1. An assortment of implicit surfaces rendered using our technique. Left: an extruded text string, rotated and rendered as a heightmap. Center: a bear head sculpted using smooth blending operations, with normals found by automatic differentiation. Right: a complex architectural model rendered with screen-space ambient occlusion and perspective. All models are rendered directly from their mathematical representations, without triangulation or raytracing.

We present a new method for directly rendering complex closed-form implicit surfaces on modern GPUs, taking advantage of their massive parallelism. Our model representation is unambiguously solid, can be sampled at arbitrary resolution, and supports both constructive solid geometry (CSG) and more unusual modeling operations (e.g. smooth blending of shapes). The rendering strategy scales to large-scale models with thousands of arithmetic operations in their underlying mathematical expressions. Our method only requires  $C^0$  continuity, allowing for warping and blending operations which break Lipschitz continuity.

To render a model, its underlying expression is evaluated in a shallow hierarchy of spatial regions, using a high branching factor for efficient parallelization. Interval arithmetic is used to both skip empty regions and construct reduced versions of the expression. The latter is the optimization that makes our algorithm practical: in one benchmark, expression complexity decreases by two orders of magnitude between the original and reduced expressions. Similar algorithms exist in the literature, but tend to be deeply recursive with heterogeneous workloads in each branch, which makes them GPU-unfriendly; our evaluation and expression reduction both run efficiently as massively parallel algorithms, entirely on the GPU.

The resulting system renders complex implicit surfaces in high resolution and at interactive speeds. We examine how performance scales with computing power, presenting performance results on hardware ranging from older laptops to modern data-center GPUs, and showing significant improvements at each stage.

CCS Concepts: • **Computing methodologies** → **Rasterization; Volumetric models.**

Additional Key Words and Phrases: implicit surface, signed distance field, freps, octrees, rasterization, gpu, cuda

Author's address: Matthew J. Keeter, Independent researcher, matt.j.keeter@gmail.com.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).  
© 2020 Copyright held by the owner/author(s). Publication rights licensed to ACM.  
0730-0301/2020/7-ART141 \$15.00  
<https://doi.org/10.1145/3386569.3392429>

## ACM Reference Format:

Matthew J. Keeter. 2020. Massively Parallel Rendering of Complex Closed-Form Implicit Surfaces. *ACM Trans. Graph.* 39, 4, Article 141 (July 2020), 10 pages. <https://doi.org/10.1145/3386569.3392429>

## 1 INTRODUCTION

Implicit surfaces and functional representations are a powerful way to represent solid models [Bloomenthal and Wyvill 1997; Gomes et al. 2009]. Compared to boundary representations (e.g. triangle meshes or NURBS surfaces), they offer unambiguous inside-outside checking, easy constructive solid geometry (CSG) operations, and arbitrary resolution. In recent years, functional representations (f-reps) have been used as the kernel of both commercial [Courter 2019] and open-source [Keeter 2019] CAD packages. They are a fundamental building block in the demoscene community [Burger et al. 2002; Quilez 2008], used as a representation for generative art [Moen 2019], and even as the underlying technology for a recent PlayStation 4 game [Evans 2015].

Unlike boundary representations, implicit surfaces cannot easily be rendered in their native forms. This paper presents a new method for rendering the family of implicit surfaces represented by arbitrary closed-form arithmetic expressions, i.e., representing a sphere as

$$f(x, y, z) < 0 \text{ where } f(x, y, z) = \sqrt{x^2 + y^2 + z^2} - 1$$

This representation is particularly flexible and can be treated as an “assembly language for shapes” which is targeted by higher-level tools. The space of higher-level tools spans the gamut from advanced solid modeling packages [Allen 2019] to user-friendly content generation tools [Keeter 2015].

Our rendering strategy runs in both 2D and 3D, making efficient use of modern GPU hardware and APIs. Unlike previous work, it scales to complex expressions, maintaining interactive framerates while rendering models built from hundreds or thousands of arithmetic operations. It requires no continuity higher than  $C^0$ , which allows for extremely flexible modeling and unusual spatial transformations. Finally, it scales well with GPU power; as GPU performance

# DEFINING IMPLICIT SURFACES

$f(x, y, z) < 0 \rightarrow$  inside the shape (“filled”)

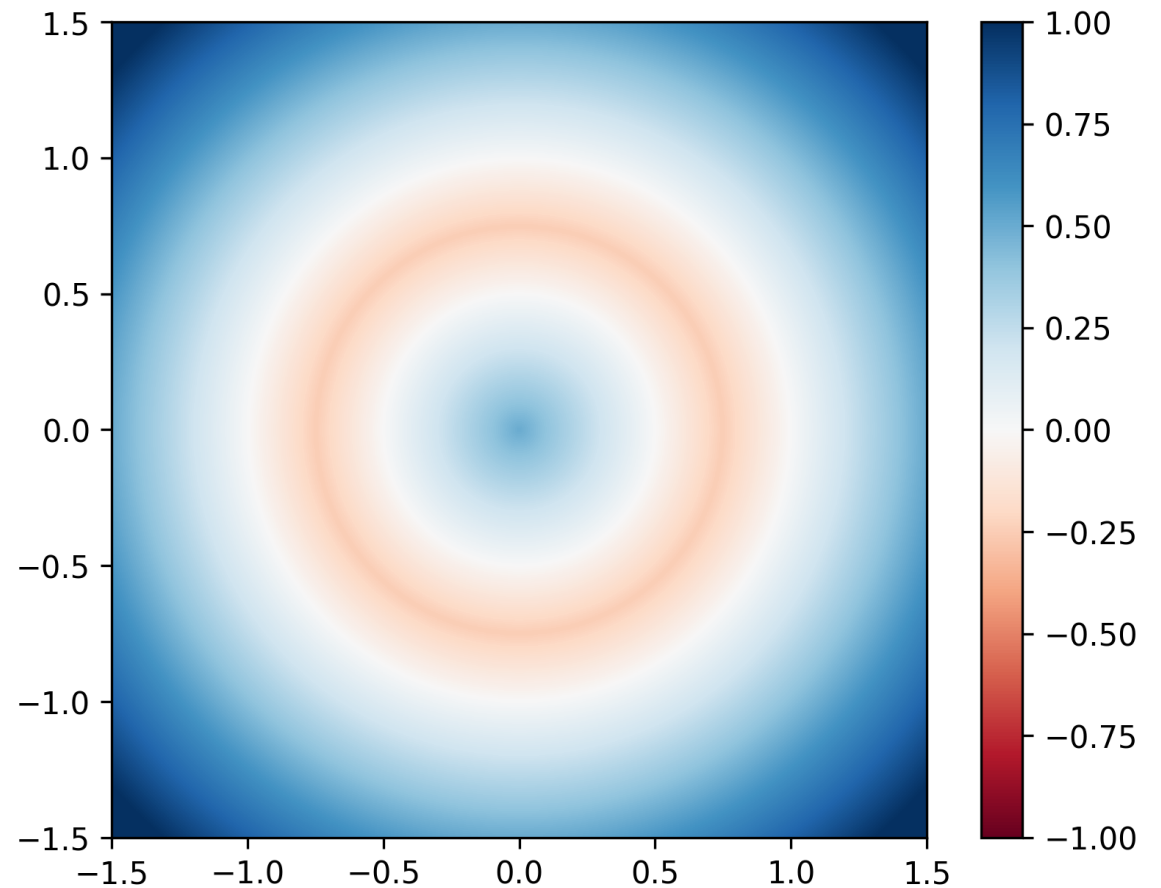
$f(x, y, z) > 0 \rightarrow$  outside the shape (“empty”)

$f(x, y, z) = 0 \rightarrow$  surface of the shape



# CLOSED-FORM IMPLICIT SURFACES

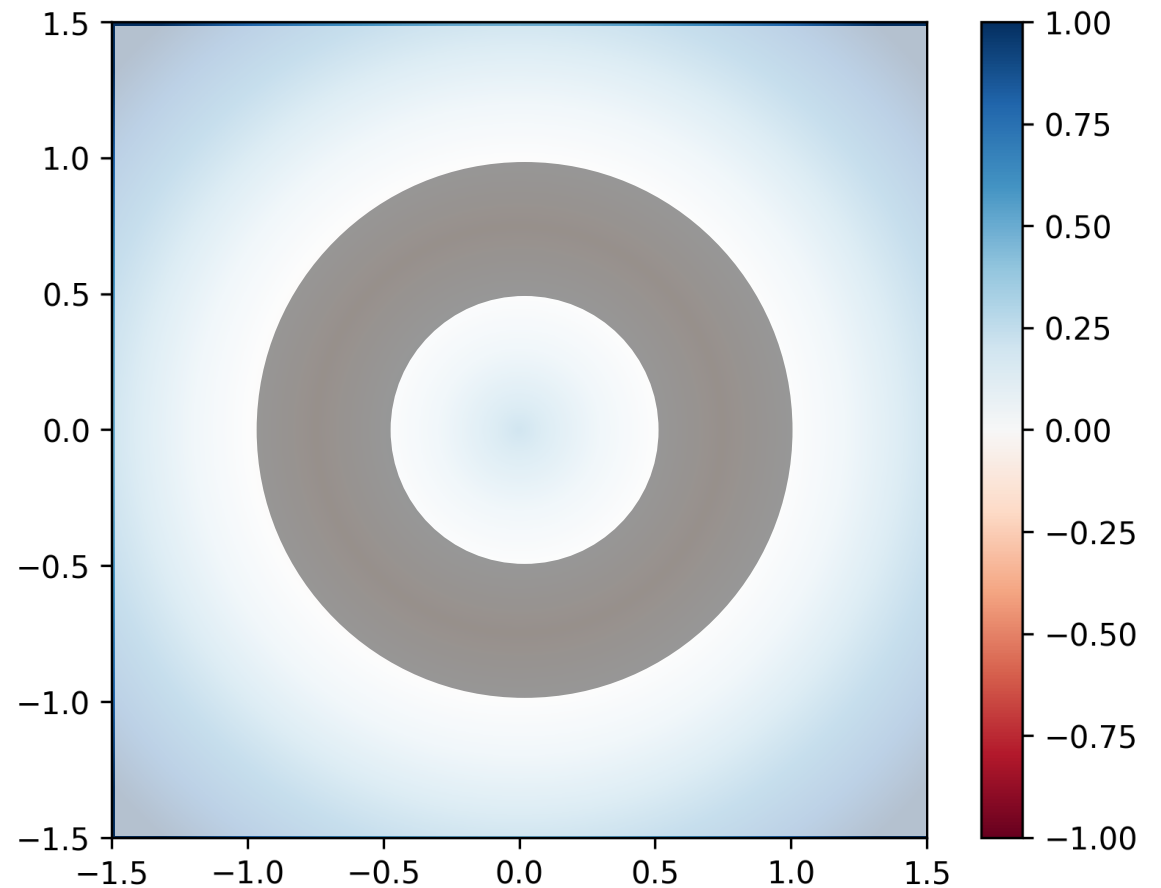
- $f(x, y, z)$  is some sequence of mathematical operations
- Our reference implementation supports  $+$ ,  $-$ ,  $\times$ ,  $/$ ,  $\text{sqrt}$ ,  $\text{sin}$ ,  $\text{cos}$ ,  $\text{asin}$ ,  $\text{acos}$ ,  $\text{atan}$ ,  $\text{exp}$ ,  $\text{log}$ ,  $\text{abs}$ ,  $\text{min}$ ,  $\text{max}$
- CSG is easy:
  - Intersection  $\rightarrow \text{max}$
  - Union  $\rightarrow \text{min}$
  - Inverse  $\rightarrow \text{negation}$



$$\max \left( 0.5 - \sqrt{x^2 + y^2}, \sqrt{x^2 + y^2} - 1 \right) < 0$$

# CLOSED-FORM IMPLICIT SURFACES

- $f(x, y, z)$  is some sequence of mathematical operations
- Our reference implementation supports  $+$ ,  $-$ ,  $\times$ ,  $/$ ,  $\text{sqrt}$ ,  $\text{sin}$ ,  $\text{cos}$ ,  $\text{asin}$ ,  $\text{acos}$ ,  $\text{atan}$ ,  $\text{exp}$ ,  $\text{log}$ ,  $\text{abs}$ ,  $\text{min}$ ,  $\text{max}$
- CSG is easy:
  - Intersection  $\rightarrow \text{max}$
  - Union  $\rightarrow \text{min}$
  - Inverse  $\rightarrow \text{negation}$

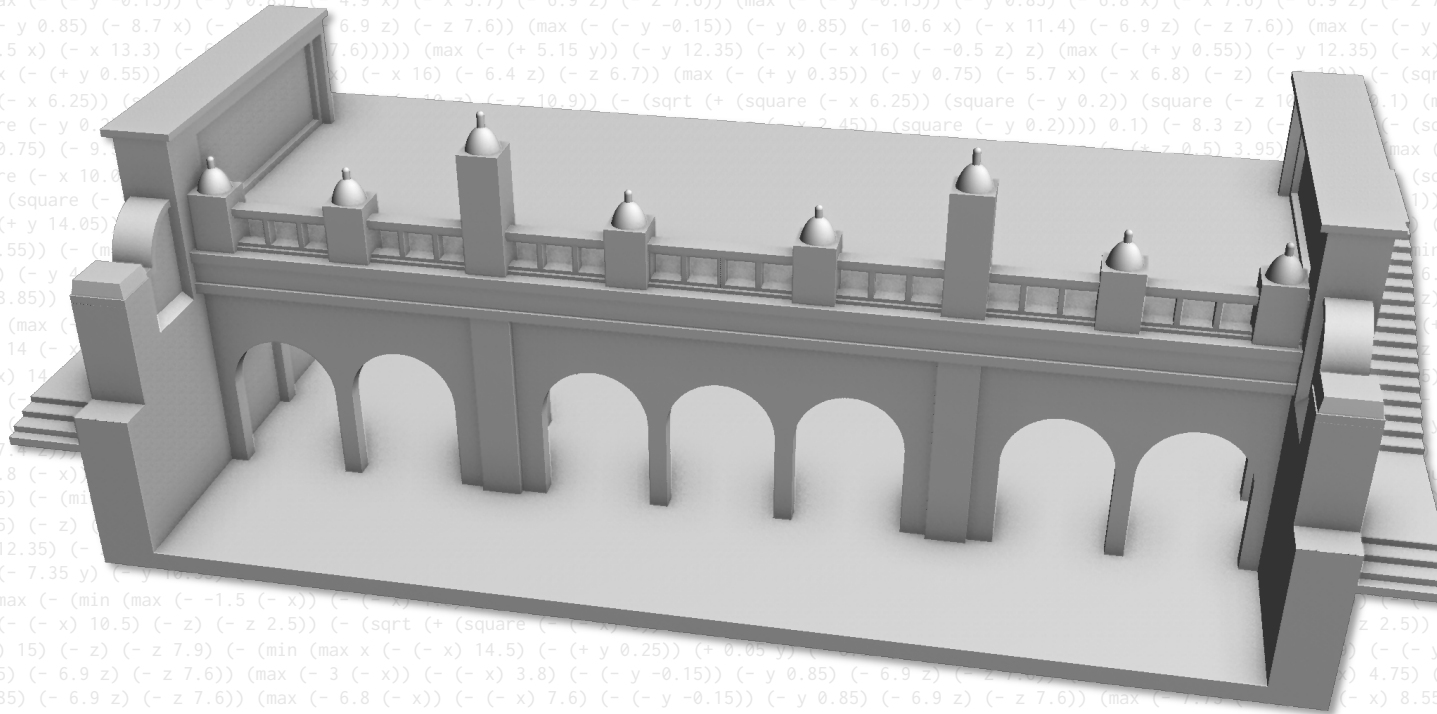


$$\max \left( 0.5 - \sqrt{x^2 + y^2}, \sqrt{x^2 + y^2} - 1 \right) < 0$$

# **COMPLEX CLOSED-FORM IMPLICIT SURFACES**

# COMPLEX CLOSED-FORM IMPLICIT SURFACES

# COMPLEX CLOSED-FORM IMPLICIT SURFACES





**“ASSEMBLY LANGUAGE FOR SHAPES”**

## Benefits

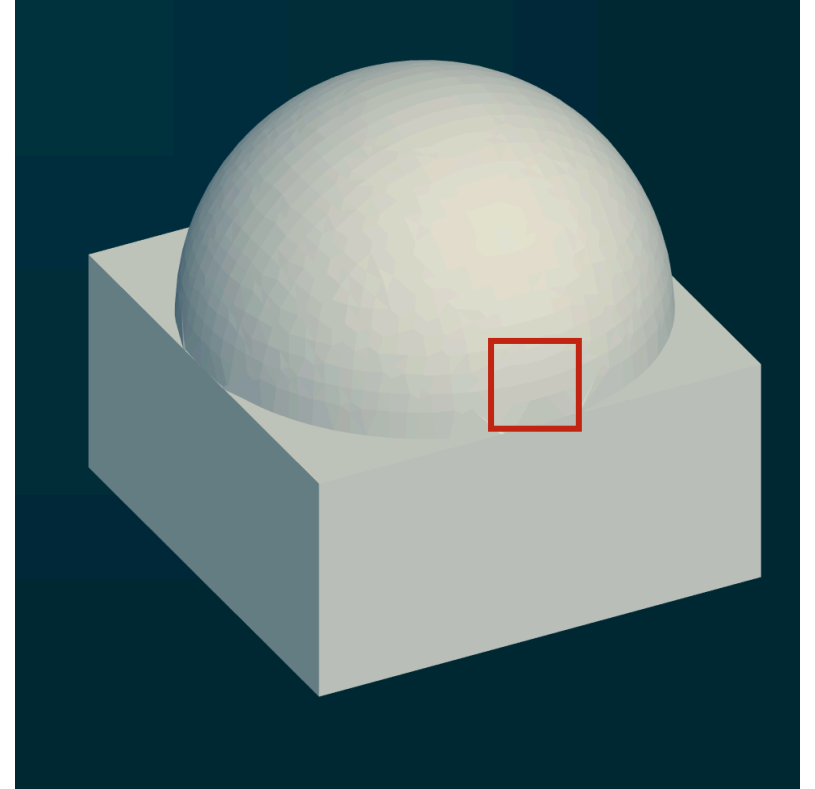
- Easy CSG and solid modeling!
- Unambiguous inside-outside checking
- Supports arbitrary resolution
- Compact representation

## Applications

- CAD/CAM
- Demoscene graphics
- Art and exploration
- User-generated content

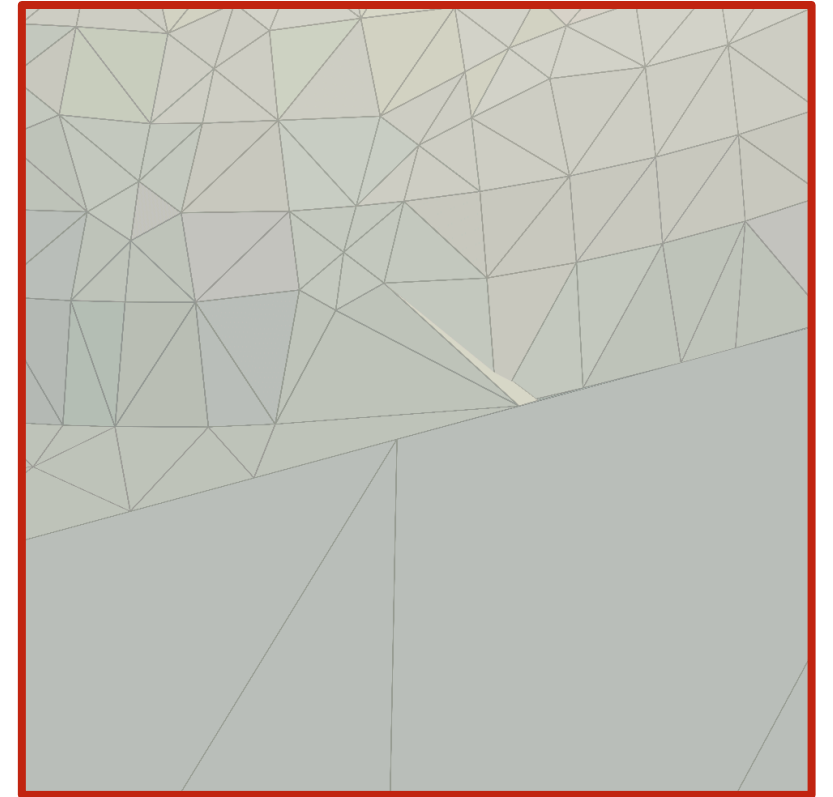
# RENDERING IMPLICIT SURFACES IS HARD

- Meshing
  - Hard to get right for arbitrary shapes
  - Philosophically unsatisfying



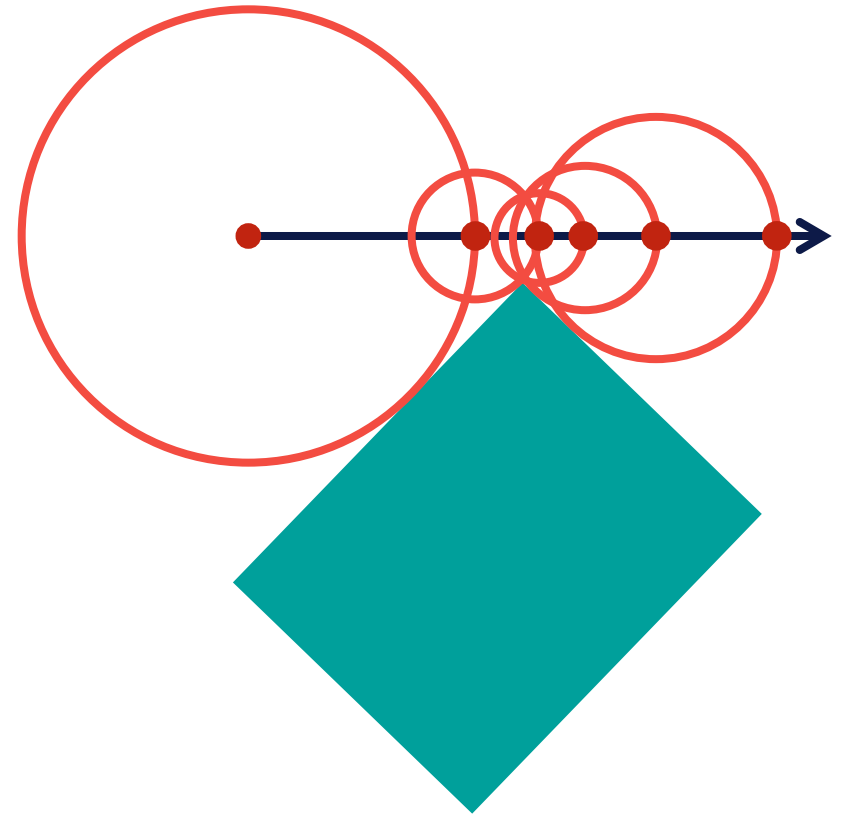
# RENDERING IMPLICIT SURFACES IS HARD

- Meshing
  - Hard to get right for arbitrary shapes
  - Philosophically unsatisfying



# RENDERING IMPLICIT SURFACES IS HARD

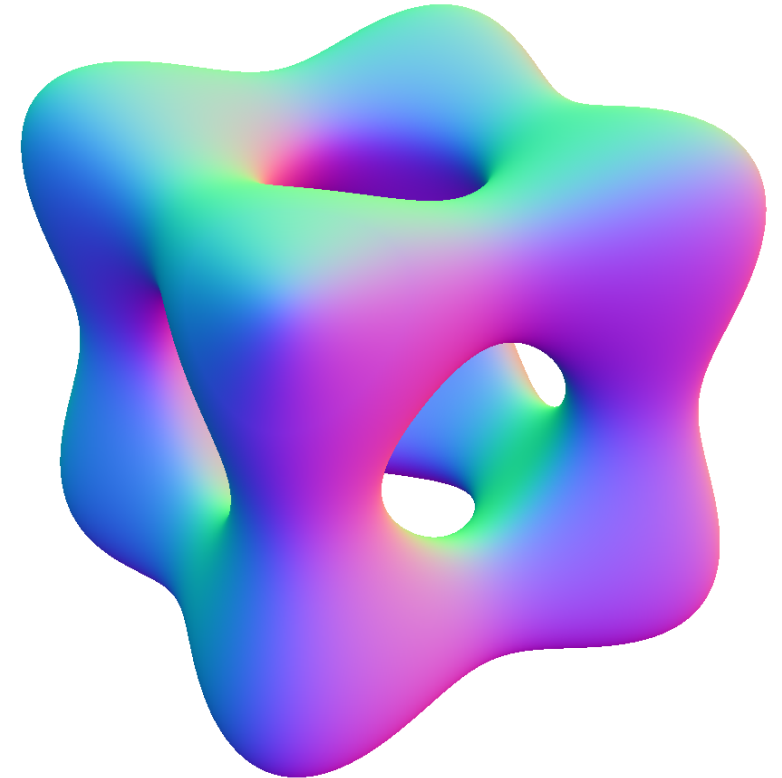
- Meshing
  - Hard to get right for arbitrary shapes
  - Philosophically unsatisfying
- Raytracing
  - Often requires  $C^1$  or Lipschitz continuity
  - Limits the kind of transforms that you can apply



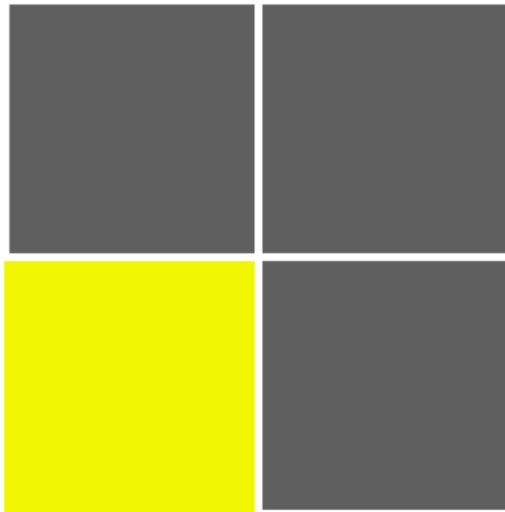


# RENDERING IMPLICIT SURFACES IS HARD

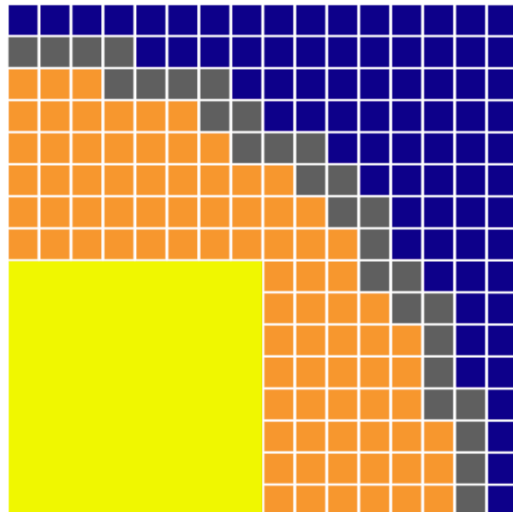
- Meshing
  - Hard to get right for arbitrary shapes
  - Philosophically unsatisfying
- Raytracing
  - Often requires  $C^1$  or Lipschitz continuity
  - Limits the kind of transforms that you can apply
- Our strategy
  - Interval arithmetic + subdivision + recursion (Duff '92)
  - Per-tile command lists on the GPU (Ganacim et al. '14)
  - Only requires  $C^0$  continuity, and scales to large models
  - Extremely philosophically satisfying



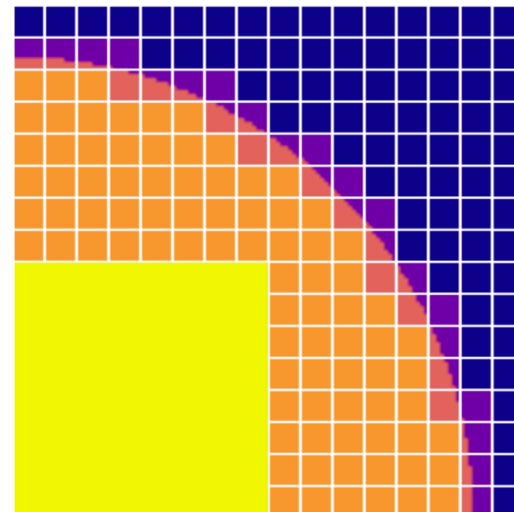
First pass



Second pass



Per-pixel evaluation



64x64 filled tile



8x8 filled tile



8x8 empty tile



8x8 ambiguous tile

Interpreter  
on the GPU

Interval  
arithmetic

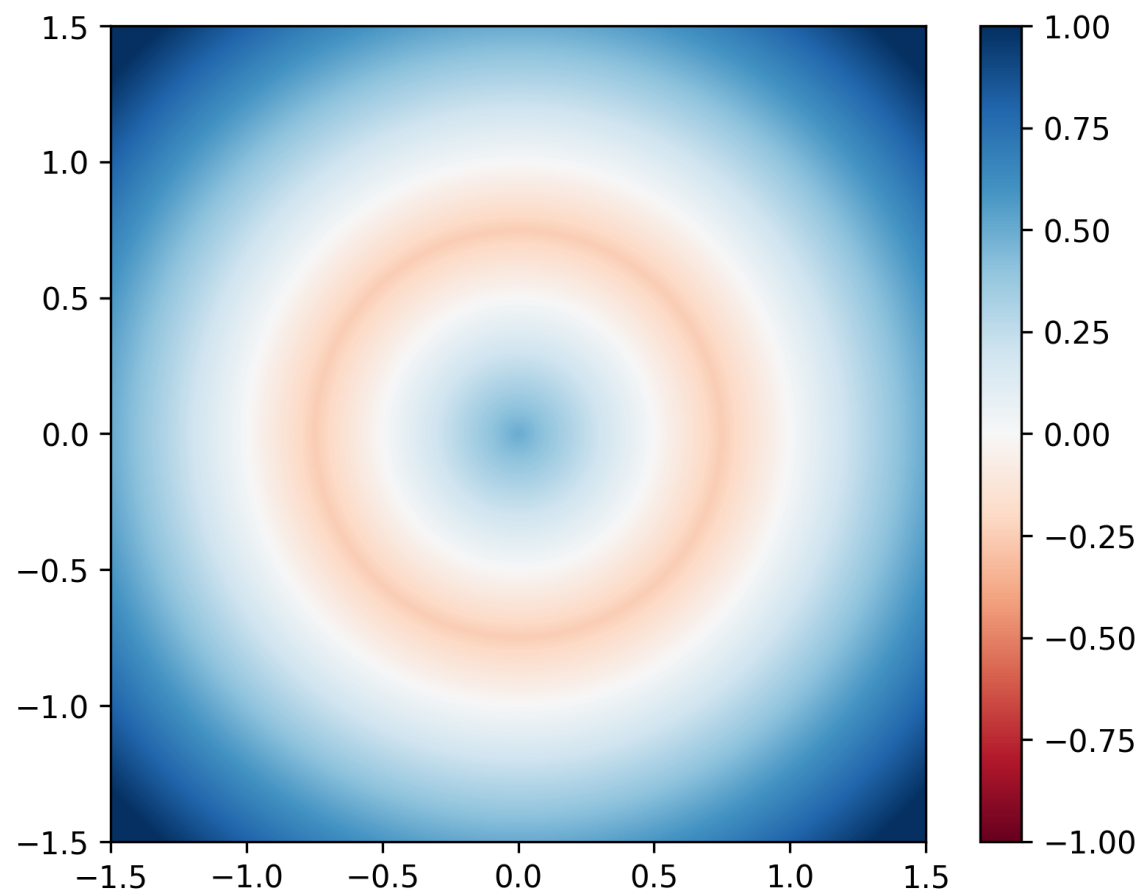
Tape  
shortening

Interpreter  
on the GPU

Interval  
arithmetic

Tape  
shortening

$$\max \left( 0.5 - \sqrt{x^2 + y^2}, \sqrt{x^2 + y^2} - 1 \right) < 0$$



## Compiled kernel

- Uses hardware registers
- Very fast!
- Longer initial startup (must compile one shader per shape)
- Harder to specialize

## Interpreted shape

- Uses soft registers (“slots”)
- Not as fast!
- Faster startup (build tape and send data to GPU)
- Easy to specialize at runtime

## Compiled kernel

- Uses hardware registers
- Very fast!
- Longer initial startup (must compile one shader per shape)
- Harder to specialize

## Interpreted shape

- Uses soft registers (“slots”)
- Not as fast!
- Faster startup (build tape and send data to GPU)
- **Easy to specialize at runtime**

# THE INSTRUCTION TAPE

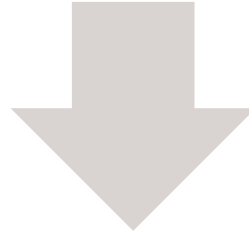
$$\max \left( 0.5 - \sqrt{x^2 + y^2}, \sqrt{x^2 + y^2} - 1 \right)$$

opcode	lhs	rhs	out
X	–	–	slot 0
Y	–	–	slot 1
SQUARE	slot 0	–	slot 0
SQUARE	slot 1	–	slot 1
ADD	slot 0	slot 1	slot 1
SQRT	slot 1	–	slot 1
SUB	slot 1	1.0f	slot 0
SUB	0.5f	slot 1	slot 1
MAX	slot 0	slot 1	slot 1



# TAPE PACKING

opcode	out	lhs	rhs	immediate
u8	u8	u8	u8	f32



8 bytes per clause

# THE INTERPRETER LOOP

```
def run(tape):
    for (opcode, out, lhs, rhs, imm) in tape:
        match opcode:
            # LHS & RHS
            ADD_LHS_RHS: slots[out] = slots[lhs] + slots[rhs]
            SUB_LHS_RHS: slots[out] = slots[lhs] - slots[rhs]
            MUL_LHS_RHS: slots[out] = slots[lhs] * slots[rhs]
            ...
            # LHS & immediate
            ADD_LHS_IMM: slots[out] = slots[lhs] + imm
            ...and so on...

    return slots[tape[-1].out]
```

# INTERPRETER OVERHEAD

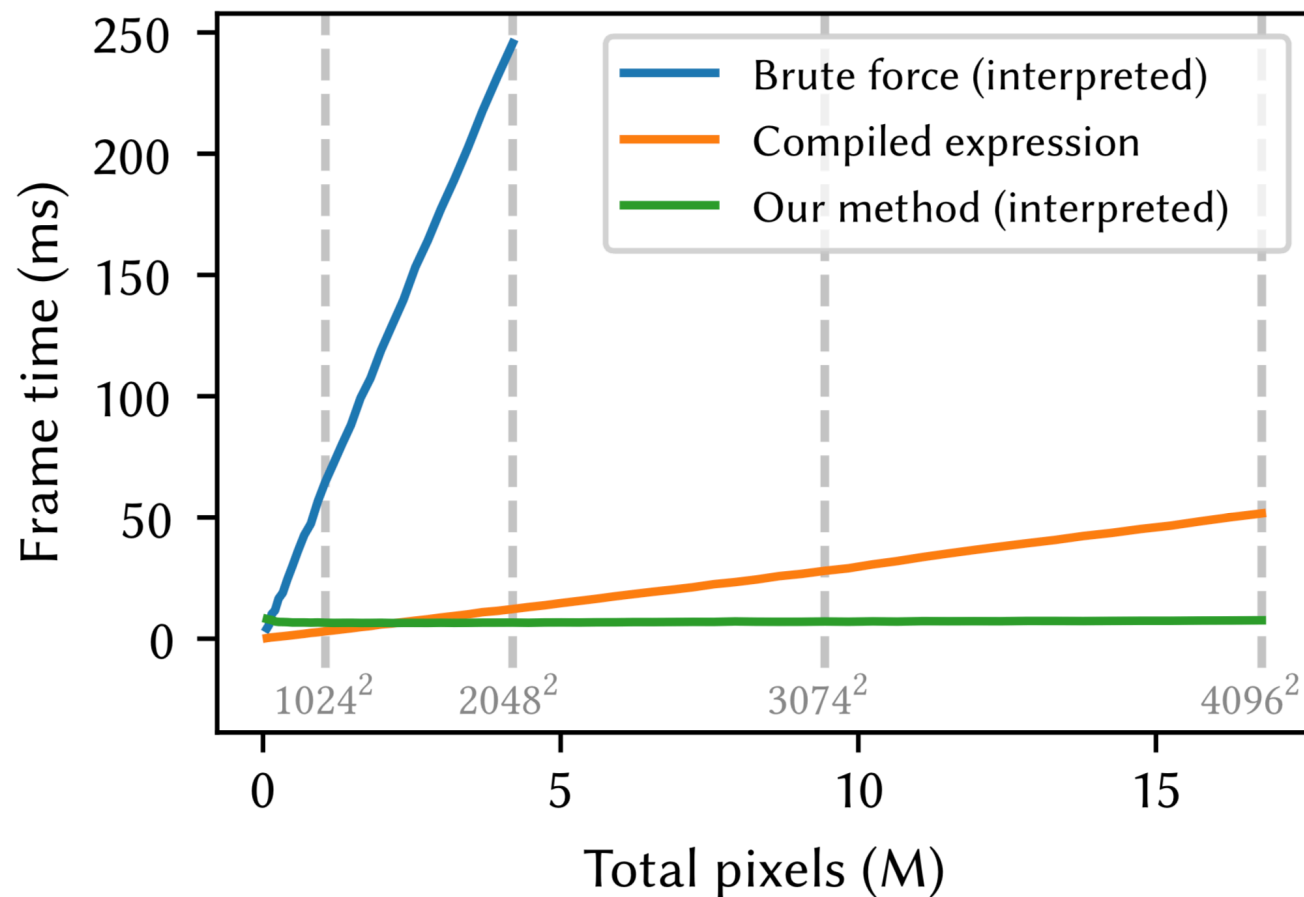
But this rough magic I  
here abjure, and when  
I have required some  
heavenly music, which even  
now I do, to work mine  
end upon their senses that  
this airy charm is for, I'll  
break my staff, bury it  
certain fathoms in the  
earth, and deeper than did  
ever plummet sound  
I'll drown my book.

6056 clauses

# INTERPRETER OVERHEAD

But this rough magic I  
here abjure, and when  
I have required some  
heavenly music, which even  
now I do, to work mine  
end upon their senses that  
this airy charm is for, I'll  
break my staff, bury it  
certain fathoms in the  
earth, and deeper than did  
ever plummet sound  
I'll drown my book.

6056 clauses



# GPU PERFORMANCE NOTES

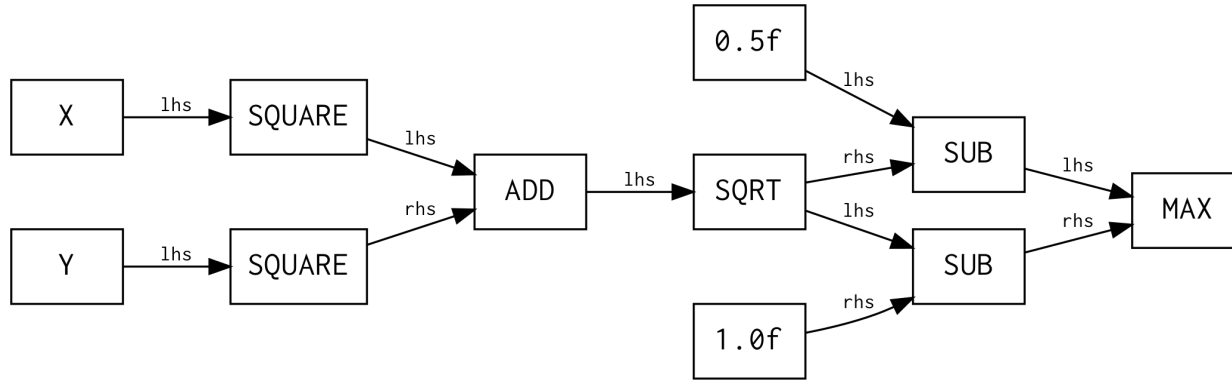
- Implemented in CUDA C/C++, compiled with `nvcc`
- Must use fewer than 32 registers to reach 100% occupancy!
  - Keep inner loop very simple
  - Consider moving setup logic to a separate kernel
  - Prefer C-style over C++-style CUDA
  - Recompile and check after even small changes; the compiler works in mysterious ways
- Avoid thread divergence within a warp
  - Ensure that all threads in a warp are executing the same **tape**
  - 32 threads per warp → 64x subdivision at each level is a convenient multiple

# WHERE DOES THE TAPE COME FROM?

$$\max \left( 0.5 - \sqrt{x^2 + y^2}, \sqrt{x^2 + y^2} - 1 \right)$$

opcode	lhs	rhs	out
X	–	–	slot 0
Y	–	–	slot 1
SQUARE	slot 0	–	slot 0
SQUARE	slot 1	–	slot 1
ADD	slot 0	slot 1	slot 1
SQRT	slot 1	–	slot 1
SUB	slot 1	1.0f	slot 0
SUB	0.5f	slot 1	slot 1
MAX	slot 0	slot 1	slot 1

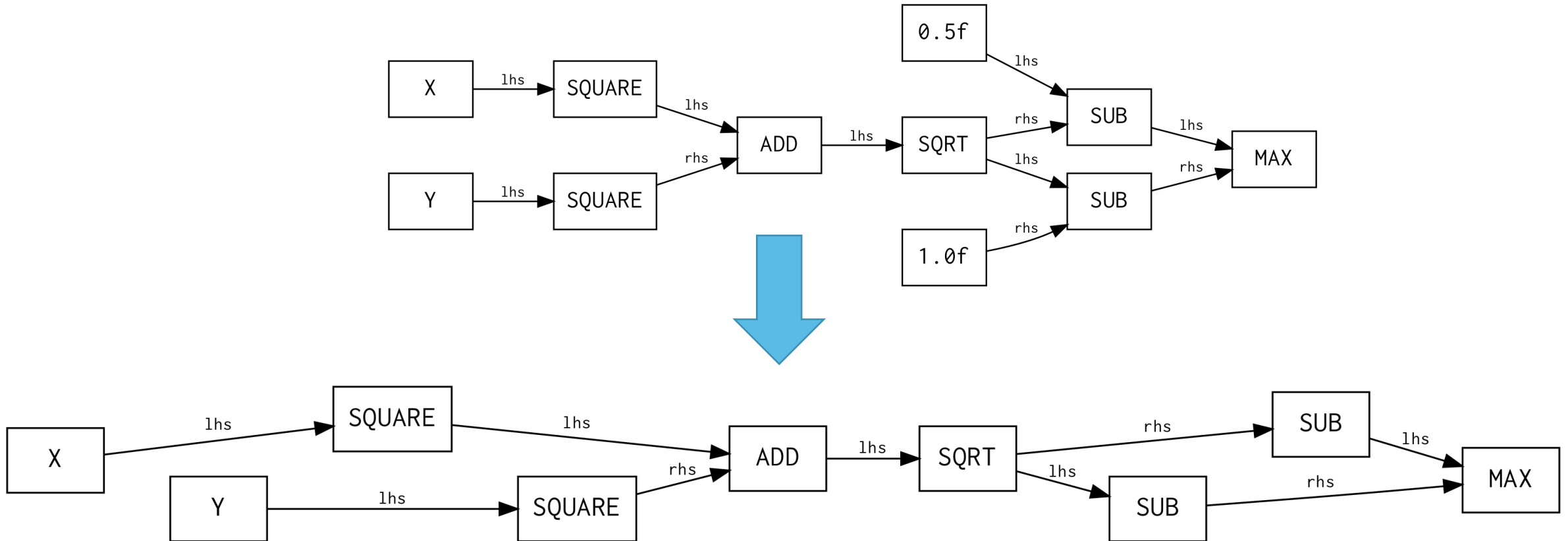
# DAG → INSTRUCTION TAPE



- The equation is encoded as a directed acyclic graph
- Common subexpressions are merged into one node

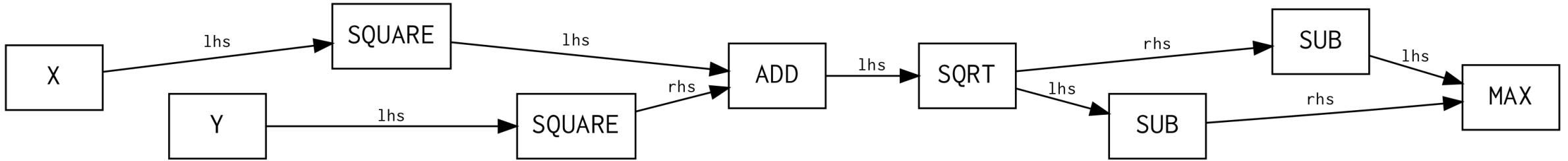
opcode	lhs	rhs	out
X	—	—	slot 0
Y	—	—	slot 1
SQUARE	slot 0	—	slot 0
SQUARE	slot 1	—	slot 1
ADD	slot 0	slot 1	slot 1
SQRT	slot 1	—	slot 1
SUB	slot 1	1.0f	slot 0
SUB	0.5f	slot 1	slot 1
MAX	slot 0	slot 1	slot 1

# TOPOLOGICAL SORT





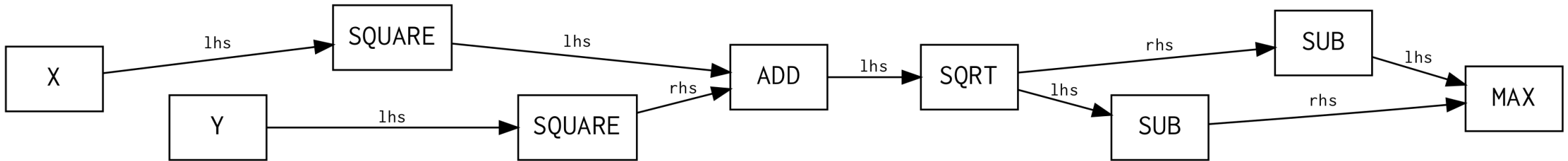
# SLOT ALLOCATION



- Same problem as **register allocation** in compiler design
- Simple greedy algorithm:
  - Maintain a pool of available slots
  - After a slot is used for the last time, release it to the pool
  - Choose the top slot in the pool
  - If pool is empty, create a new slot

opcode	lhs	rhs	out
X	–	–	\$1
Y	–	–	\$2
SQUARE	\$1	–	\$3
SQUARE	\$2	–	\$4
ADD	\$3	\$4	\$5
SQRT	\$5	–	\$6
SUB	\$6	1.0f	\$7
SUB	0.5f	\$6	\$8
MAX	\$7	\$8	\$9

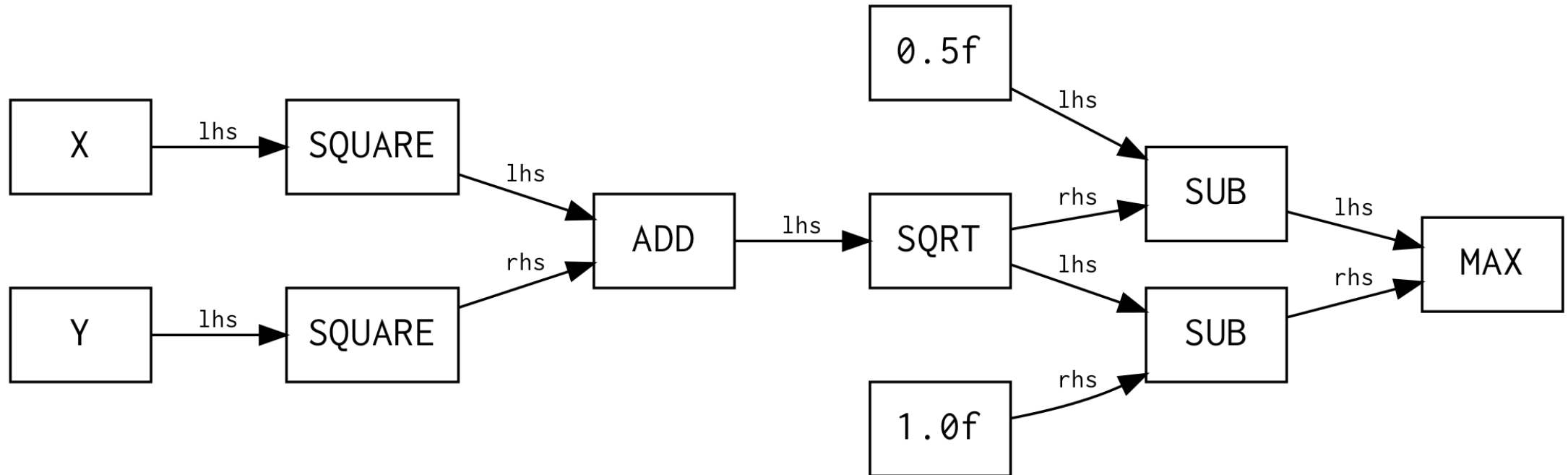
# SLOT ALLOCATION



- Same problem as **register allocation** in compiler design
- Simple greedy algorithm:
  - Maintain a pool of available slots
  - After a slot is used for the last time, release it to the pool
  - Choose the top slot in the pool
  - If pool is empty, create a new slot

opcode	lhs	rhs	out
X	–	–	slot 0
Y	–	–	slot 1
SQUARE	slot 0	–	slot 0
SQUARE	slot 1	–	slot 1
ADD	slot 0	slot 1	slot 1
SQRT	slot 1	–	slot 1
SUB	slot 1	1.0f	slot 0
SUB	0.5f	slot 1	slot 1
MAX	slot 0	slot 1	slot 1

# WHERE DOES THE DAG COME FROM?



## Studio

Desktop design tool

## Bindings

Scheme and Python

## Standard library

Shapes and geometric operations

## C API

```
#include <libfive.h>
```

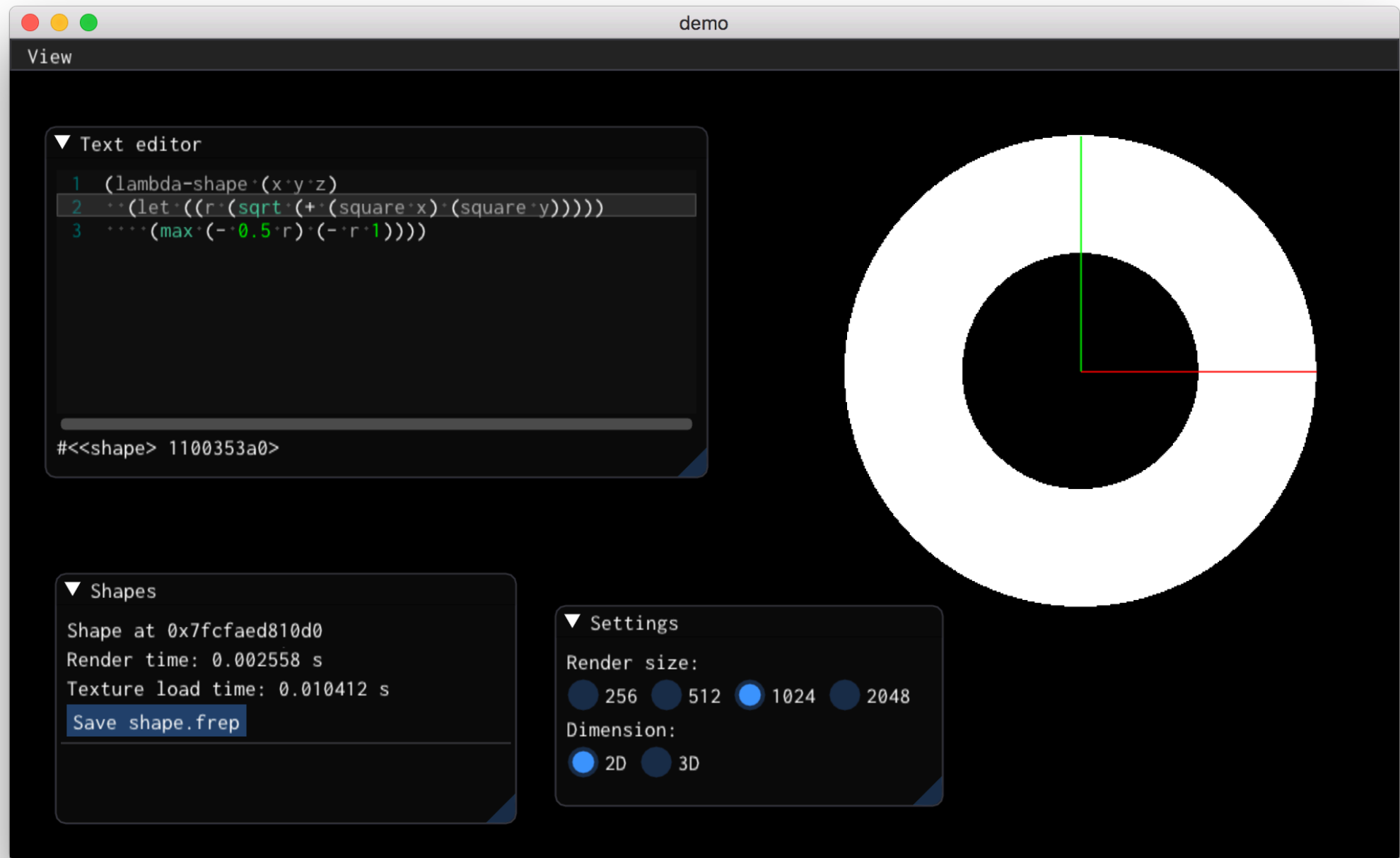
## libfive

C++ library



# libfive

<https://libfive.com>



Interpreter  
on the GPU

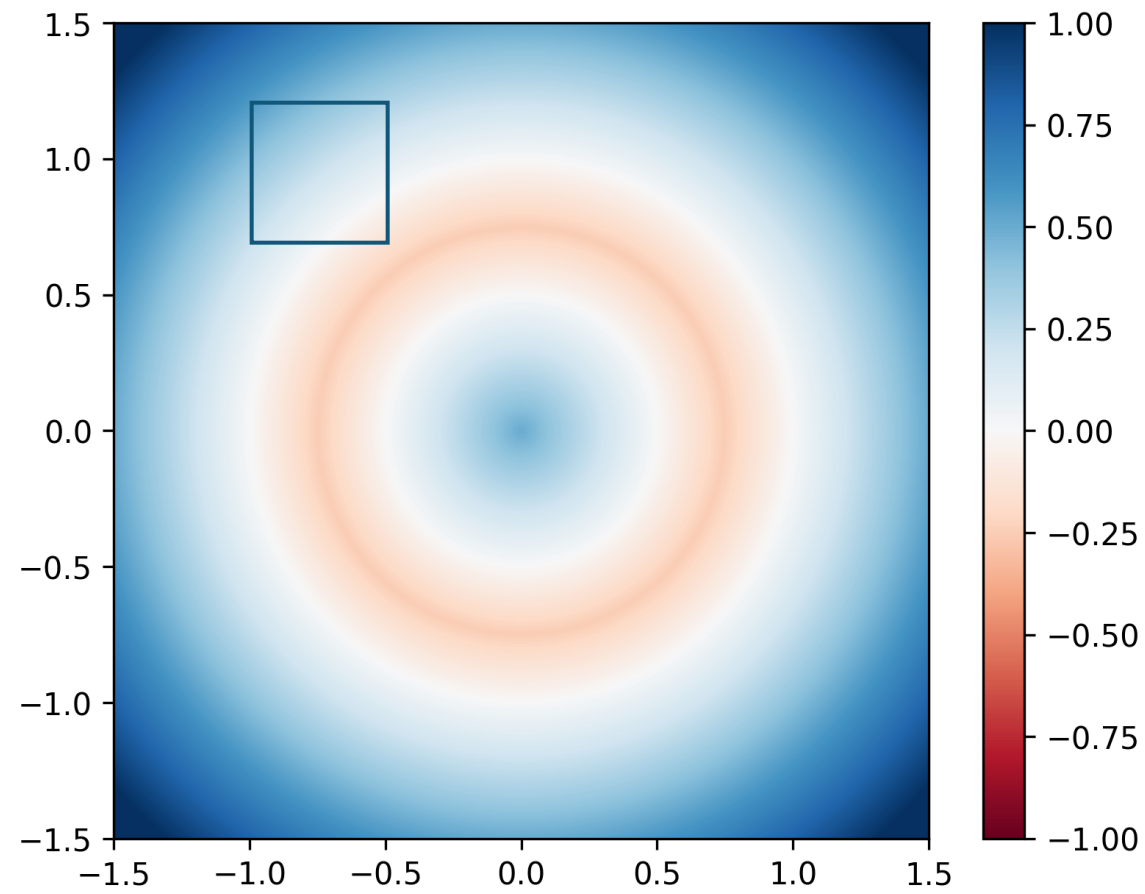
Interval  
arithmetic

Tape  
shortening

Interpreter  
on the GPU

Interval  
arithmetic

Tape  
shortening





# INTERVAL ARITHMETIC

opcode	lhs	rhs	out
X	—	—	$[-1.0, -0.5]$
Y	—	—	$[0.7, 1.2]$
SQUARE		—	
SQUARE		—	
ADD			
SQRT		—	
SUB		$1.0f$	
SUB	$0.5f$		
MAX			

# INTERVAL ARITHMETIC

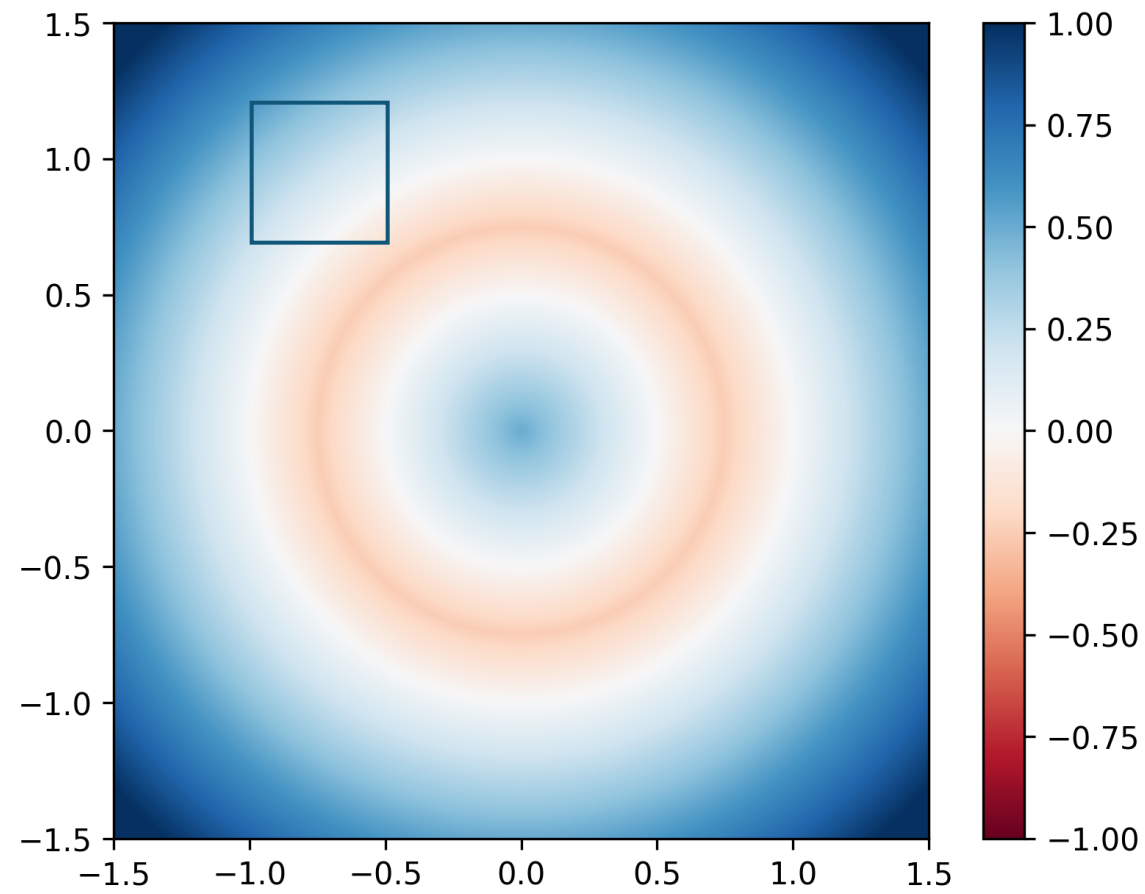
opcode	lhs	rhs	out
X	–	–	$[-1.0, -0.5]$
Y	–	–	$[0.7, 1.2]$
SQUARE	$[-1.0, -0.5]$	–	
SQUARE		–	
ADD			
SQRT		–	
SUB		$1.0f$	
SUB	$0.5f$		
MAX			

# INTERVAL ARITHMETIC

opcode	lhs	rhs	out
X	—	—	$[-1.0, -0.5]$
Y	—	—	$[0.7, 1.2]$
SQUARE	$[-1.0, -0.5]$	—	$[0.25, 1.0]$
SQUARE		—	
ADD			
SQRT		—	
SUB		$1.0f$	
SUB	$0.5f$		
MAX			

# INTERVAL ARITHMETIC

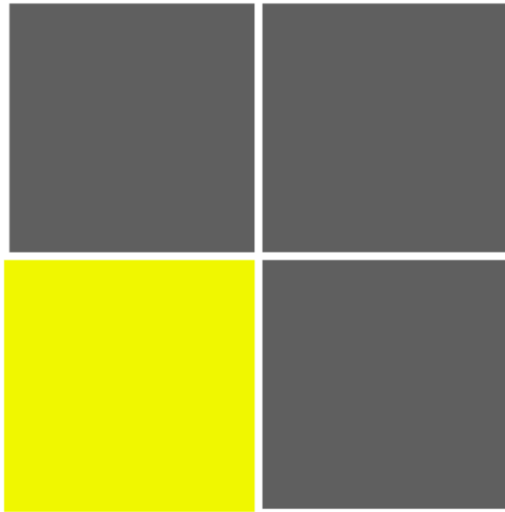
opcode	lhs	rhs	out
X	–	–	$[-1.0, -0.5]$
Y	–	–	$[0.7, 1.2]$
SQUARE	$[-1.0, -0.5]$	–	$[0.25, 1.0]$
SQUARE	$[0.7, 1.2]$	–	$[0.49, 1.44]$
ADD	$[0.25, 1.0]$	$[0.49, 1.44]$	$[0.74, 2.44]$
SQRT	$[0.74, 2.44]$	–	$[0.86, 1.56]$
SUB	$[0.86, 1.56]$	$1.0f$	$[-0.14, 0.56]$
SUB	$0.5f$	$[0.86, 1.56]$	$[-1.06, -0.36]$
MAX	$[-0.14, 0.56]$	$[-1.06, -0.36]$	$[-0.14, 0.56]$



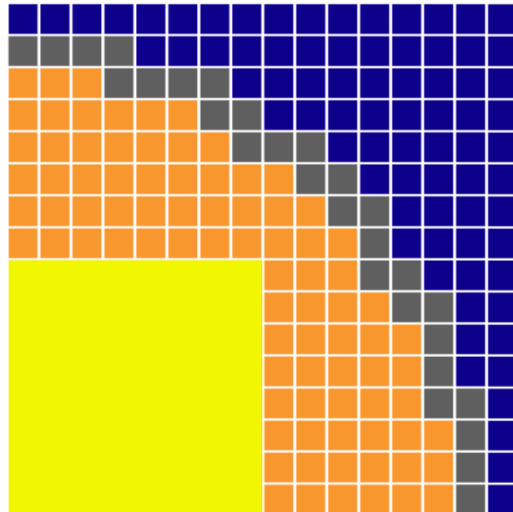
# INTERVAL ARITHMETIC RESULTS

$f(X, Y, Z).upper < 0 \rightarrow$  inside the shape (“filled”)  
 $f(X, Y, Z).lower > 0 \rightarrow$  outside the shape (“empty”)  
otherwise  $\rightarrow$  ambiguous

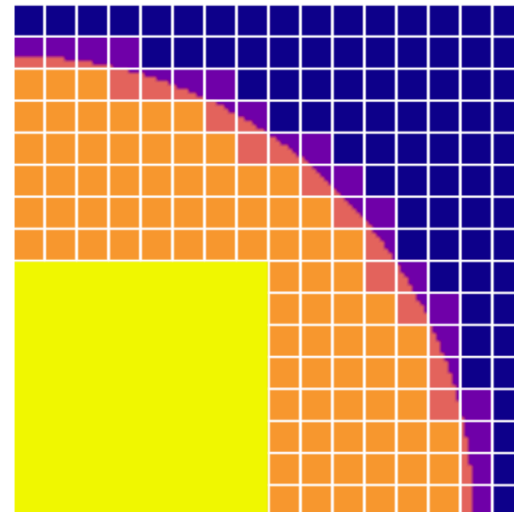
First pass



Second pass



Per-pixel evaluation



64x64 filled tile



8x8 filled tile



8x8 empty tile



8x8 ambiguous tile

# GPU IMPLEMENTATION NOTES

- Implemented a GPU interval class based on `boost::interval`
- Using compiler intrinsics to ensure proper rounding modes

```
__device__ inline Interval operator+(const Interval& x,  
                                     const Interval& y)  
{  
    return {__fadd_rd(x.lower(), y.lower()),  
            __fadd_ru(x.upper(), y.upper())};  
}
```

- Intrinsics are not available on all GPU compute platforms (but it may not matter)



Interpreter  
on the GPU

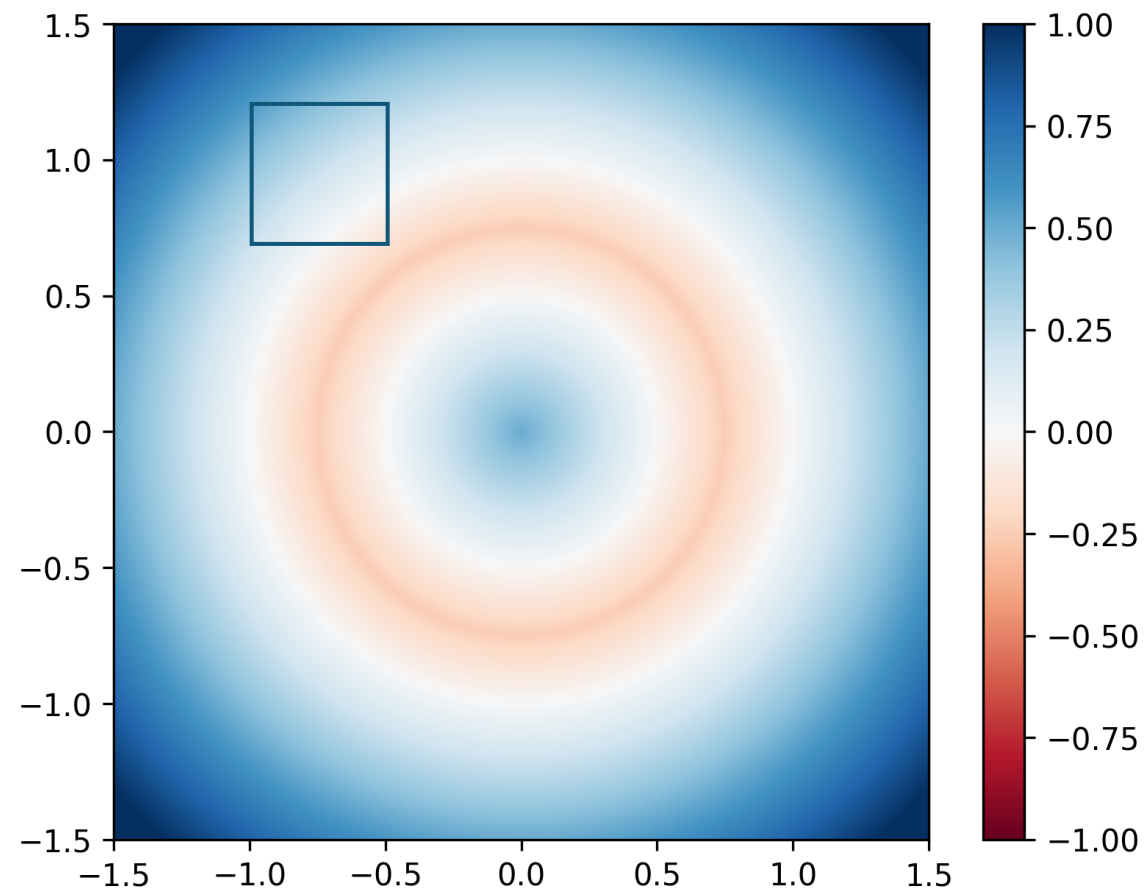
Interval  
arithmetic

Tape  
shortening

Interpreter  
on the GPU

Interval  
arithmetic

Tape  
shortening



$$\max \left( 0.5 - \sqrt{x^2 + y^2}, \sqrt{x^2 + y^2} - 1 \right)$$

# MIN AND MAX ARE SPECIAL

opcode	lhs	rhs	out
X	—	—	$[-1.0, -0.5]$
Y	—	—	$[0.7, 1.2]$
SQUARE	$[-1.0, -0.5]$	—	$[0.25, 1.0]$
SQUARE	$[0.7, 1.2]$	—	$[0.49, 1.44]$
ADD	$[0.25, 1.0]$	$[0.49, 1.44]$	$[0.74, 2.44]$
SQRT	$[0.74, 2.44]$	—	$[0.86, 1.56]$
SUB	$[0.86, 1.56]$	$1.0f$	$[-0.14, 0.56]$
SUB	$0.5f$	$[0.86, 1.56]$	$[-1.06, -0.36]$
MAX	$[-0.14, 0.56]$	$[-1.06, -0.36]$	$[-0.14, 0.56]$

# MIN AND MAX ARE SPECIAL



# MIN AND MAX ARE SPECIAL



# REPLACING WITH COPY

opcode	lhs	rhs	out
X	–	–	slot 0
Y	–	–	slot 1
SQUARE	slot 0	–	slot 0
SQUARE	slot 1	–	slot 1
ADD	slot 0	slot 1	slot 1
SQRT	slot 1	–	slot 1
SUB	slot 1	1.0f	slot 0
SUB	0.5f	slot 1	slot 1
MAX	slot 0	slot 1	slot 1

# REPLACING WITH COPY

opcode	lhs	rhs	out
X	–	–	slot 0
Y	–	–	slot 1
SQUARE	slot 0	–	slot 0
SQUARE	slot 1	–	slot 1
ADD	slot 0	slot 1	slot 1
SQRT	slot 1	–	slot 1
SUB	slot 1	1.0f	slot 0
SUB	0.5f	slot 1	slot 1
COPY	slot 0	–	slot 1



# NOTICING INACTIVE CLAUSES

opcode	lhs	rhs	out
X	–	–	slot 0
Y	–	–	slot 1
SQUARE	slot 0	–	slot 0
SQUARE	slot 1	–	slot 1
ADD	slot 0	slot 1	slot 1
SQRT	slot 1	–	slot 1
SUB	slot 1	1.0f	slot 0
<del>SUB</del>	<del>0.5f</del>	<del>slot 1</del>	<del>slot 1</del>
COPY	slot 0	–	slot 1

# NOTICING INACTIVE CLAUSES

opcode	lhs	rhs	out
X	–	–	slot 0
Y	–	–	slot 1
SQUARE	slot 0	–	slot 0
SQUARE	slot 1	–	slot 1
ADD	slot 0	slot 1	slot 1
SQRT	slot 1	–	slot 1
SUB	slot 1	1.0f	slot 0
COPY	slot 0	–	slot 1

# TAPE SHORTENING

- Use interval results at **min** / **max** clauses to detect clauses that are inactive in a particular region
- Construct a shortened tape without those clauses
- That tape is valid for all subregions contained within the parent region

---

**Algorithm 1:** Evaluate a tape, recording which branch is taken at every min and max node

---

```
choices ← an empty stack
foreach clause in tape do
  lhs ← getValue(clause.lhs)
  rhs ← getValue(clause.rhs)
  switch clause.opcode do
    case OP_MIN
      if lhs.upper < rhs.lower then
        | choices.push(CHOICE_LHS)
      else if rhs.upper < lhs.lower then
        | choices.push(CHOICE_RHS)
      else
        | choices.push(CHOICE_BOTH)
      slots[clause.out] ← min(lhs, rhs)
    case OP_MAX
      | Similar logic to push a choice slots[clause.out]
      | ← max(lhs, rhs)
    case OP_ADD
      | slots[clause.out] ← lhs + rhs
    case OP_SUB
      | ...and so on for other opcodes
  end
end
clause ← the last clause in the tape
return (slots[clause.out], choices)
```

---

---

**Algorithm 2:** Use the choice data from Alg. 1 to construct a shorter tape which only contains active clauses

---

```
output ← an empty tape
active ← an array of all false
active[final output slot] ← true
foreach clause in tape.reversed() do
  if clause.opcode ∈ [OP_MIN, OP_MAX] then
    | choice ← choices.pop()
  else
    | choice ← CHOICE_BOTH
  if active[clause.out] then
    active[clause.out] ← false
    if choice == CHOICE_LHS then
      | active[clause.lhs] ← true
      | clause.rhs ← clause.lhs
    else if choice == CHOICE_RHS then
      | active[clause.rhs] ← true
      | clause.lhs ← clause.rhs
    else
      | active[clause.lhs] ← true
      | active[clause.rhs] ← true
    output.push_back(clause)
end
return output
```

---

# TAPE SHORTENING

- Analogous to *mark-and-sweep* garbage collection
- Mark output slot as **active**, all other slots as **inactive**
- Walk *backwards* through the tape, skipping clauses with **inactive** output slots
- For **MIN** and **MAX** clauses:
  - Mark the output slot as **inactive**
  - If unambiguous, only mark *one* input output slot as **active**, and replace **MIN/MAX** with **COPY**
  - Otherwise, mark *all* input slots as **active**
- For all **other** clauses
  - Mark the output slot as **inactive**
  - Mark *all* input slots as **active**

# HOW WELL DOES THIS WORK?

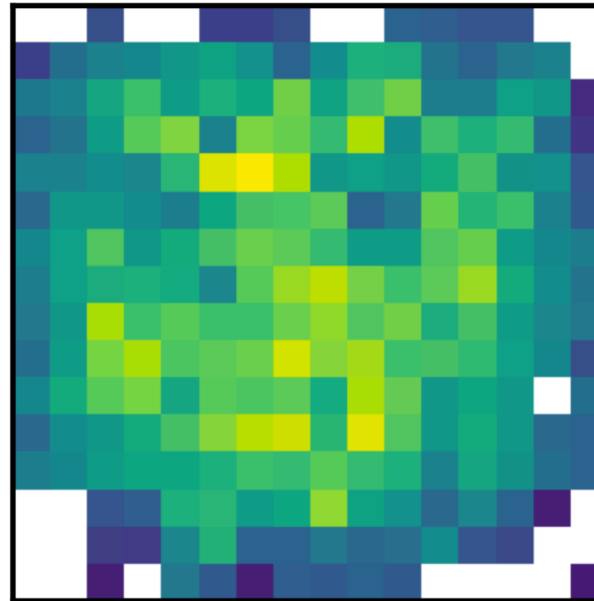
But this rough magic I  
here abjure, and when  
I have required some  
heavenly music, which even  
now I do, to work mine  
end upon their senses that  
this airy charm is for, I'll  
break my staff, bury it  
certain fathoms in the  
earth, and deeper than did  
ever plummet sound  
I'll drown my book.

Original: 6056 clauses  
(rendered at 1024x1024)

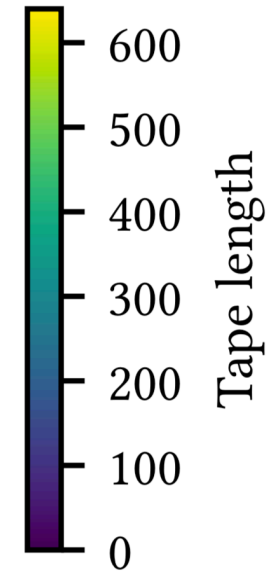
# HOW WELL DOES THIS WORK?

But this rough magic I  
here abjure, and when  
I have required some  
heavenly music, which even  
now I do, to work mine  
end upon their senses that  
this airy charm is for, I'll  
break my staff, bury it  
certain fathoms in the  
earth, and deeper than did  
ever plummet sound  
I'll drown my book.

Original: 6056 clauses  
(rendered at 1024x1024)



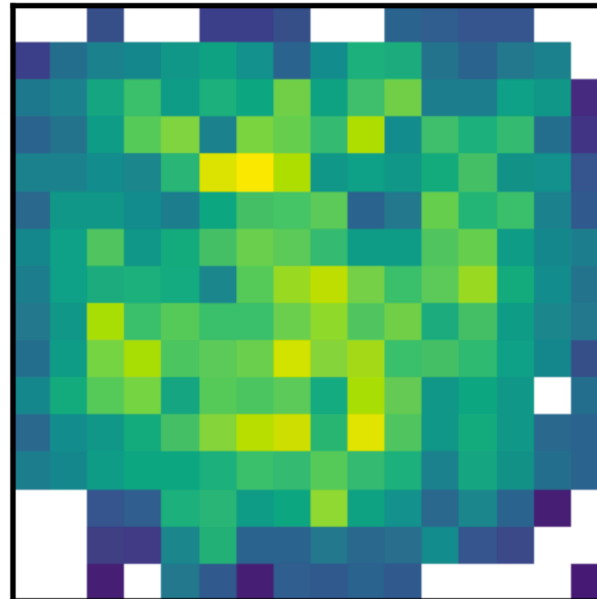
$356 \pm 125$  clauses



# HOW WELL DOES THIS WORK?

But this rough magic I  
here abjure, and when  
I have required some  
heavenly music, which even  
now I do, to work mine  
end upon their senses that  
this airy charm is for, I'll  
break my staff, bury it  
certain fathoms in the  
earth, and deeper than did  
ever plummet sound  
I'll drown my book.

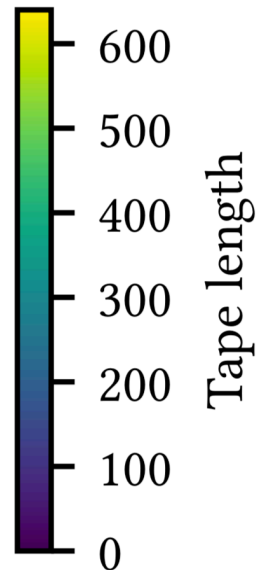
Original: 6056 clauses  
(rendered at 1024x1024)



$356 \pm 125$  clauses



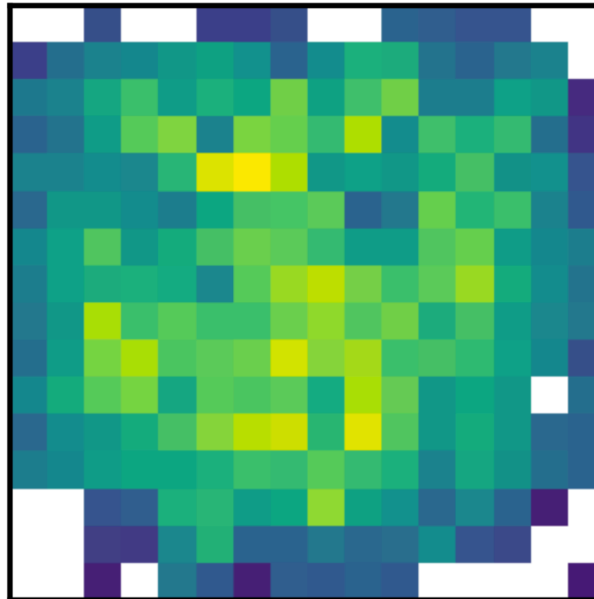
$28 \pm 13$  clauses



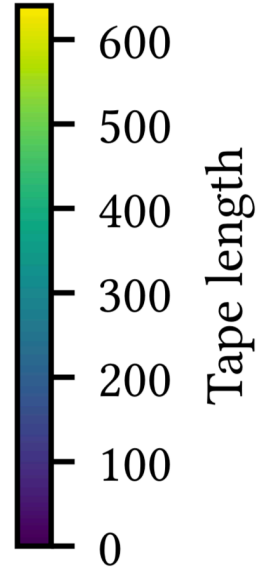
# SHORTENING TAPES IN PRACTICE

- We no longer have a single input tape!
- Each region can build a shortened tape
- Dynamic memory allocation on the GPU?

But this rough magic I  
here abjure, and when  
I have required some  
heavenly music, which even  
now I do, to work mine  
end upon their senses that  
this airy charm is for, I'll  
break my staff, bury it  
certain fathoms in the  
earth, and deeper than did  
ever plummet sound  
I'll drown my book.



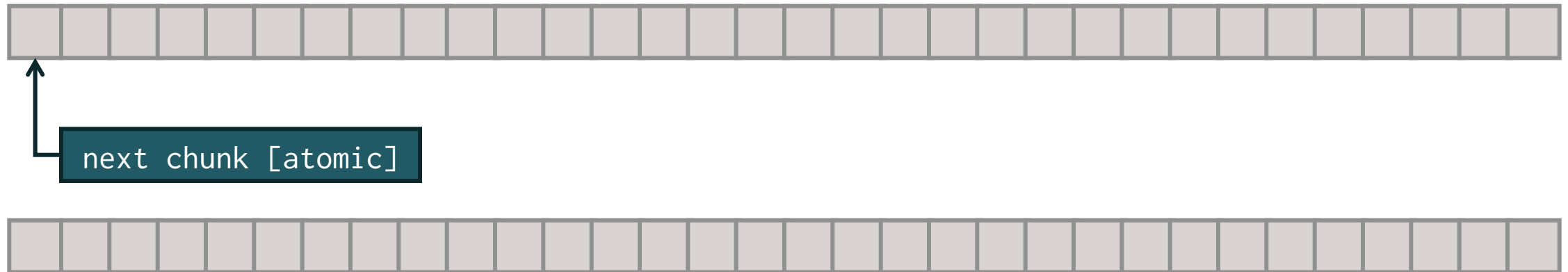
But this rough magic I  
here abjure, and when  
I have required some  
heavenly music, which even  
now I do, to work mine  
end upon their senses that  
this airy charm is for, I'll  
break my staff, bury it  
certain fathoms in the  
earth, and deeper than did  
ever plummet sound  
I'll drown my book.





# STORING TAPES ON THE GPU

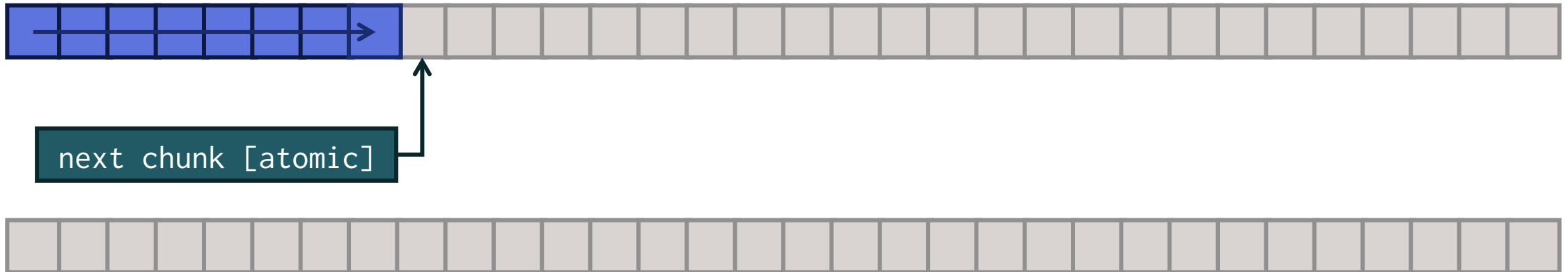
- Unrolled linked lists (64 clauses per chunk)
- Stored in a large (> 1 GB) pre-allocated array in **global** memory
- Threads claim memory by bumping an atomic integer



# STORING TAPES ON THE GPU

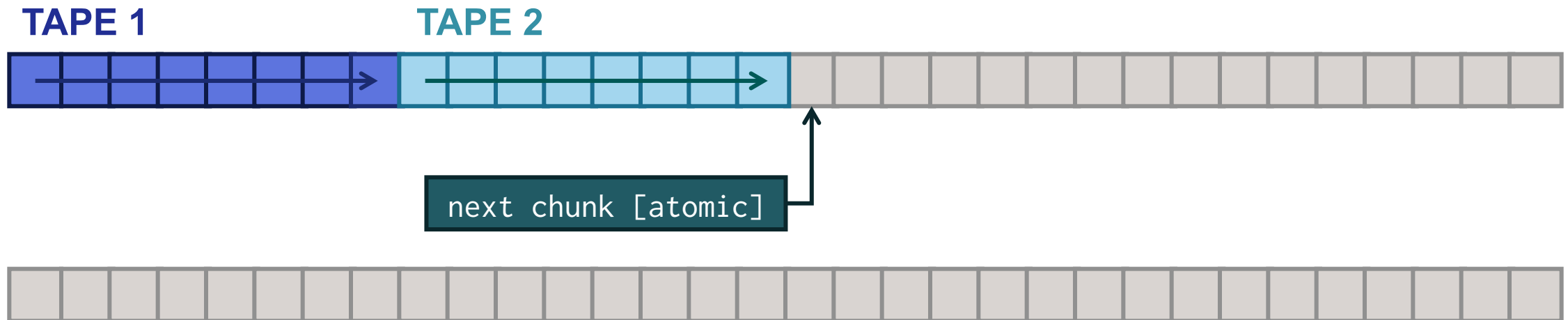
- Unrolled linked lists (64 clauses per chunk)
- Stored in a large (> 1 GB) pre-allocated array in **global** memory
- Threads claim memory by bumping an atomic integer

## TAPE 1



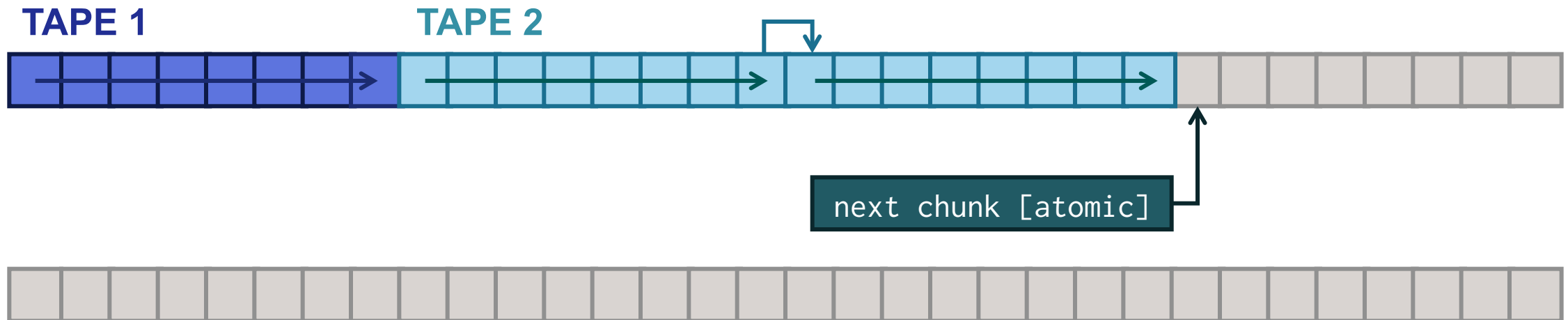
# STORING TAPES ON THE GPU

- Unrolled linked lists (64 clauses per chunk)
- Stored in a large (> 1 GB) pre-allocated array in **global** memory
- Threads claim memory by bumping an atomic integer



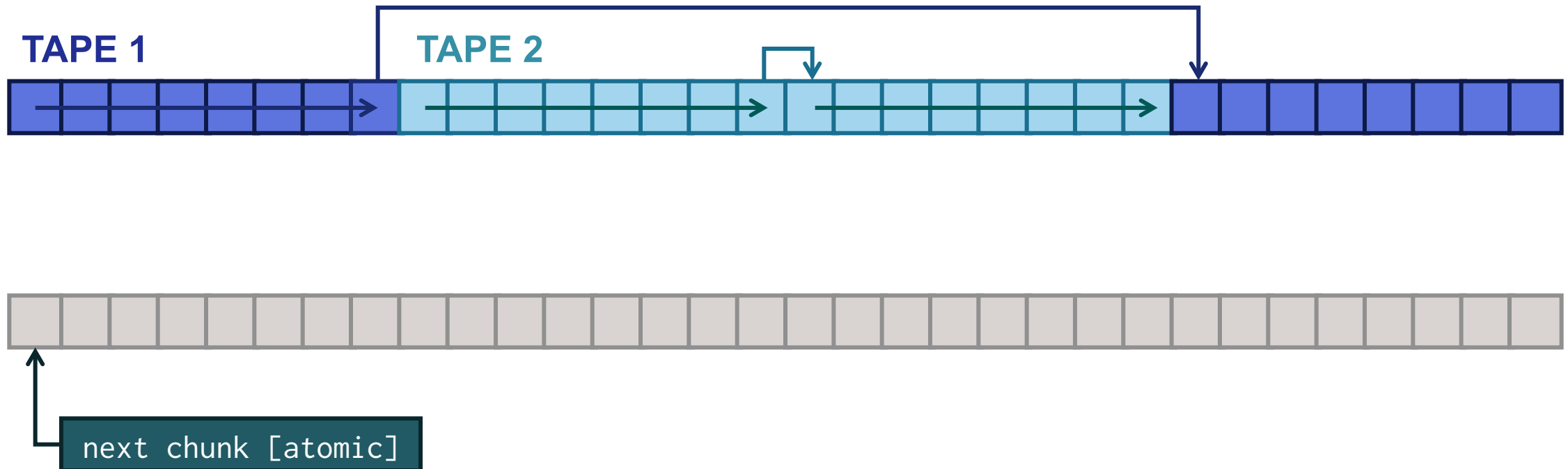
# STORING TAPES ON THE GPU

- Unrolled linked lists (64 clauses per chunk)
- Stored in a large (> 1 GB) pre-allocated array in **global** memory
- Threads claim memory by bumping an atomic integer



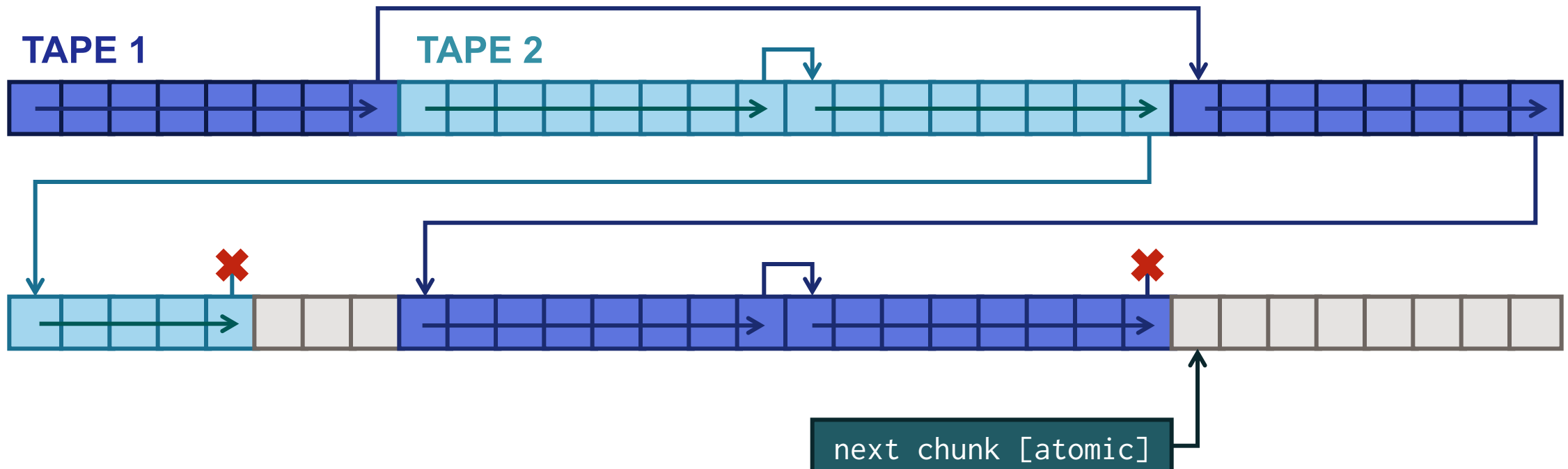
# STORING TAPES ON THE GPU

- Unrolled linked lists (64 clauses per chunk)
- Stored in a large (> 1 GB) pre-allocated array in **global** memory
- Threads claim memory by bumping an atomic integer



# STORING TAPES ON THE GPU

- Unrolled linked lists (64 clauses per chunk)
- Stored in a large (> 1 GB) pre-allocated array in **global** memory
- Threads claim memory by bumping an atomic integer



opcode	out	lhs	rhs	immediate <i>or</i> jump
u8	u8	u8	u8	f32 / i32
JUMP	—	—	—	128

# THE INTERPRETER LOOP

```
def run(tape):  
    i = 0  
    loop:  
        (opcode, lhs, rhs, out, imm) = tape[i]  
        match opcode:  
            DONE: break  
            JUMP: i += imm  
            ADD_LHS_RHS: slots[out] = slots[lhs] + slots[rhs]  
            ...and so on...  
        i += 1  
    ...
```



Interpreter  
on the GPU

Interval  
arithmetic

Tape  
shortening

Interpreter  
on the GPU

Interval  
arithmetic

Tape  
shortening

**PUTTING IT ALL TOGETHER!**

Shape

Encode &  
allocate slots

Instruction tape

Send to GPU

Interval  
evaluation

Tape shortening

Pixel  
evaluation

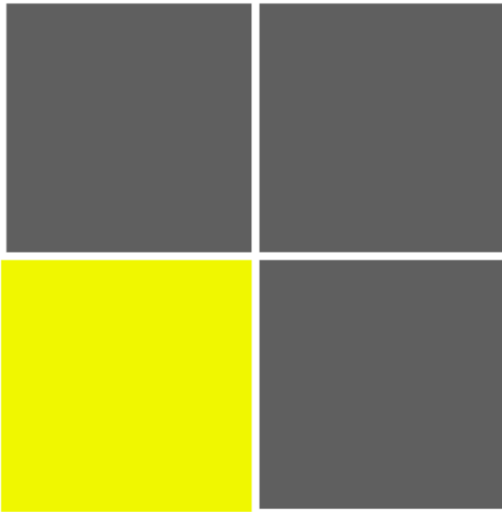
Output image

Subdivide & recurse

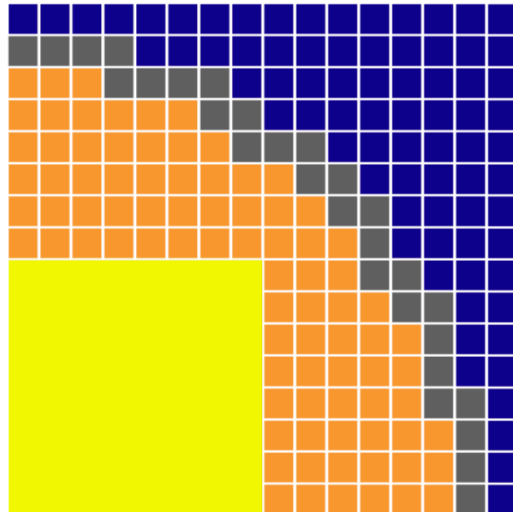
$64^2$ ,  $8^2$  pixel regions

GPU

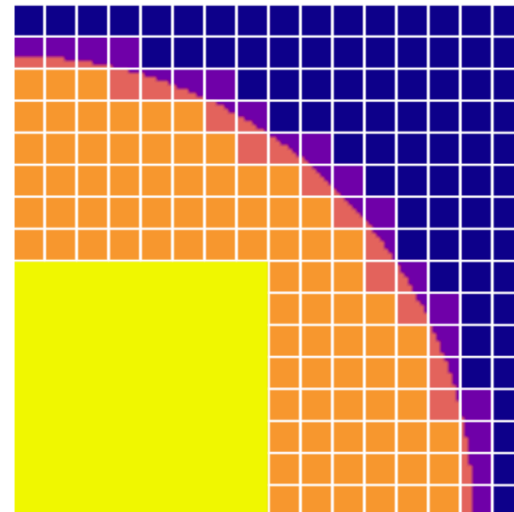
First pass



Second pass



Per-pixel evaluation



64x64 filled tile



8x8 filled tile



8x8 empty tile

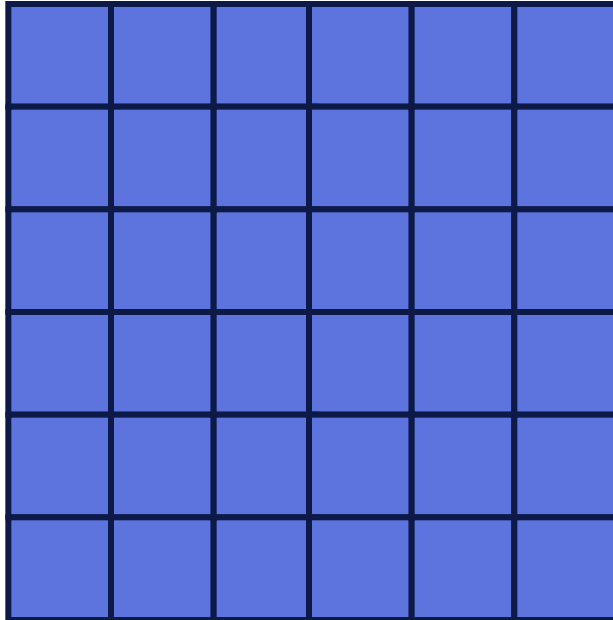


8x8 ambiguous tile

# 2D RENDERING PROCEDURE

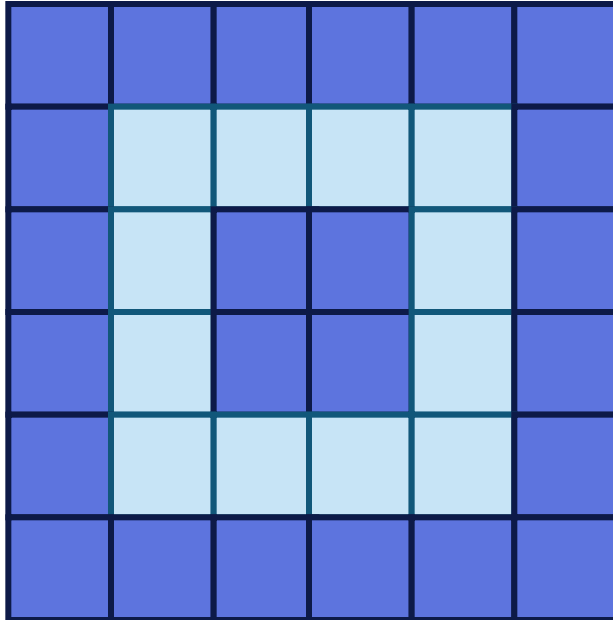
- **Interval evaluation** of 64x64-pixel **tiles**
  - Record empty and filled tiles, but do nothing more with them
  - For ambiguous tiles, **calculate shortened tape**
  - Collect ambiguous tiles into a dense list for further evaluation
- Split ambiguous tiles into 64 **subtiles** (8x8-pixels) each
- **Interval evaluation** of subtiles
  - Record empty and filled subtiles, but do nothing more with them
  - For ambiguous subtiles, **calculate shortened tape**
  - Collect ambiguous subtiles into a dense list for pixel evaluation
- **Per-pixel** evaluations of ambiguous subtiles
  - Same interpreter loop as interval evaluation, but using single values instead

# MEMORY LAYOUT



- Largest tiles are **densely packed**
  - $(N / 64)^D$  tiles, where N is image size and D is dimension
- Each tile contains three values
  - Tile position (packed 32-bit int, with 10 bits each for x/y/z)
  - Index of specialized tape for this tile (-1 by default)
  - Index of subdivided region for this tile (-1 by default)

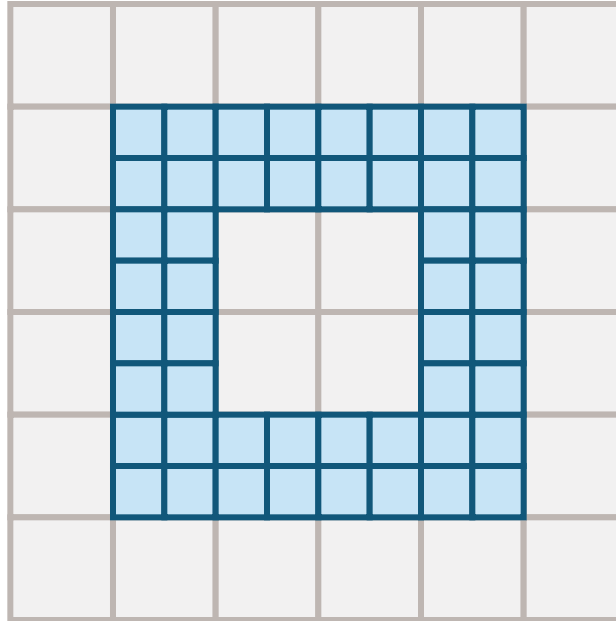
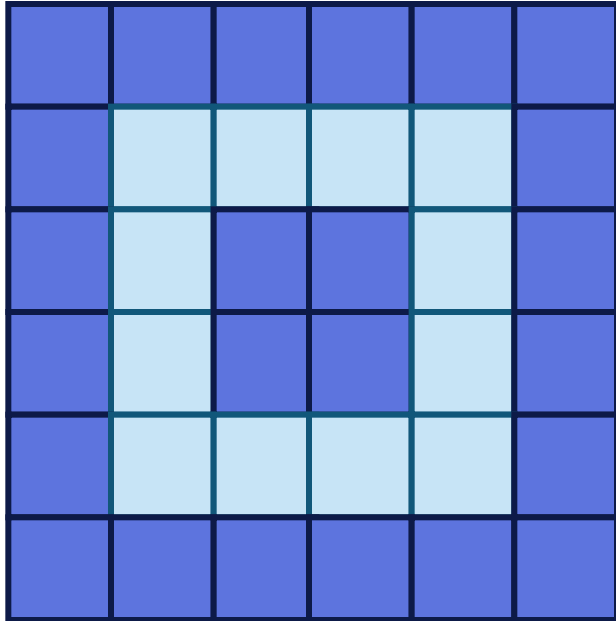
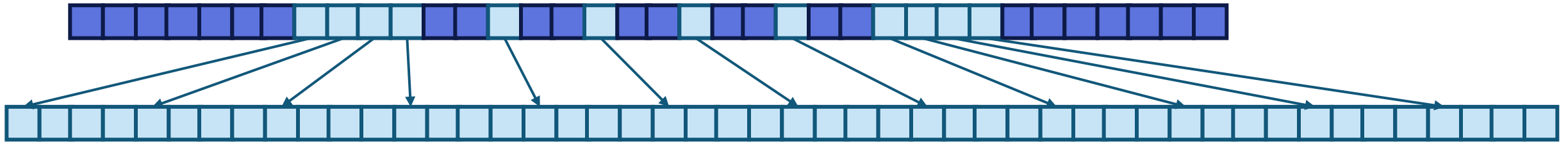
# MEMORY LAYOUT



- Interval evaluation picks out ambiguous tiles
- These tiles will be subdivided and evaluated again
- Subdivided tiles are **sparsely distributed** in space but **densely packed** in GPU RAM

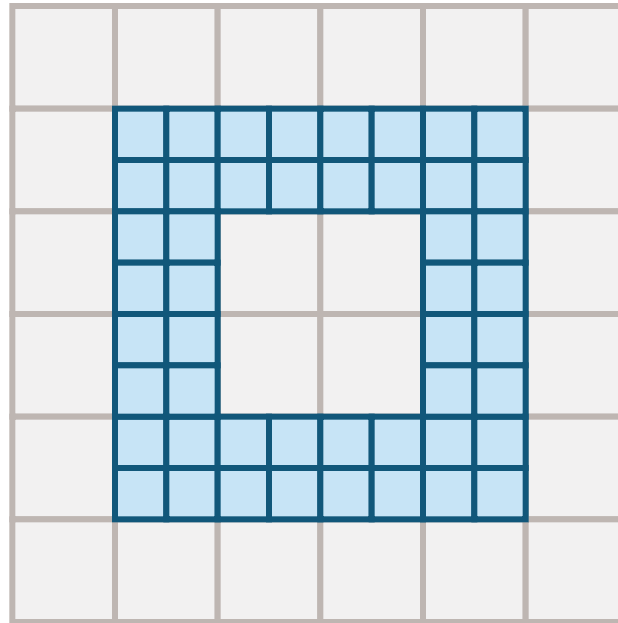
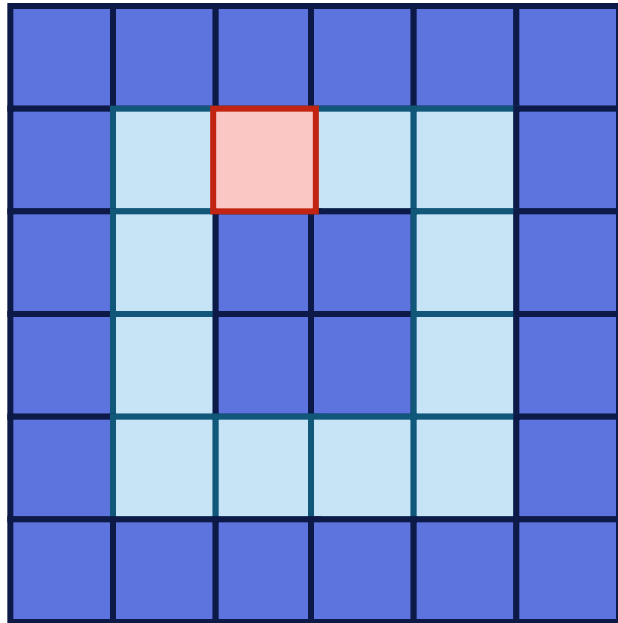
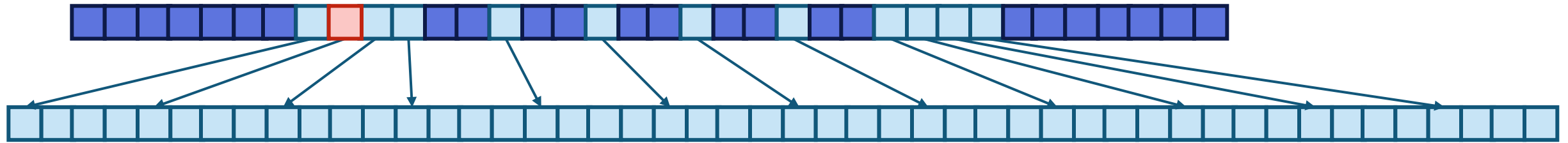


# MEMORY LAYOUT



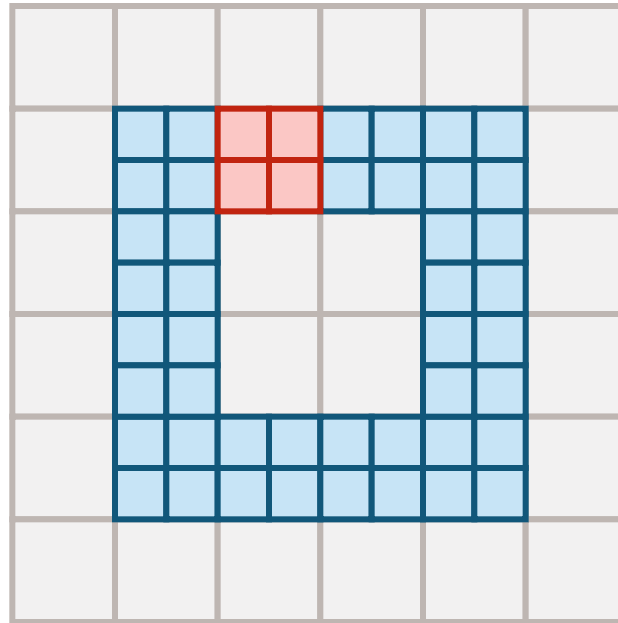
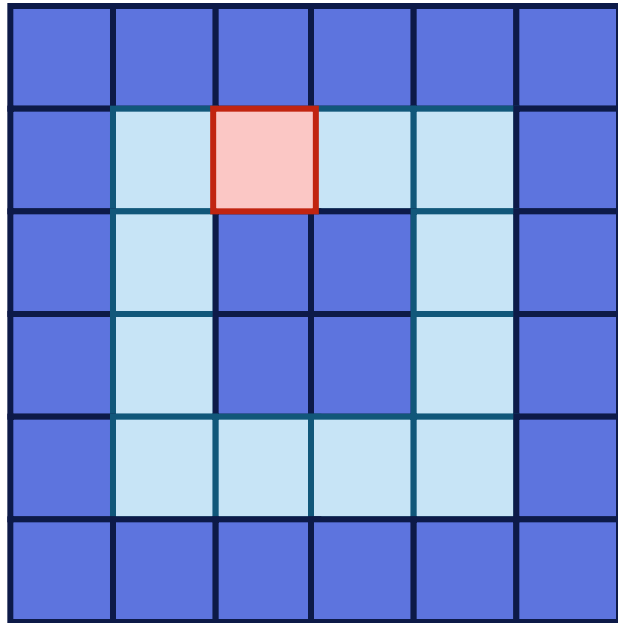
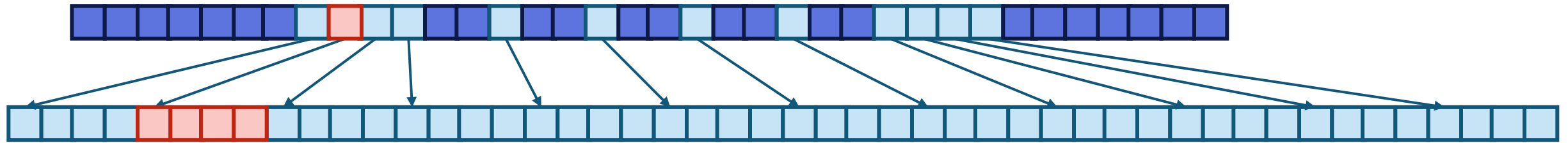
(Diagram shows 4x subdivision,  
but we use 64x in practice)

# STORING THE HIERARCHY



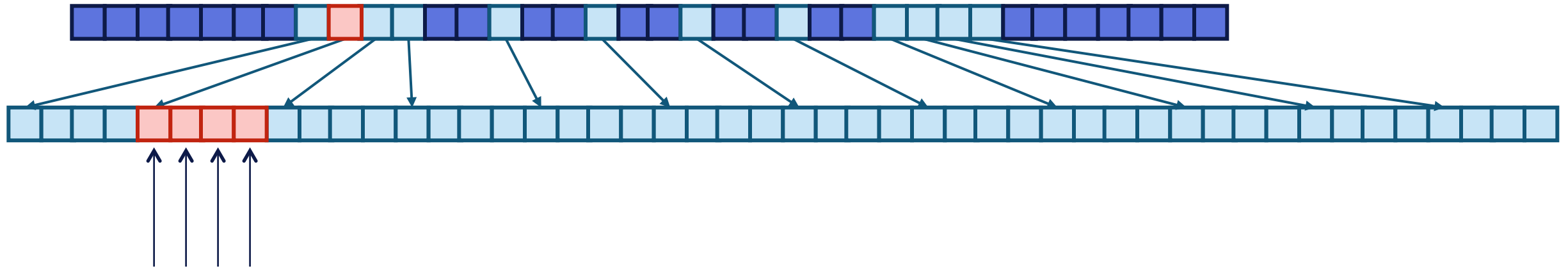
(Diagram shows 4x subdivision,  
but we use 64x in practice)

# MEMORY LAYOUT



(Diagram shows 4x subdivision,  
but we use 64x in practice)

# MEMORY LAYOUT



64 subtiles = 64 threads = 2 warps

All threads within these warps use the **same shortened tape** (from parent tile), to minimize divergence

# BENCHMARKING MACHINES

GeForce GT 750M  
2013 MacBook Pro

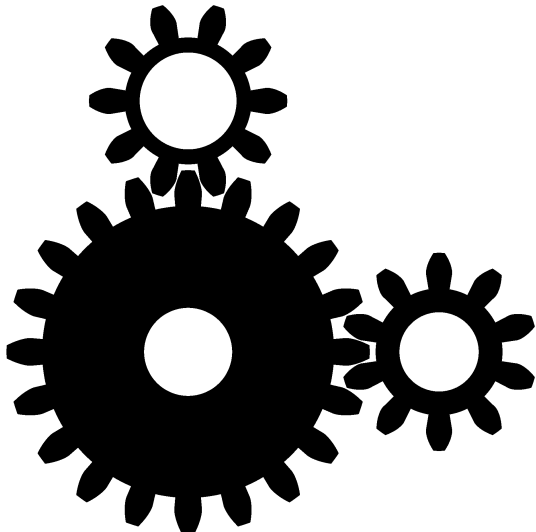
GTX 1080 Ti  
2017 VR/ML workstation

Tesla V100  
AWS p3.2xlarge

But this rough magic I  
here abjure, and when  
I have required some  
heavenly music, which even  
now I do, to work mine  
end upon their senses that  
this airy charm is for, I'll  
break my staff, bury it  
certain fathoms in the  
earth, and deeper than did  
ever plummet sound  
I'll drown my book.

## Text benchmark

- 6056 clauses (2354 min/max)
- Shows off CSG and tape pruning performance



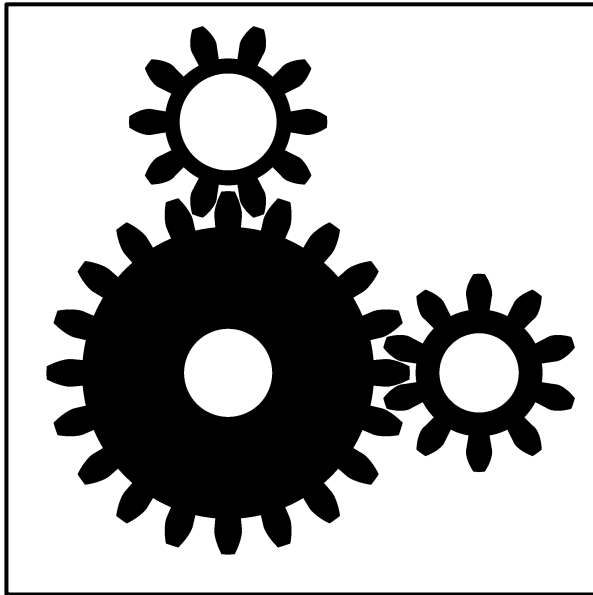
## 2D Gears

- 1735 clauses (374 min/max)
- Curves are mathematically correct involutes, which requires transcendental functions

But this rough magic I  
here abjure, and when  
I have required some  
heavenly music, which even  
now I do, to work mine  
end upon their senses that  
this airy charm is for, I'll  
break my staff, bury it  
certain fathoms in the  
earth, and deeper than did  
ever plummet sound  
I'll drown my book.

## Frame time (ms)

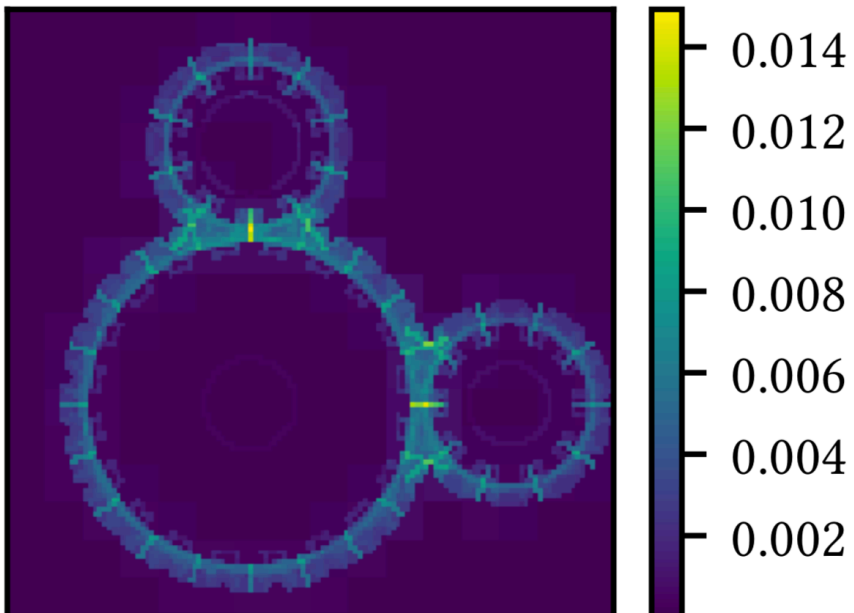
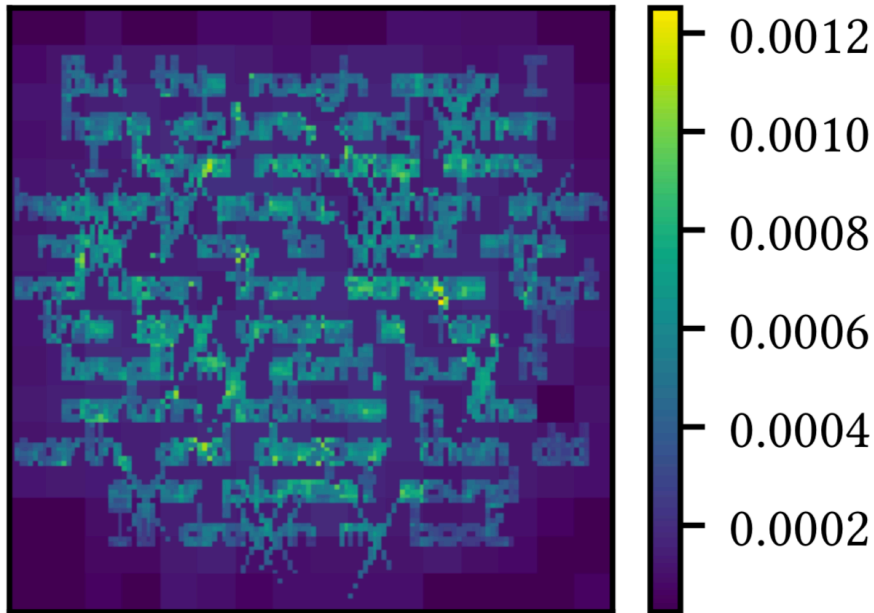
Size	GT 750M	GTX 1080 Ti	Tesla V100
256 <sup>2</sup>	17.5	8.3	5.2
512 <sup>2</sup>	14.8	6.8	4.2
1024 <sup>2</sup>	16.5	6.5	3.9
2048 <sup>2</sup>	20.7	6.6	3.9
3072 <sup>2</sup>	27.1	6.9	3.9
4096 <sup>2</sup>	35.9	7.4	4.1



## Frame time (ms)

Size	GT 750M	GTX 1080 Ti	Tesla V100
256 <sup>2</sup>	9.2	4.0	2.8
512 <sup>2</sup>	9.3	3.7	2.5
1024 <sup>2</sup>	12.1	3.4	2.2
2048 <sup>2</sup>	17.3	3.4	2.2
3072 <sup>2</sup>	23.4	3.7	2.3
4096 <sup>2</sup>	30.6	4.0	2.4





## Metric: Normalized + amortized work

- A score of 1 means that the tape is fully walked once for that pixel
- Interval evaluation work is amortized across all pixels in the interval
- Less than 1 full-tape evaluation per pixel!

# RENDERING IN 3D

Shape

Encode &  
allocate registers

Instruction tape

Send to GPU

Interval  
evaluation

Tape shortening

Voxel  
evaluation

Output  
heightmap

Subdivide & recurse

$64^3$ ,  $16^3$ ,  $4^3$  voxel regions

Calculate  
normals

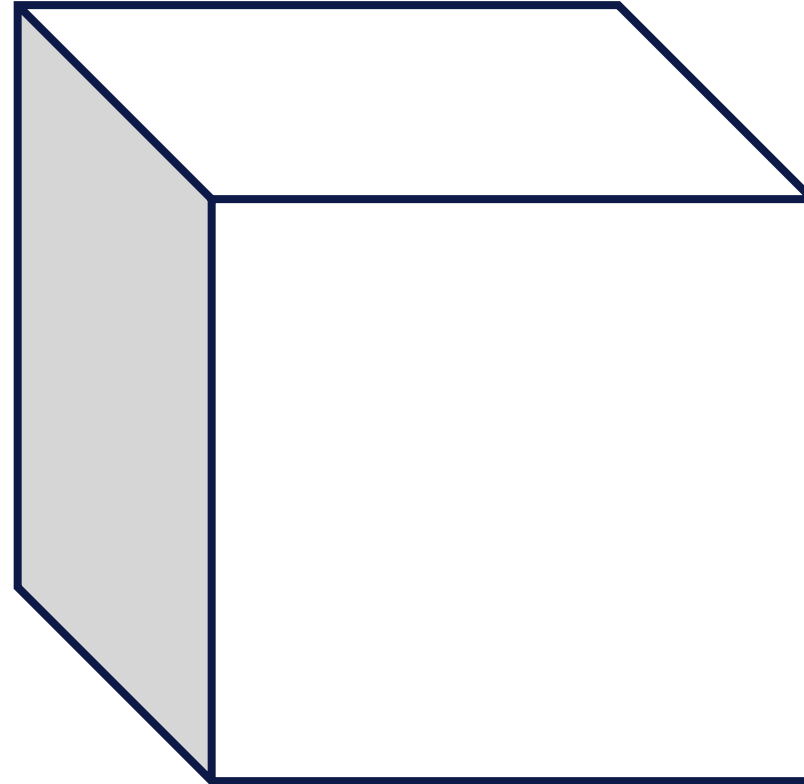
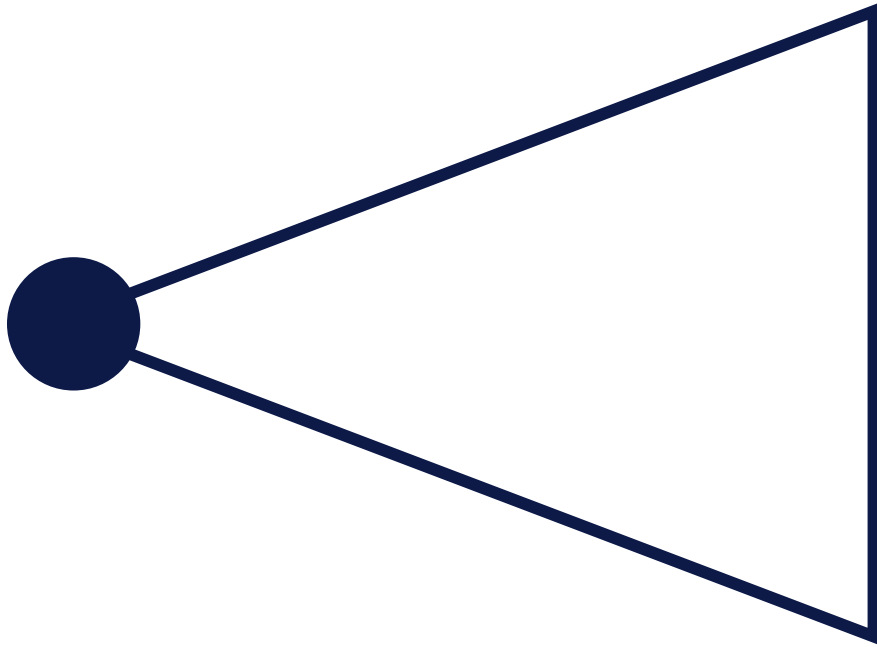
Output normals

GPU

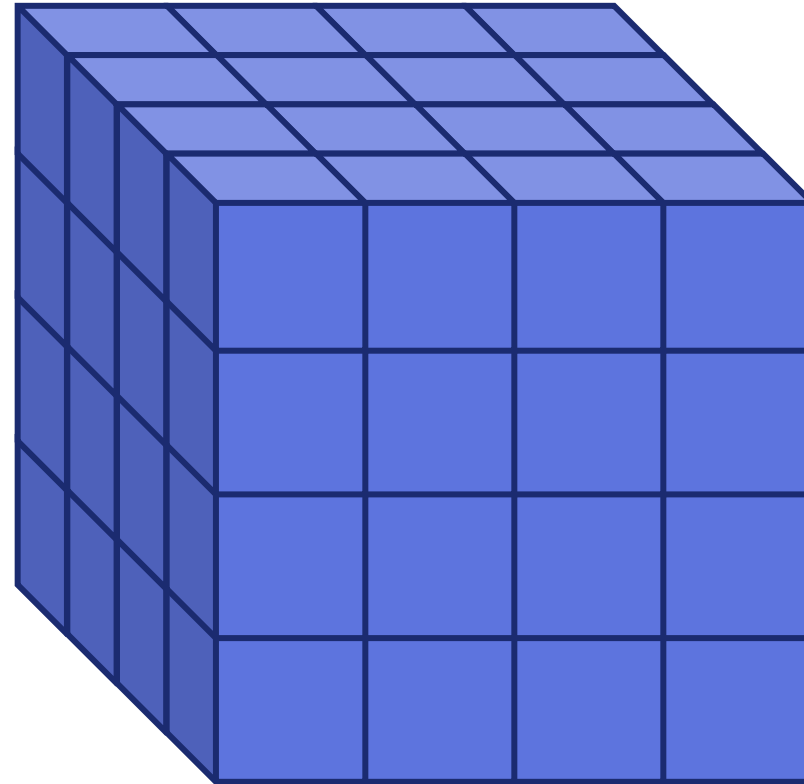
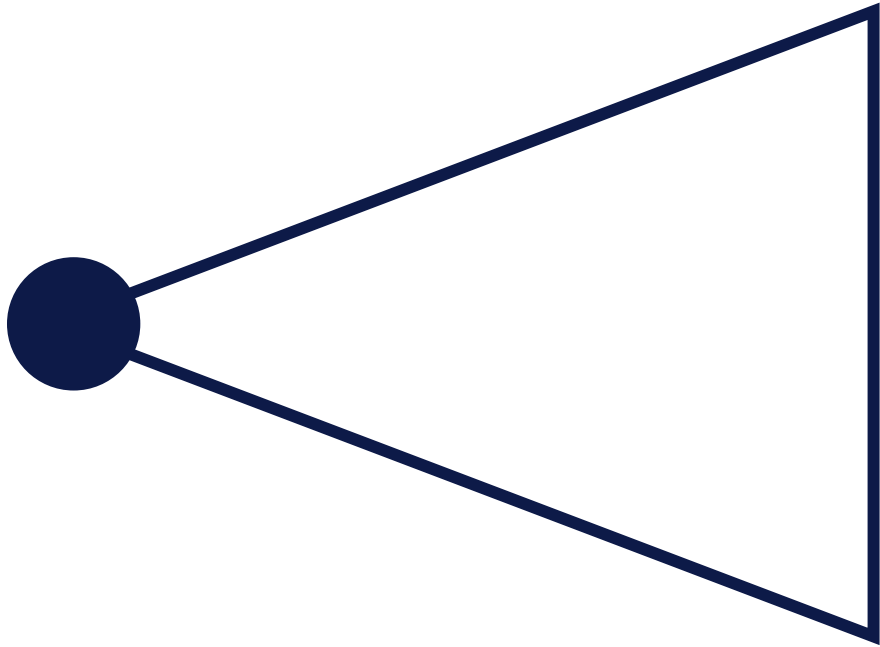
# 3D RENDERING PROCEDURE

- Interval evaluation of 64x64x64-voxel **tiles**
  - Split ambiguous tiles into 64 **subtiles** (16x16x16-voxels) each
- Interval evaluation of **subtiles**
  - Split ambiguous subtiles into 64 **microtiles** (4x4x4-voxels) each
- Interval evaluation of **microtiles**
- Per-**voxel** evaluations of remaining ambiguous **microtiles**
- Per-**pixel** evaluation, using automatic differentiation to find normals
- *Optional*: post-processing depth + normal buffer to generate final image

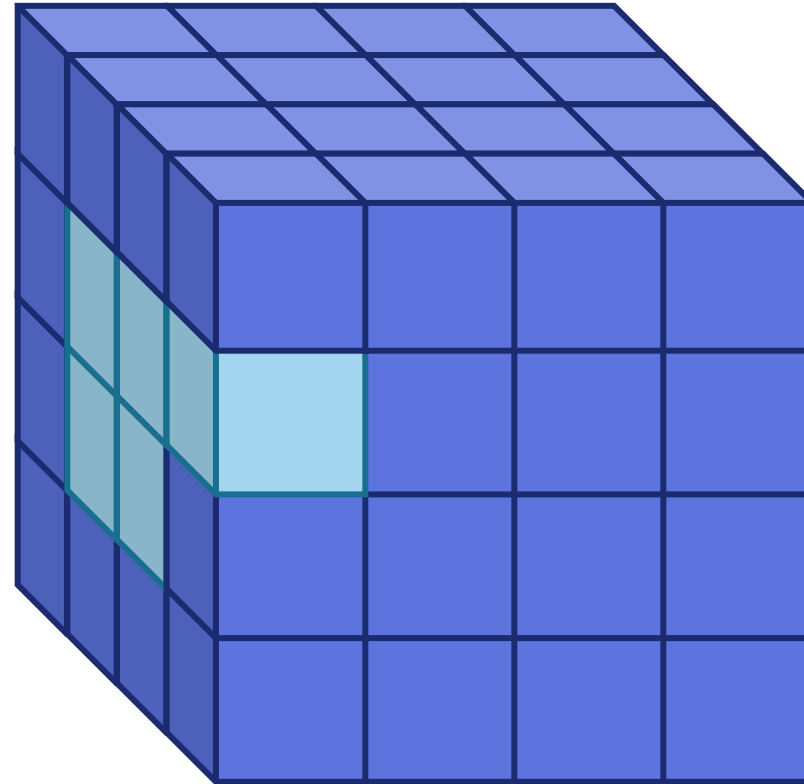
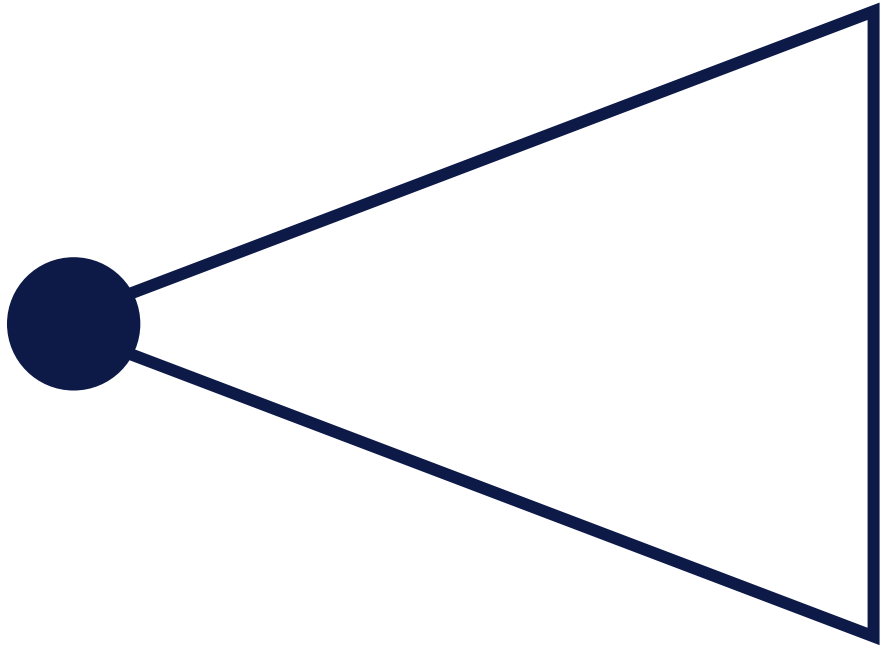
# RENDERING A HEIGHTMAP



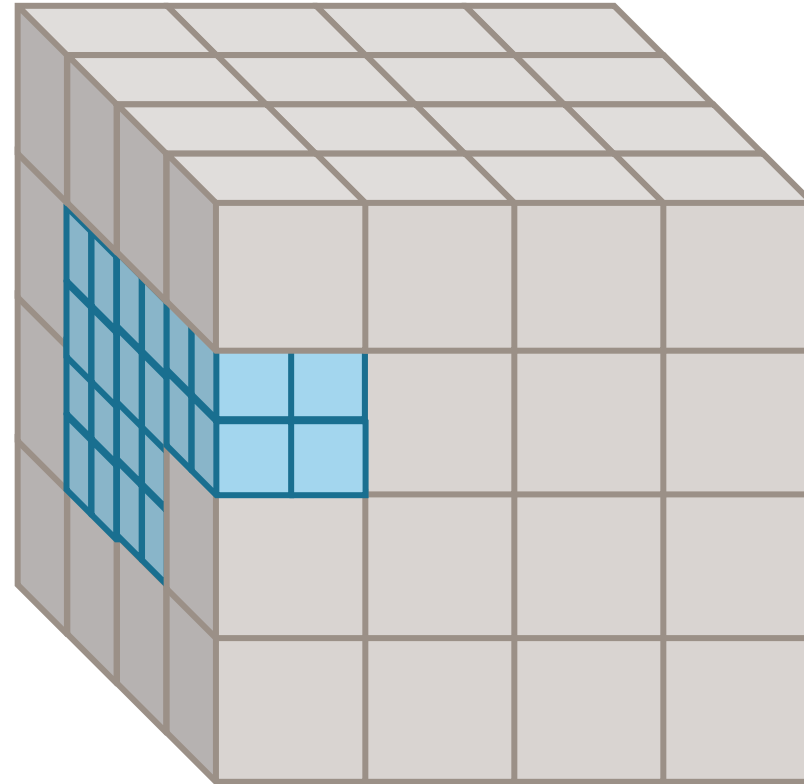
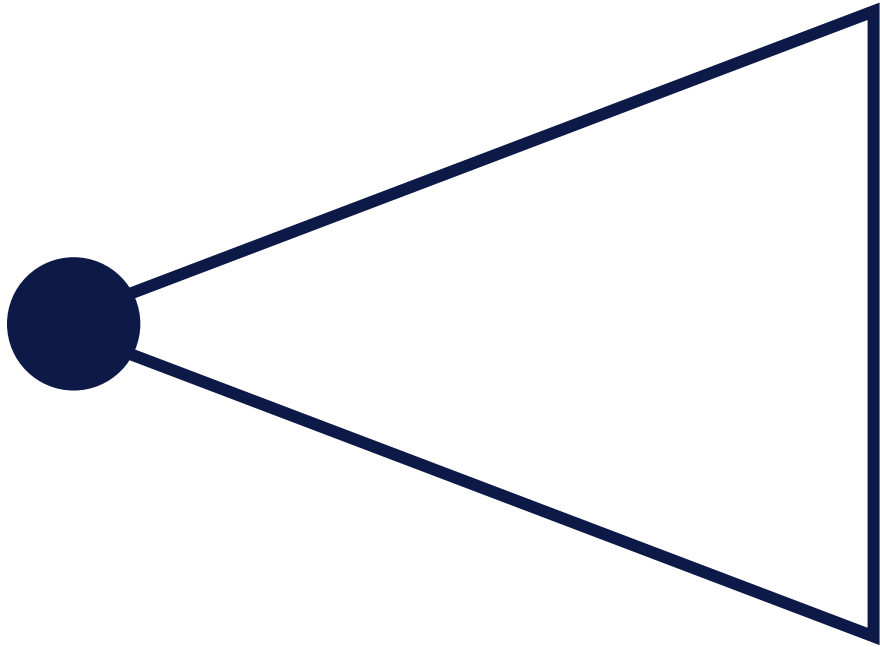
# RENDERING A HEIGHTMAP



# RENDERING A HEIGHTMAP

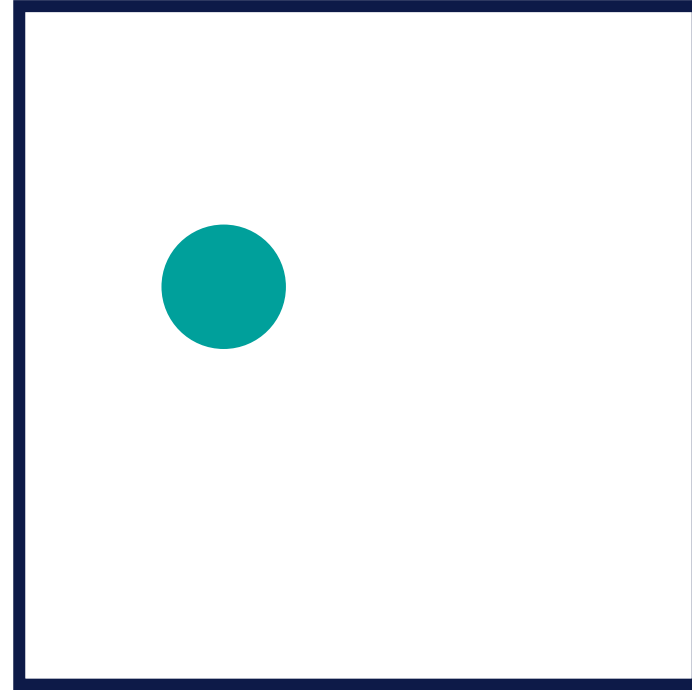
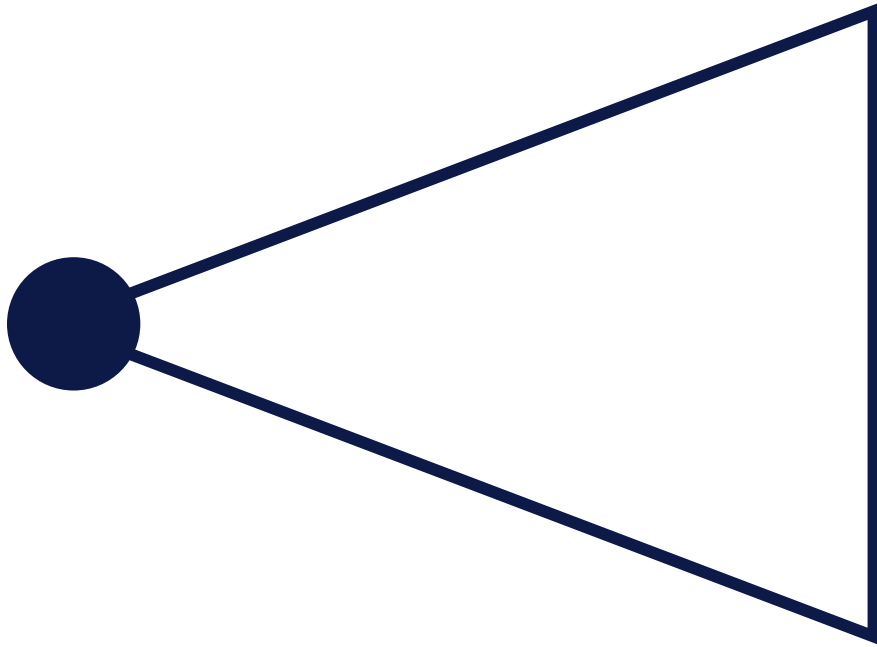


# RENDERING A HEIGHTMAP

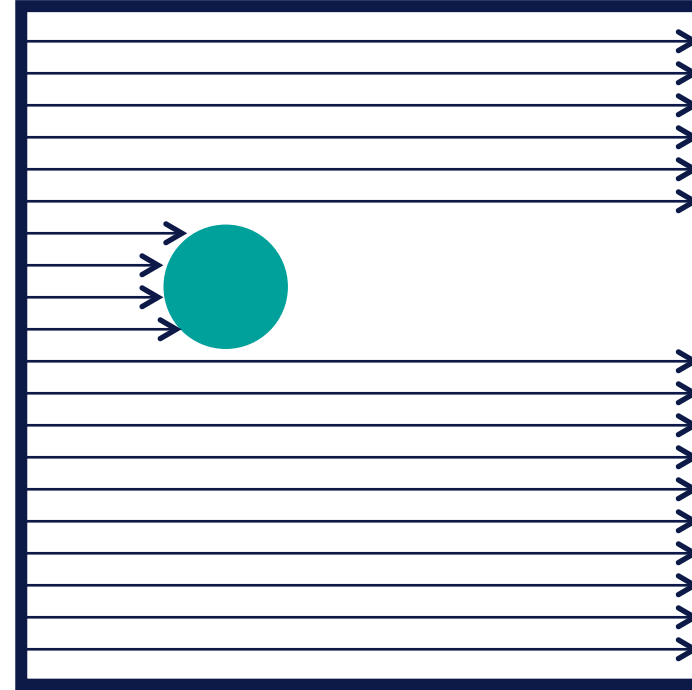
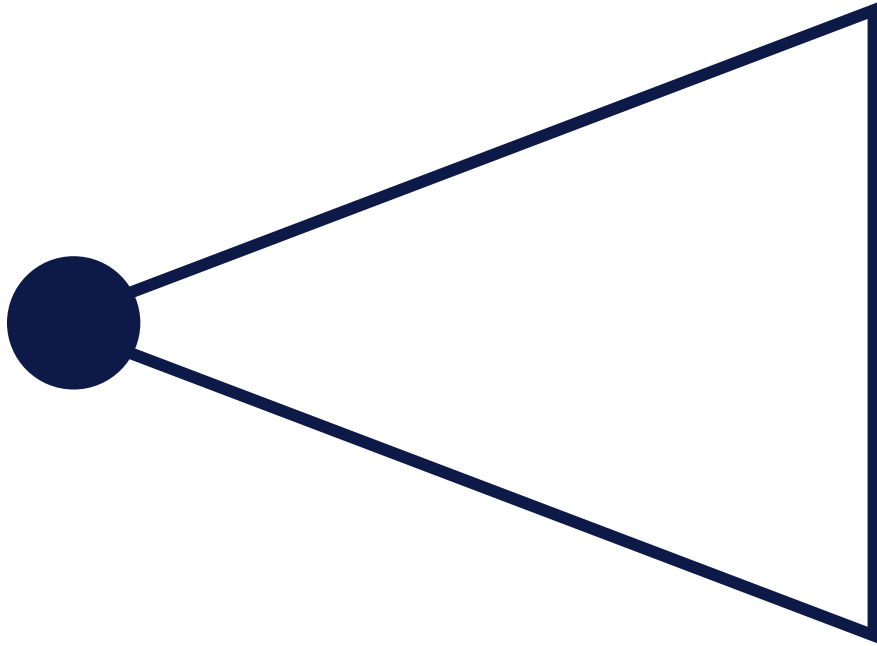




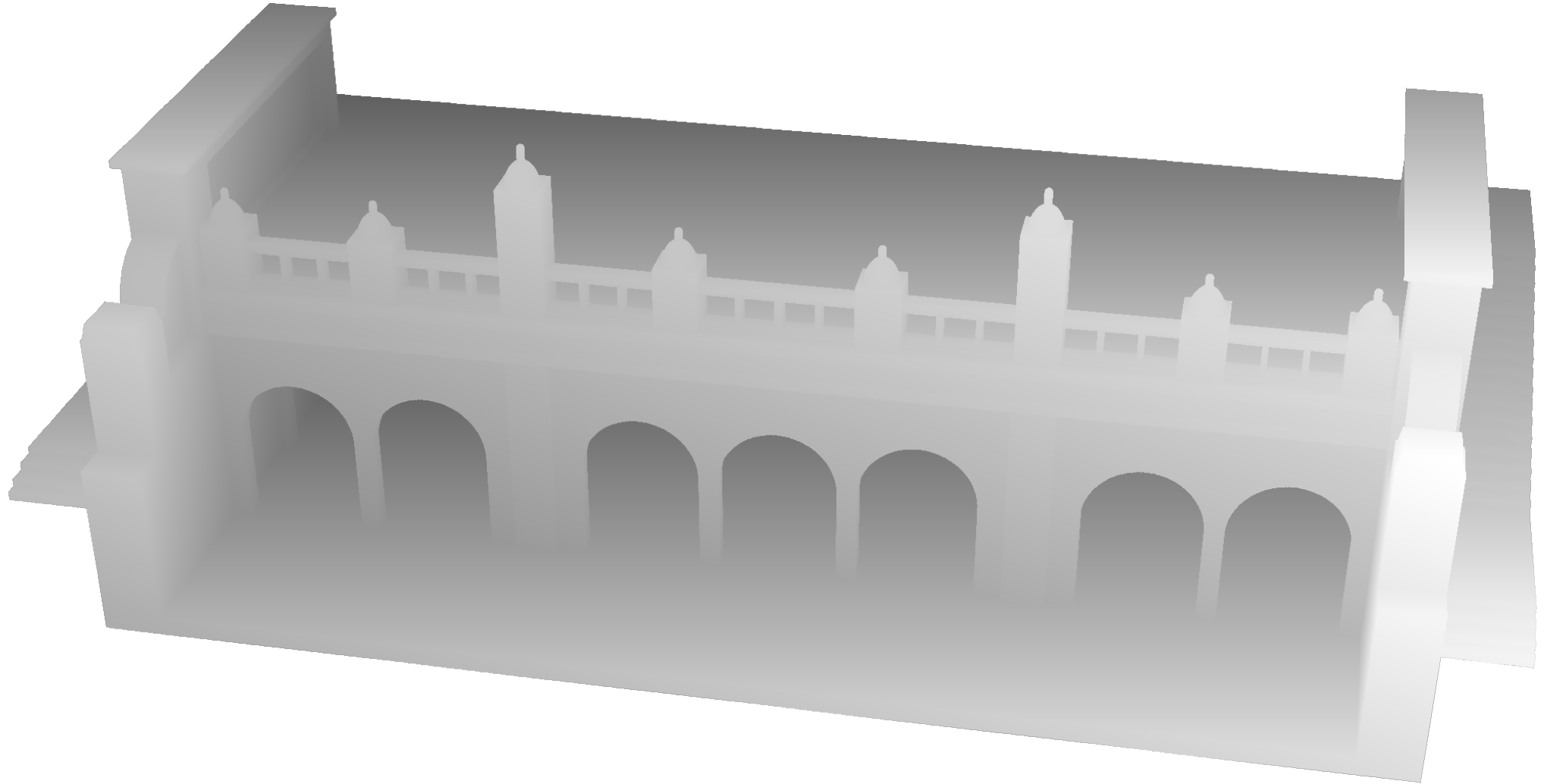
# ORTHOGRAPHIC HEIGHTMAP



# ORTHOGRAPHIC HEIGHTMAP



# HEIGHTMAP OUTPUT

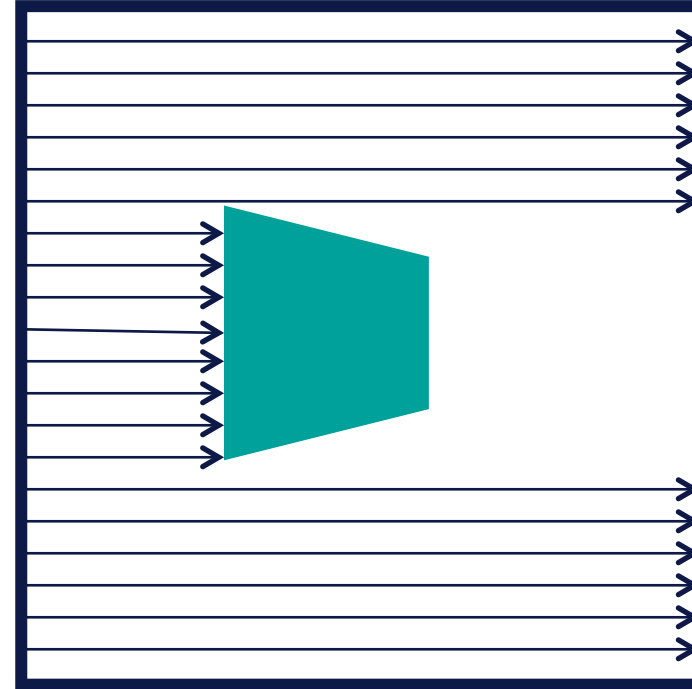
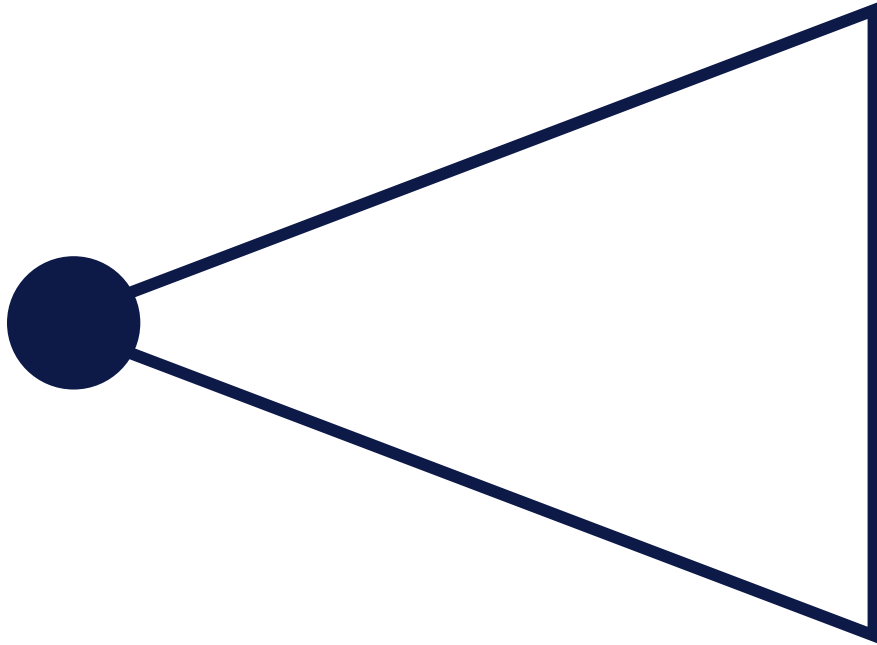


# USER-PROVIDED MATRIX

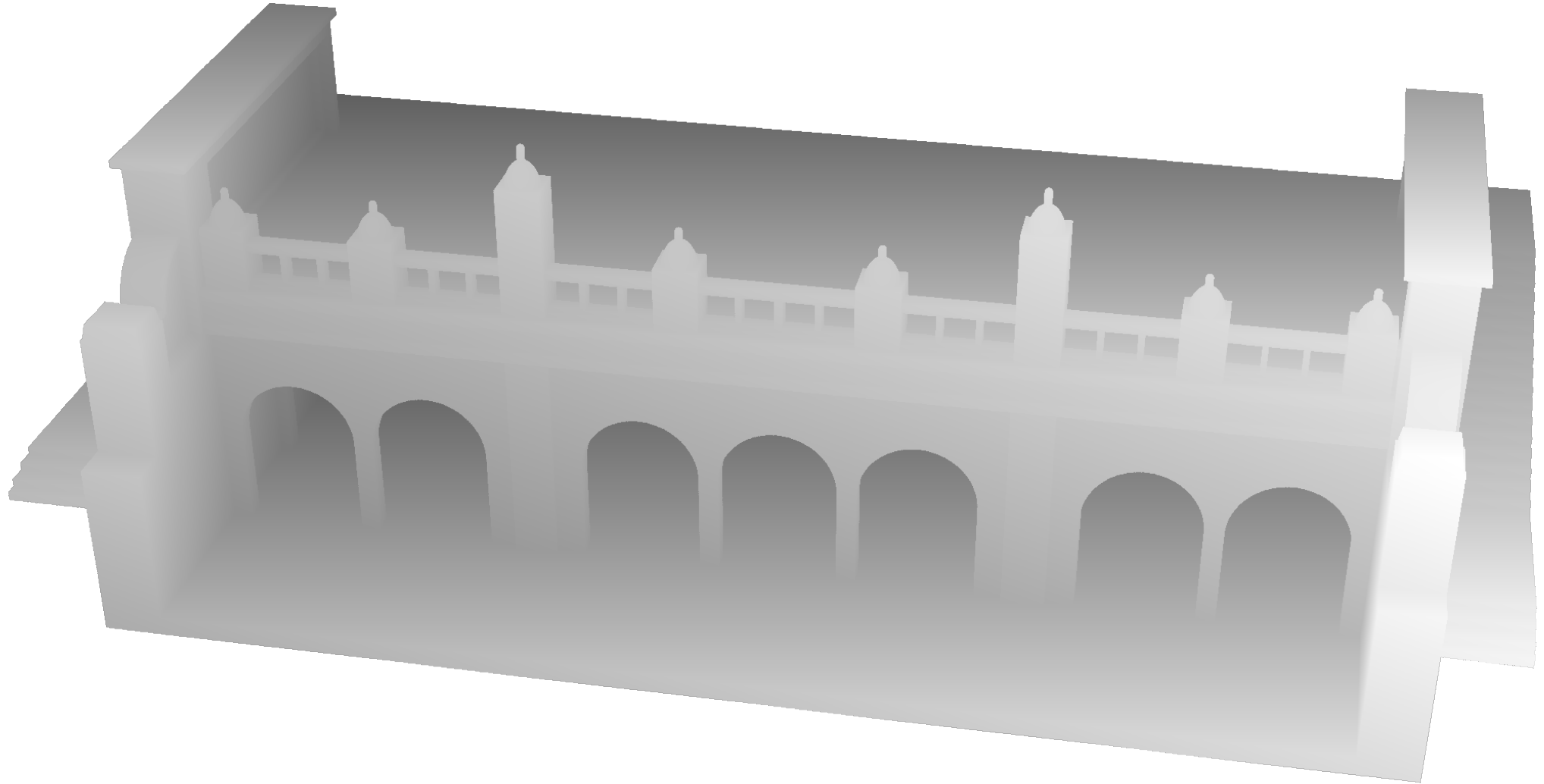
opcode	lhs	rhs	out
X	–	–	slot 0
Y	–	–	slot 1
SQUARE	slot 0	–	slot 0
SQUARE	slot 1	–	slot 1
ADD	slot 0	slot 1	slot 1
SQRT	slot 1	–	slot 1
SUB	slot 1	1.0f	slot 0
SUB	0.5f	slot 1	slot 1
MAX	slot 0	slot 1	slot 1

- XYZ opcodes return coordinates that are transformed by a 4x4 matrix
- OpenGL-style homogenous coordinates (allows for perspective)
- Also used to implementing the camera
- Faster than transforming the model and uploading a new tape

# USER-PROVIDED MATRIX

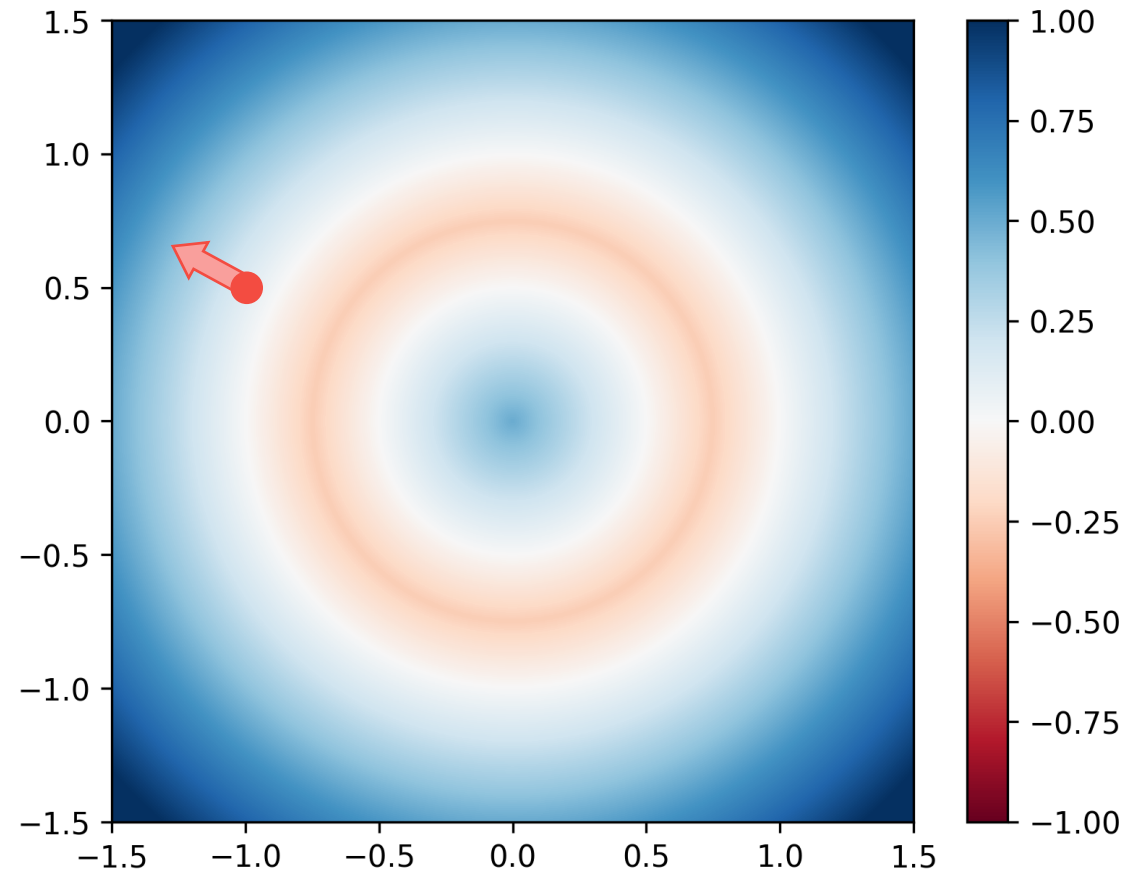


# HEIGHTMAP OUTPUT



# CALCULATING NORMALS

- The **partial derivatives**  
 $df/dx$ ,  $df/dy$ ,  $df/dz$   
approximate the normal at points  
close to the model surface
- **Forward-mode automatic differentiation** can be implemented  
by interpreting the tape with  
**value + partial derivatives**



# 2D AUTOMATIC DIFFERENTIATION

opcode	lhs	rhs	out
X	—	—	$-1.0, (1, 0)$
Y	—	—	$0.5, (0, 1)$
SQUARE		—	
SQUARE		—	
ADD			
SQRT		—	
SUB		$1.0f$	
SUB	$0.5f$		
MAX			



# 2D AUTOMATIC DIFFERENTIATION

opcode	lhs	rhs	out
X	—	—	$-1.0, (1, 0)$
Y	—	—	$0.5, (0, 1)$
SQUARE	$-1.0, (1, 0)$	—	
SQUARE		—	
ADD			
SQRT		—	
SUB		$1.0f$	
SUB	$0.5f$		
MAX			

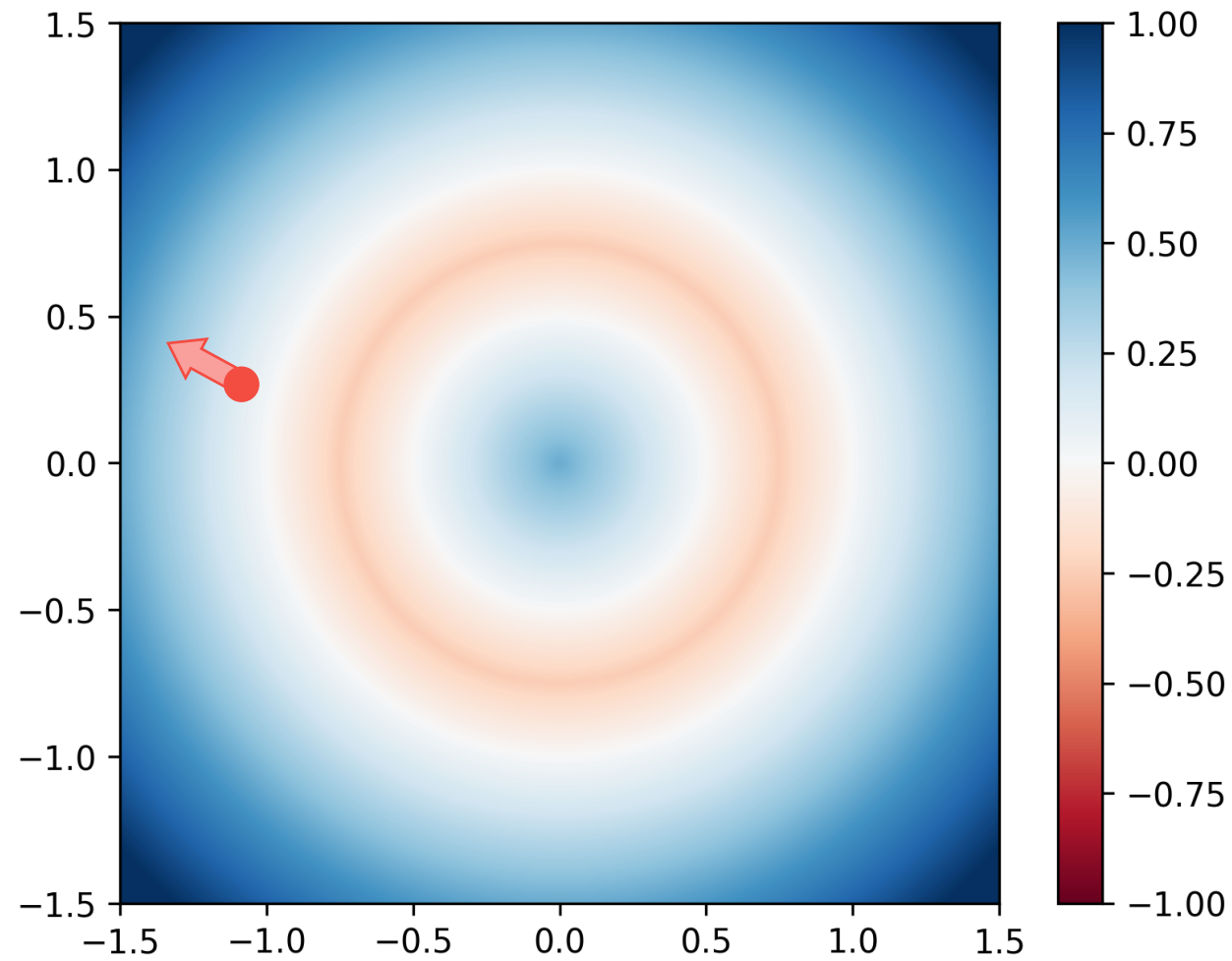
# 2D AUTOMATIC DIFFERENTIATION

opcode	lhs	rhs	out
X	—	—	-1.0, (1, 0)
Y	—	—	0.5, (0, 1)
SQUARE	-1.0, (1, 0)	—	1.0, (-2, 0)
SQUARE		—	
ADD			
SQRT		—	
SUB		1.0f	
SUB	0.5f		
MAX			

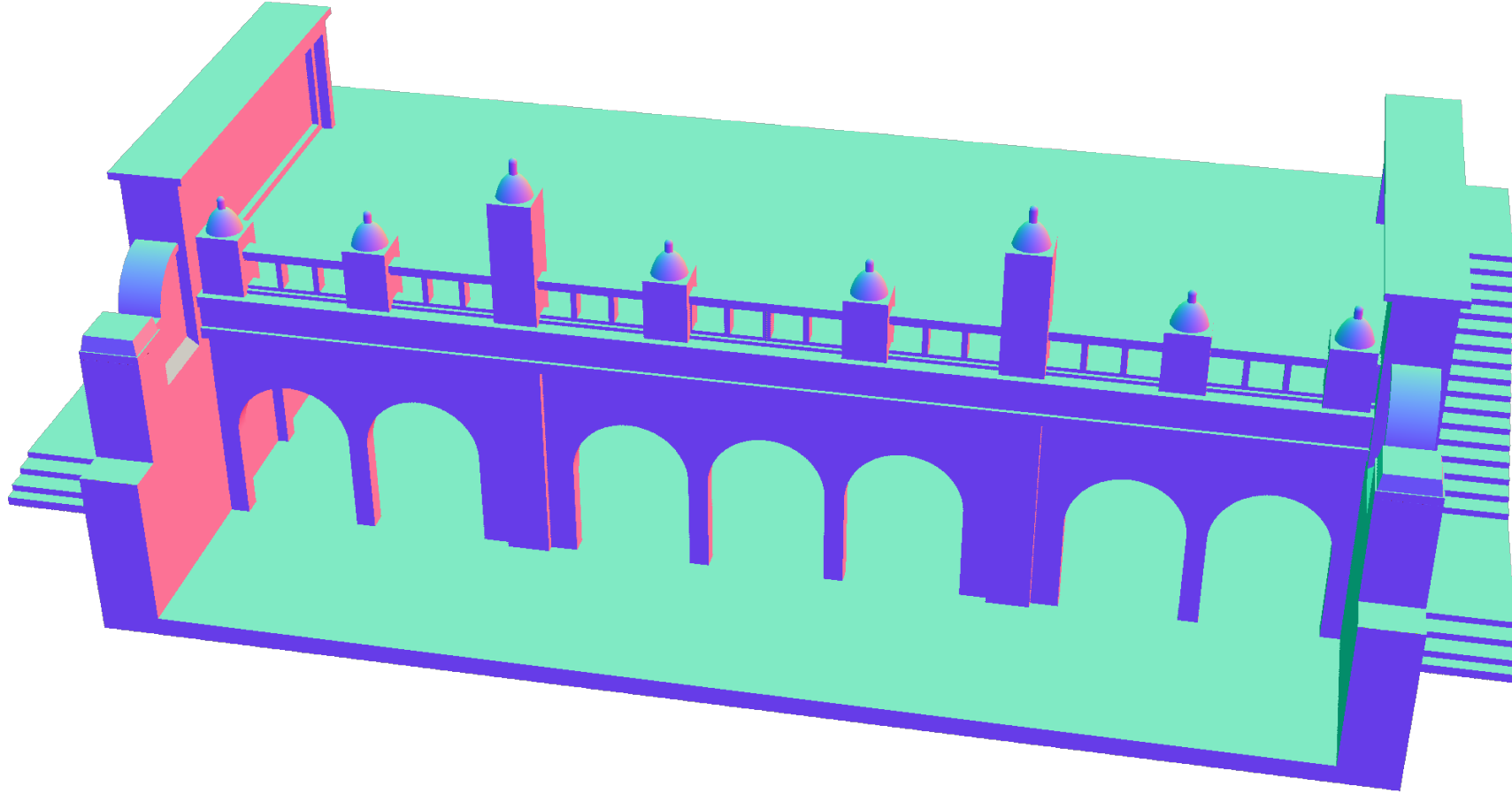
# 2D AUTOMATIC DIFFERENTIATION

opcode	lhs	rhs	out
X	—	—	-1.0, (1, 0)
Y	—	—	0.5, (0, 1)
SQUARE	-1.0, (1, 0)	—	1.0, (-2, 0)
SQUARE	0.5, (0, 1)	—	0.25, (0, 1)
ADD	1.0, (-2, 0)	0.25, (0, 1)	1.25, (-2, 1)
SQRT	1.25, (-2, 1)	—	1.12, (-0.89, 0.45)
SUB	1.12, (-0.89, 0.45)	1.0f	0.12, (-0.89, 0.45)
SUB	0.5f	1.12, (-0.89, 0.45)	-0.62, (0.89, -0.45)
MAX	0.12, (-0.89, 0.45)	-0.62, (0.89, -0.45)	0.12, (-0.89, 0.45)

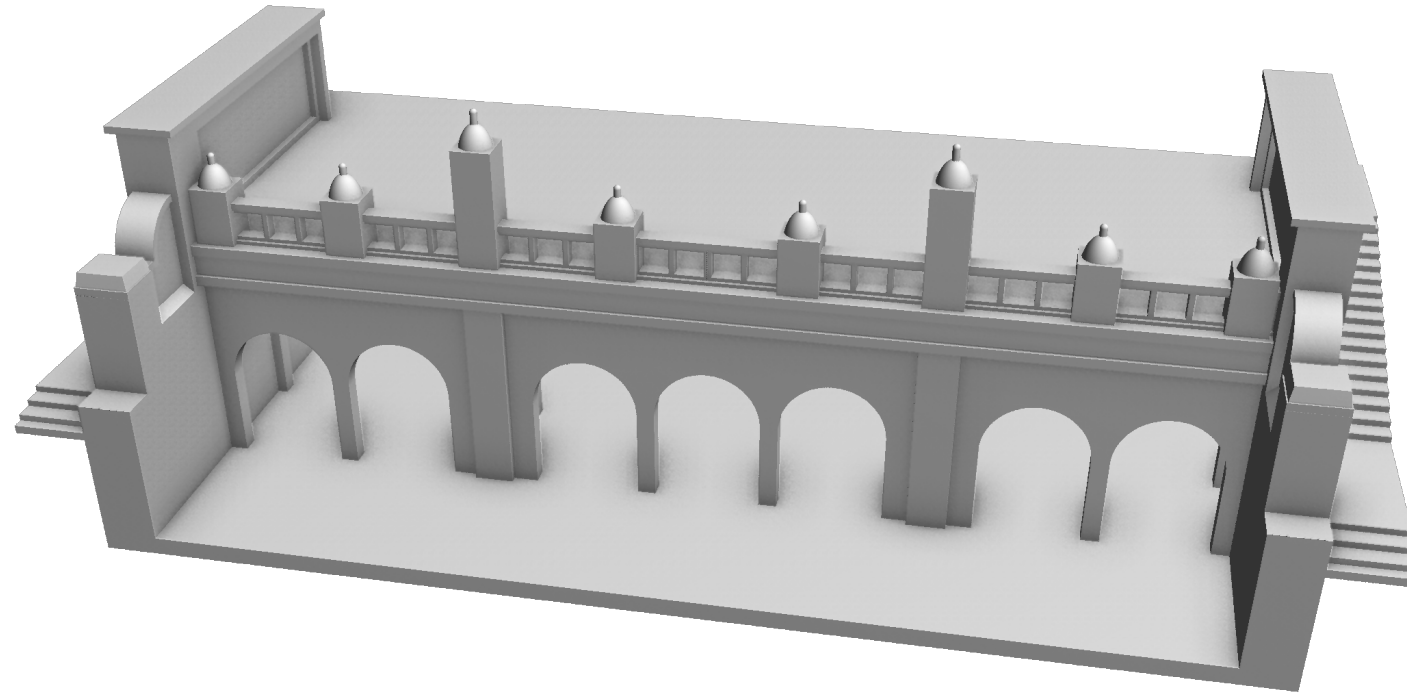
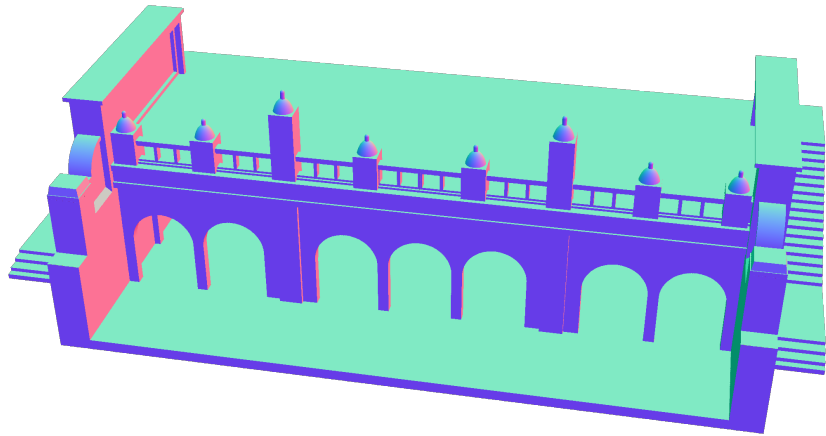
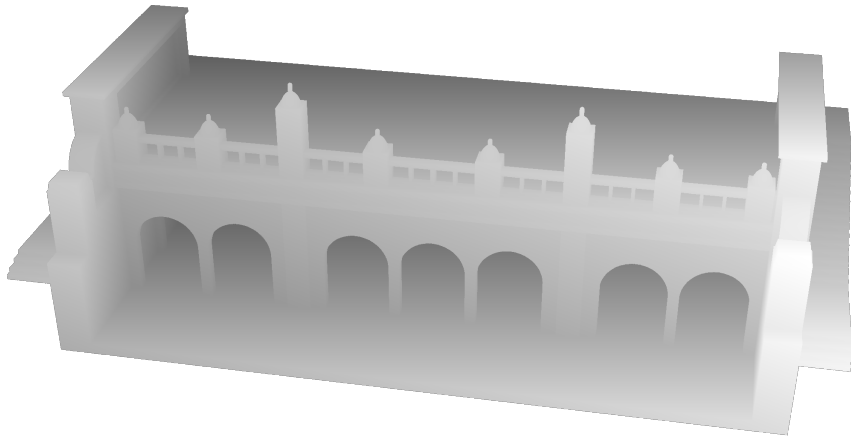
# CALCULATING NORMALS



# NORMALS OUTPUT



# POST-PROCESSING



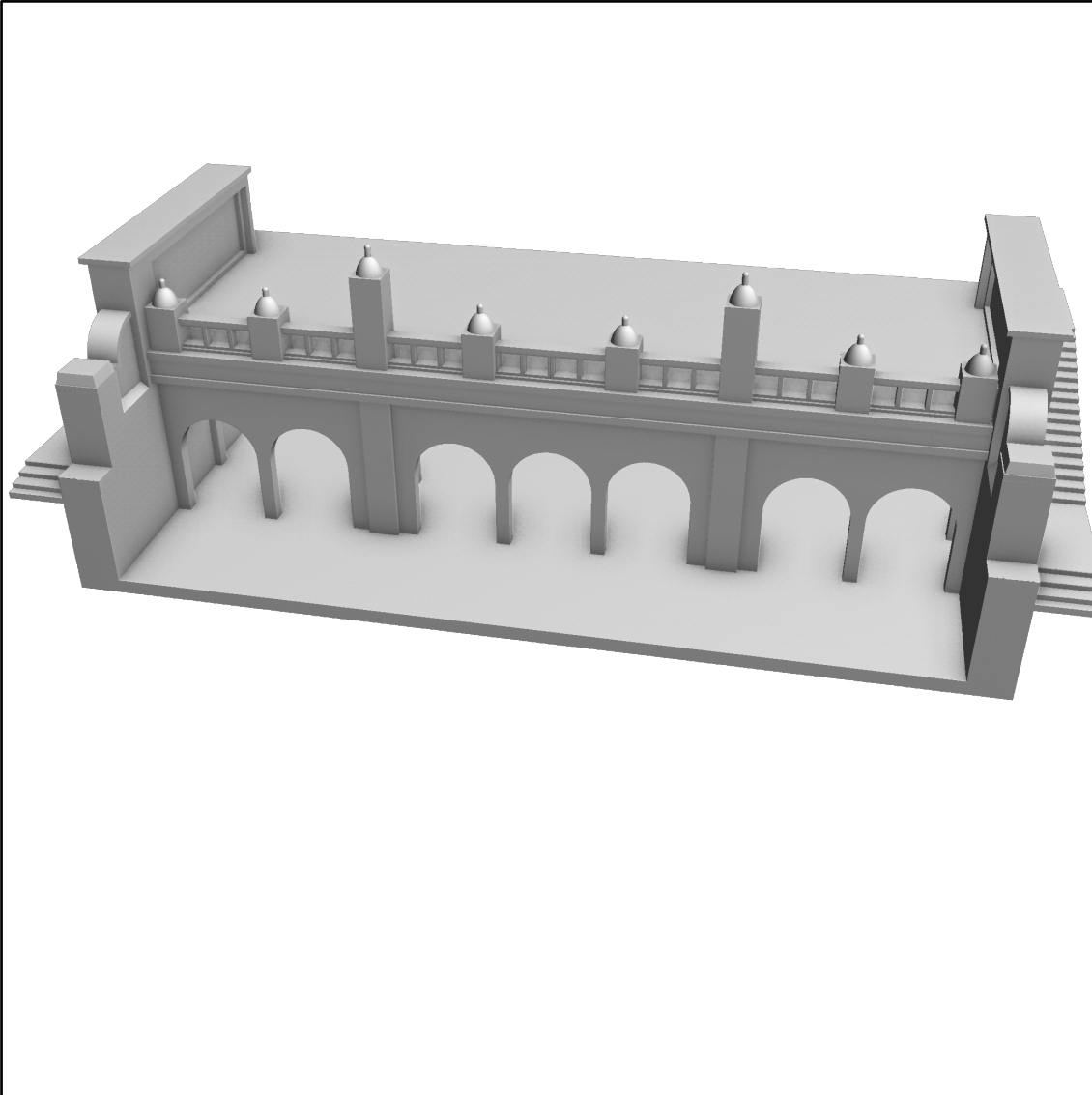
# BENCHMARKING MACHINES

GeForce GT 750M  
2013 MacBook Pro

GTX 1080 Ti  
2017 VR/ML workstation

Tesla V100  
AWS p3.2xlarge

# ARCHITECTURE MODEL



## Frame time (ms)

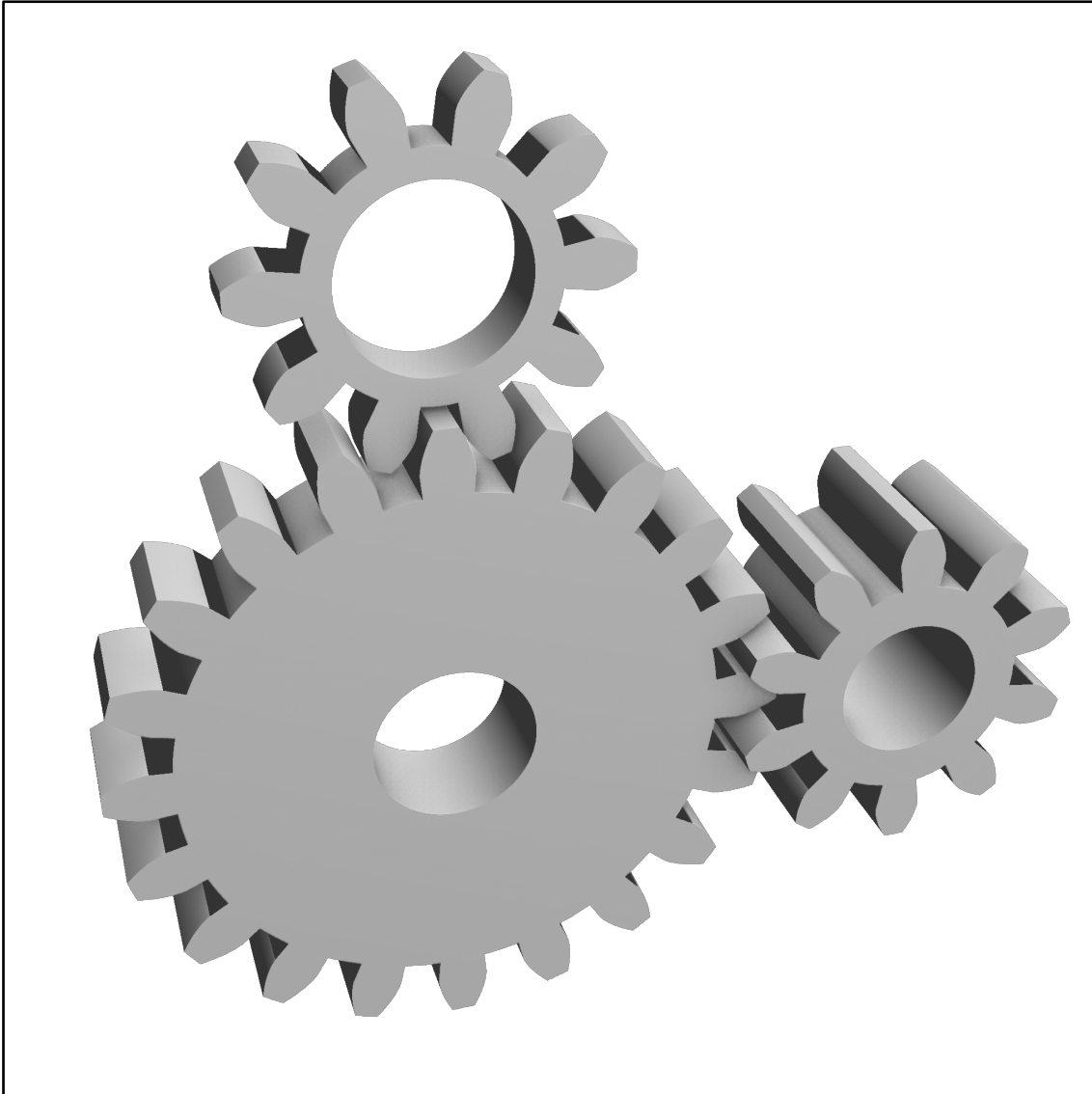
Size	GT 750M	GTX 1080 Ti	Tesla V100
256 <sup>3</sup>	34.3	5.5	3.2
512 <sup>3</sup>	73.9	9.9	5.3
1024 <sup>3</sup>	189.9	22.6	12.2
1536 <sup>3</sup>	331.9	39.3	20.8
2048 <sup>3</sup>	510.7	60.6	31.9

- 961 clauses (465 min/max)
- CSG heavy model
- Best case for our algorithm

Based on a model by  
Jennifer Keeter



# GEARS MODEL



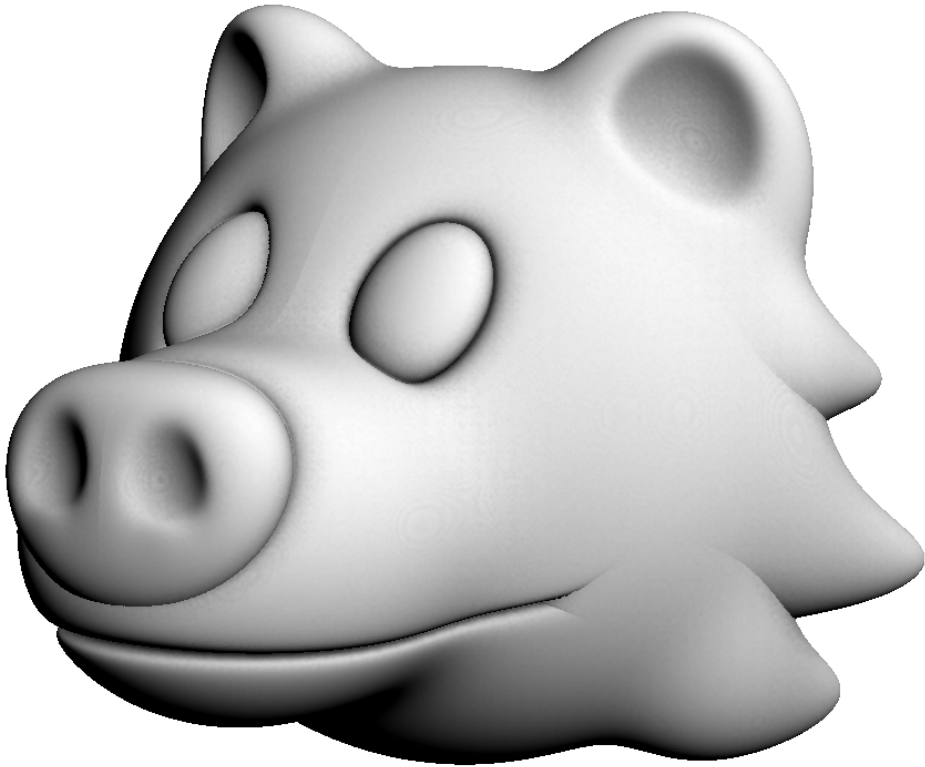
## Frame time (ms)

Size	GT 750M	GTX 1080 Ti	Tesla V100
256 <sup>3</sup>	65.0	9.4	6.2
512 <sup>3</sup>	154.5	16.6	9.2
1024 <sup>3</sup>	426.2	40.3	23.1
1536 <sup>3</sup>	930.3	72.0	39.5
2048 <sup>3</sup>	—	115.4	62.0

- 1735 clauses (374 min/max)
- Medium amount of CSG
- Transcendental functions

Clever involute curve derivation by  
Peter Fedak

# BEAR HEAD MODEL



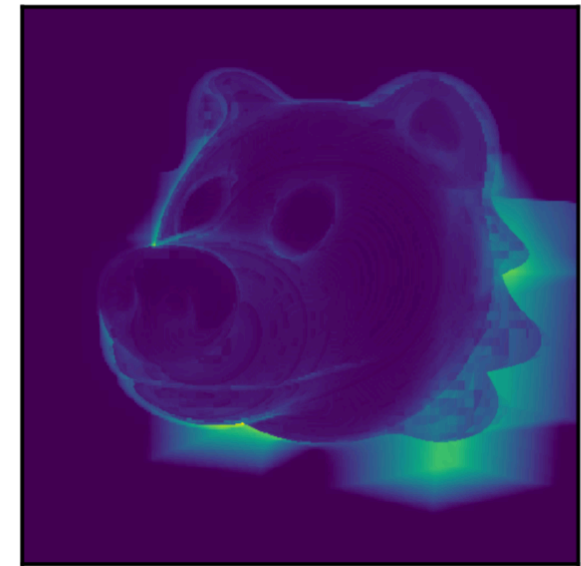
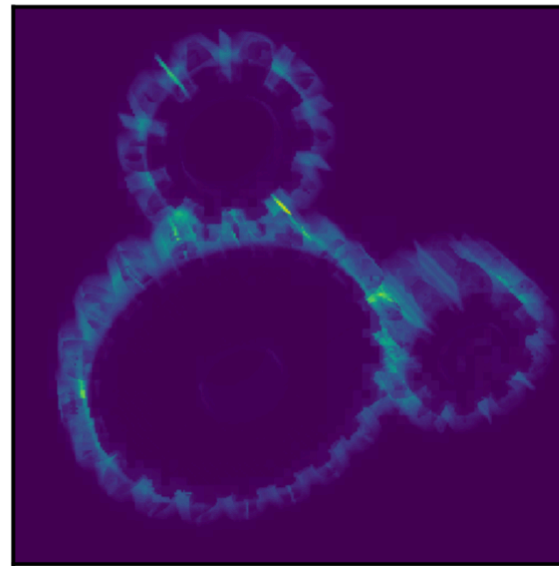
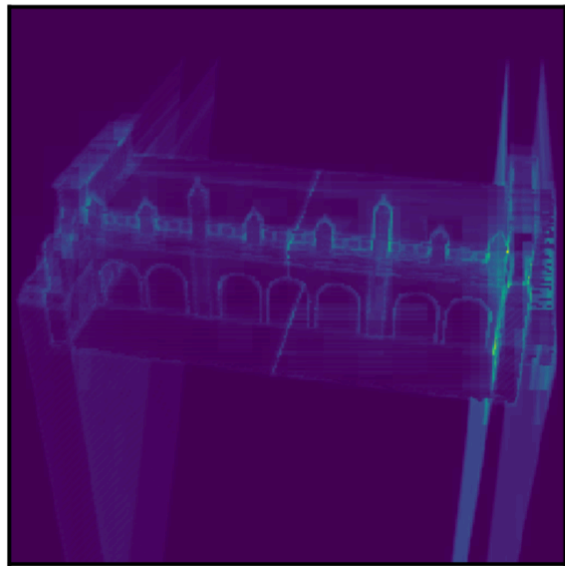
## Frame time (ms)

Size	GT 750M	GTX 1080 Ti	Tesla V100
256 <sup>3</sup>	111.3	11.3	5.2
512 <sup>3</sup>	503.6	41.1	20.3
1024 <sup>3</sup>	2352.1	191.0	88.3
1536 <sup>3</sup>	—	504.2	228.3
2048 <sup>3</sup>	—	1053.2	437.3

- 541 clauses (27 min/max)
- Very little CSG
- Many smooth blends (exp/log)

Based on a design by  
Hazel Fraticelli and Anthony Taconi

# AMORTIZED WORK ( $1024^3$ VOXELS)



**WHERE DO WE GO FROM HERE?**

# SPEED IS LIMITED BY RAM ACCESS

HPC

Feb 11, 2015

## GPU Pro Tip: Fast Dynamic Indexing of Private Arrays in CUDA

By Maxim Milakov

Tags: CUDA, CUDA C/C++, Pro Tip

Discuss (8)

Using local memory is slower than keeping array elements directly in registers, but if you have sufficient math instructions in your kernel and enough threads to hide the latency, the local load/store instructions may be a minor cost. Empirically, a 4:1 to 8:1 ratio of math to memory operations should be enough, the exact number depends on your particular kernel and GPU architecture. Your kernel should also have occupancy high enough to hide local memory access latencies.

# CISC INSTRUCTIONS

- Improve the math-to-computation ratio by adding bigger opcodes
- For example:
  - Every 2D circle computes  $\text{sqrt}((x - x_0)^2 + (y - y_0)^2) - r$
  - `RADIUS_IMM`  $\rightarrow$  `slots[out] = sqrt(slots[lhs]2 + slots[rhs]2) - imm`
  - Before:
    - 5 tape clauses
    - 5 slot reads
    - 5 slot writes
  - After:
    - 1 tape clause
    - 2 slot reads
    - 1 slot write

# VARIABLE INSTRUCTION LENGTH

opcode	out	lhs	rhs	immediate / jump
u8	u8	u8	u8	f32 / i32

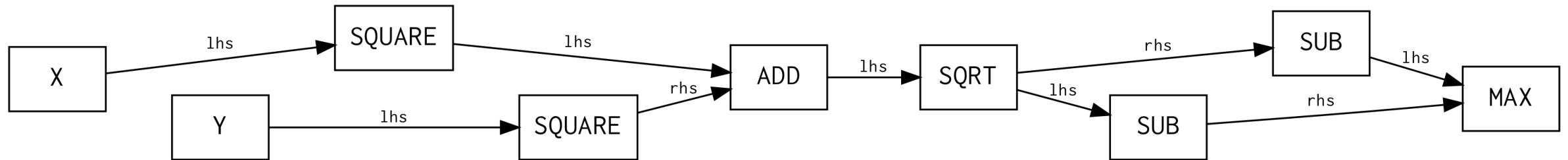


50% of clause size  
Not always used!

- Using variable-length instructions would reduce memory traffic
- This is harder than it sounds!
- Tapes are decoded in both directions
  - Forward during evaluation
  - Backwards during shortening
- Need a strategy for unambiguous encoding

# TAPE SHORTENING WITHOUT INTERPRETER

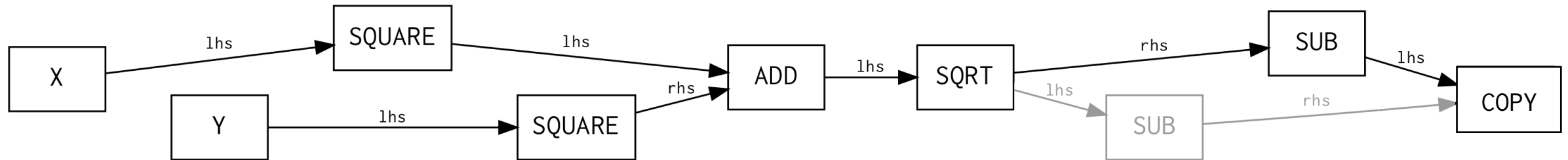
- Instead of building shortened tapes, compile a kernel with sections that are skipped based on flags





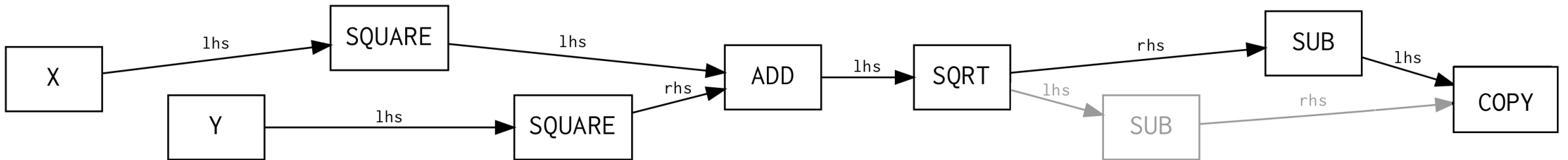
# TAPE SHORTENING WITHOUT INTERPRETER

- Instead of building shortened tapes, compile a kernel with sections that are skipped based on flags



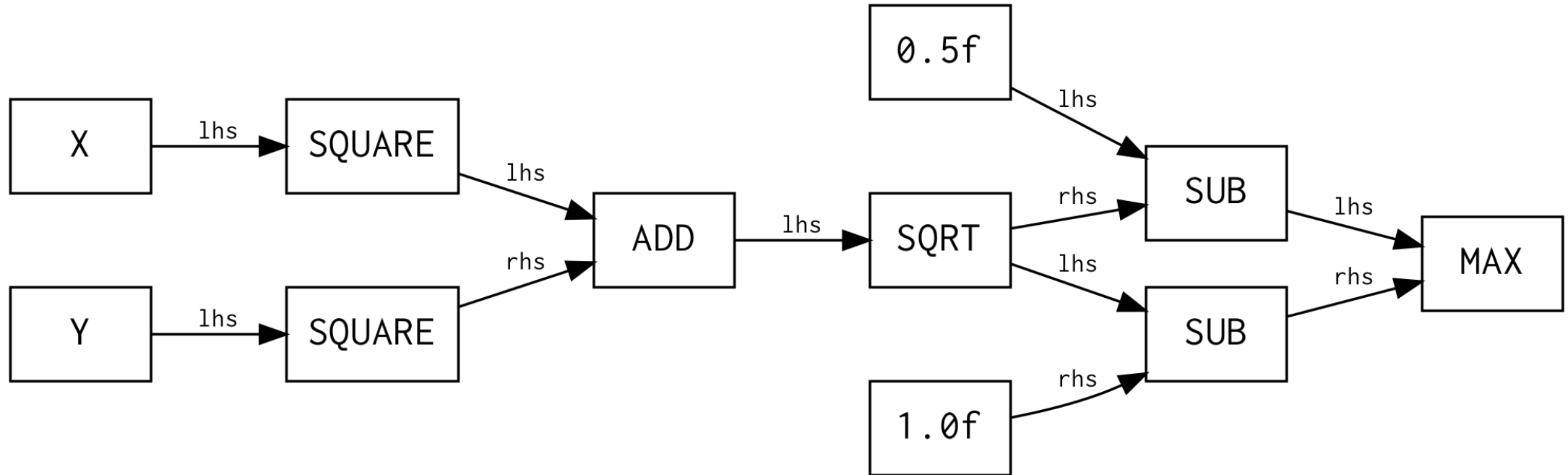
# TAPE SHORTENING WITHOUT INTERPRETER

- Instead of building shortened tapes, compile a kernel with sections that are skipped based on flags

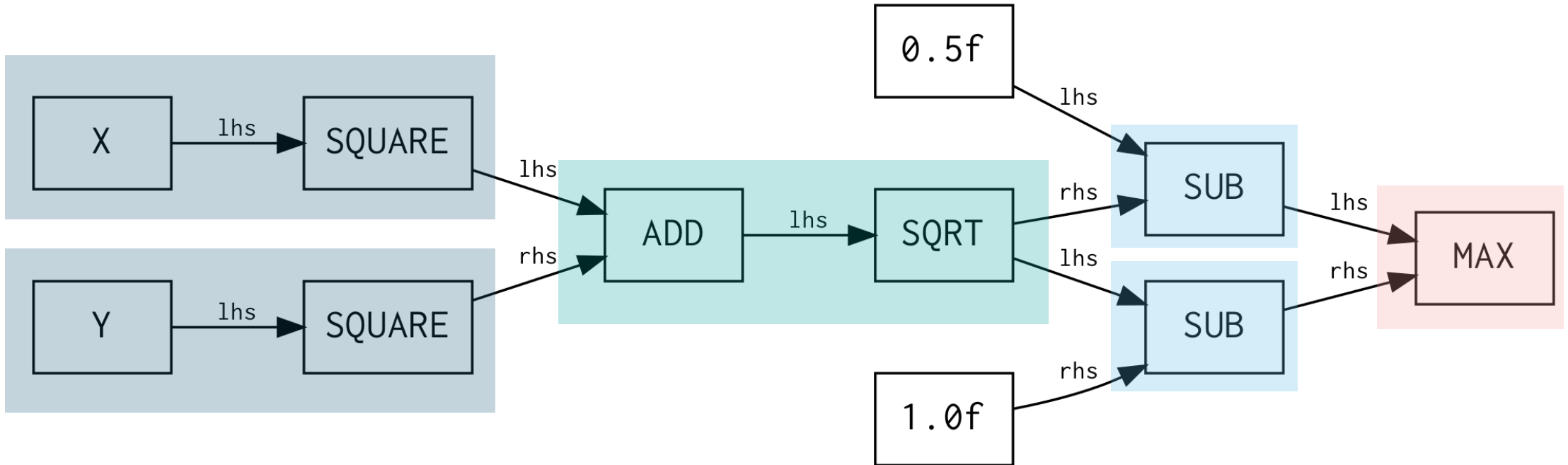


- How do you build the splitting points?
  - $O(n^3)$  possibilities with min/max nodes
  - Interesting graph problem!
  - Draw from compiler literature?

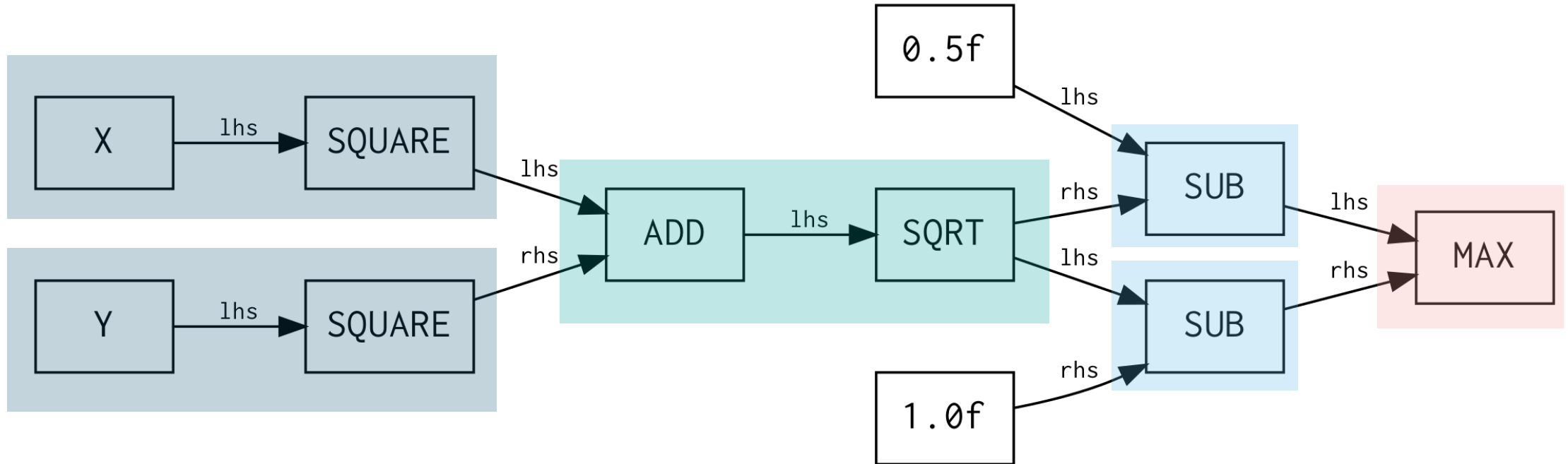
# MORE PARALLELISM



# MORE PARALLELISM



# MORE PARALLELISM



- Also an interesting graph problem!
  - Postdominance frontiers?
  - Synchronization overhead?
  - Effectively scale to **even larger** GPUs


# OTHER POTENTIAL FUTURE WORK

- Use this technique for fast voxelization, then render with other GPU-first algorithms for rendering voxel data
  - Sparse Voxel Octrees (Laine and Karras '10)
  - GVDB (Hoetzlein '16)
- Use reduced affine arithmetic for better intervals (Fryazinov et al. '10)
- Better depth culling of tiles and voxels
- Implementing black-box “oracles” for interoperability with other model formats
  - `libfive` already does this for CPU-side evaluation

# CONCLUSIONS

- New method for rendering complex closed-form implicit surfaces on the GPU
- No triangulation or conventional raytracing!
- Limited to rendering pure mathematical shapes
- Scales well with GPU power
- Works best on hard-surface CSG
- Interpreter speed is gated by global RAM access
  - Being able to generate executable code on the GPU would be great!
- Lots of RAM required for storing shortened tapes

# REFERENCE IMPLEMENTATION

 [mkeeter](#) / [mpr](#)

Unwatch6

Star52

Fork5

<> Code

Issues

Pull requests

Actions

Projects

Wiki

Security

Insights


...

master

Go to file

Add file

Code

 **mkeeter** committed 6519406 on May 8

518 commits2 branches0 tags

benchmark	Update license details and dylib name	2 months ago
gui	Update license details and dylib name	2 months ago
inc	Rename everything (libfive-cuda -> mpr)	3 months ago
<a href="#">libfive @ 8c6638c</a>	Optimize tape construction and bump submodule	6 months ago
src	Update license details and dylib name	2 months ago
.gitignore	Add a 'local' folder for files that shouldn't be shar...	3 months ago
.gitmodules	Use HTTPS for submodule (rather than git)	2 months ago
CMakeLists.txt	Rename everything (libfive-cuda -> mpr)	3 months ago
README.md	Fix links and update preprint status	2 months ago
run_benchmarks.sh	Shuffle lots of files around	3 months ago

README.md

About

Reference implementation for "Massively Parallel Rendering of Complex Closed-Form Implicit Surfaces"

Readme

Releases

No releases published  
[Create a new release](#)

Packages

No packages published  
[Publish your first package](#)

Languages

C++ 82.5% C 11.8% Cuda 5.1% Other 0.6%

[github.com/mkeeter/mpr](https://github.com/mkeeter/mpr)

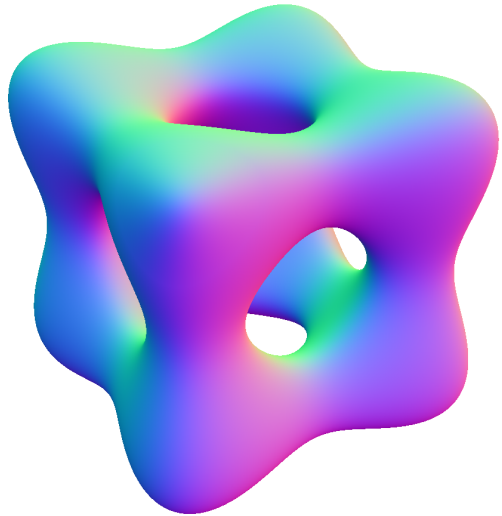
Released under MPL v2  
(weak copyleft)

Reproduce my results on  
AWS for about \$5!



# ACKNOWLEDGEMENTS AND THANKS

- **Beta readers:** Jonathan Bachrach, Blake Courter, Martin Galese, Neil Gershenfeld, Raph Levien, Brian Merchant, Doug Moen, and Amira Abdel Rahman.
- **Benchmarking:** Martin Galese
- **Models:**
  - The bear head is based on a design by Hazel Fraticelli and Anthony Taconi
  - The involute curve math was derived by Peter Fedak
  - The architecture model is based on a design by Jennifer Keeter
- The anonymous reviewers for their feedback and insight
- My colleagues at **Formlabs**, for encouraging this work
- **nTopology**, for their support of the libfive kernel



# THANK YOU!

Matt Keeter  
Independent researcher

 [mattkeeter.com/research/mpr](http://mattkeeter.com/research/mpr)

 [matt.j.keeter@gmail.com](mailto:matt.j.keeter@gmail.com)

 [@impraxical](https://twitter.com/impraxical)