

GOLIGHTLY:  
A GPU IMPLEMENTATION OF THE  
FINITE-DIFFERENCE TIME-DOMAIN METHOD

Approved by:

---

Dr. Marc Christensen

---

Professor Ira Greenberg

---

Dr. Nathan Huntoon

GOLIGHTLY:  
A GPU IMPLEMENTATION OF THE  
FINITE-DIFFERENCE TIME-DOMAIN METHOD

A Thesis Presented to the Graduate Faculty of the

Lyle School of Engineering

Southern Methodist University

in

Partial Fulfillment of the Requirements

for the degree of

Master of Science

with a

Major in Electrical Engineering

by

S. David Lively

(B.S.E.E, Southern Methodist University, 2008)

August 1, 2016

## **ACKNOWLEDGMENTS**

I thank my committee for their patience, insight and unfailing encouragement. Without them, this thesis would remain vaporware. Never give up, never surrender!

Lively , S. David

B.S.E.E, Southern Methodist University, 2008

GoLightly:

A GPU Implementation of the  
Finite-Difference Time-Domain Method

Advisor: Professor Marc Christensen

Master of Science degree conferred August 1, 2016

Thesis completed August 1, 2016

Traditionally, optical circuit design is tested and validated using software which implement numerical modeling techniques such as Beam Propagation, Finite Element Analysis and the Finite-Difference Time-Domain (FDTD) method.

While effective and accurate, FDTD simulations require significant computational power. Existing installations may distribute the computational requirements across large clusters of high-powered servers. This approach entails significant expense in terms of hardware, staffing and software support which may be prohibitive for some research facilities and private-sector engineering firms.

Application of modern programmable GPGPUs to problems in scientific visualization and computation has facilitated dramatically accelerated development cycles for a variety of industry segments including large dataset visualization[3][1], aerospace[11] and electromagnetic wave propagation in the context of optical circuit design. GPU-based supercomputers such as National Labs' Summit[2], co-designed by NVIDIA and IBM, provide dramatically increased compute capability while using less power and fewer computers.

The FDTD algorithm maps well to the massively-multithreaded data-parallel nature of GPUs. This thesis explores a GPU FDTD implementation and details performance gains, limitations of the GPU approach, optimization techniques and potential future enhancements.

## TABLE OF CONTENTS

LIST OF FIGURES .....	vii
LIST OF TABLES .....	viii
CHAPTER	
1. Introduction .....	1
1.1. FDTD Overview .....	2
1.1.1. Wave equation .....	2
1.1.2. Yee Cell .....	2
2. Device Architecture .....	4
2.1. CPU .....	4
2.2. GPU .....	4
2.2.1. SIMD .....	5
2.2.2. FDTD in SIMD .....	5
3. Meep .....	7
3.1. Modeling .....	7
3.2. Performance .....	7
3.3. Limitations .....	8
4. GoLightly .....	9
4.1. Goals .....	9
4.2. Architecture .....	10
4.2.1. Model Processor .....	10
4.2.2. Simulator .....	12
4.2.3. Visualizer .....	19

4.3. Modeling approach .....	22
4.4. Testing and Validation Methodology .....	24
4.4.1. Analytical Result .....	25
4.4.2. Numerical Result .....	27
5. Results.....	33
5.1. Optimization and Enhancements .....	36
6. Conclusions .....	38
6.1. Usability .....	38
6.2. Future Work .....	38
6.2.1. Genetic Algorithms .....	38
6.2.2. Arbitrary Domain Shape and PML Sinks .....	39
6.3. Final Words.....	39
REFERENCES .....	41

## LIST OF FIGURES

Figure	Page
1.1 2D $TM_Z$ Yee Cell .....	3
4.1 2D Whispering Gallery Mode Sensor .....	22
4.2 Arbitrarily-shaped source (Red pixels) .....	24
4.3 Arbitrarily-shaped source after 20 frames .....	24
4.4 Arbitrarily-shaped source after 100 frames .....	24
4.5 $TM_Z$ Test Model with Plane Wave.....	25
4.6 Snell's Law .....	26
4.7 Plane wave with $\epsilon_R = 1$ .....	28
4.8 Time-averaged (RMS) output in free space $\epsilon_R = 1$ .....	28
4.9 Steady state with $\epsilon_R = 9$ .....	29
4.10 Raw monitor output with $\epsilon_R = 9$ .....	29
4.11 RMS output with $\epsilon_R = 9$ .....	30
4.12 Normalized output with $\epsilon_R = 9$ .....	31
5.1 GoLightly: seconds for 5000 frames with the given domain size .....	33
5.2 GoLightly: Completed E and H updates per second .....	34
5.3 GoLightly vs Meep .....	35
5.4 Speedup - Meep Time / GoLightly Time.....	36

## LIST OF TABLES

Table	Page
4.1 Model processor inputs .....	10
4.2 Color component usage.....	23

## Listings

*To Audrey, Wyatt, Walter and Gwendolyn*

## Chapter 1

### Introduction

The FDTD [13] algorithm is the mechanism used by many commercial electro-magnetic wave simulation packages, as well as open source software such as MIT’s Meep[7].

Given the computationally-intensive nature of FDTD, organizations requiring simulation of large domains or complex circuits must provide significant resources. These may take the form of leased server time or utilization of an on-site high-performance cluster, amongst other options.

In this thesis, we explore an implementation of FDTD utilizing graphics processing units (GPUs) via NVIDIA’s CUDA[6] language. Initially designed to perform image generation tasks such as those required by games, cinema and related fields, modern versions are well-suited for general computation work. GPUs are now enjoying wide adoption in fields such as machine learning[8] and artificial intelligence[12], medical research[10], and other areas which require rapid analysis of large datasets.

Modern consumer-grade GPUs offer thousands or tens of thousands of processing units (“cores”), while high-end CPUs typically offer 4-8 cores. While the two are not interchangeable, some algorithms, such as FDTD, require little or no data interdependence, no branching logic (a severe performance impediment on GPUs) and consist of short cycles of simple operations. The power of the GPU lies in performing these simple operations at large scale, with thousands of threads running in parallel.

The following sections detail FDTD. Later sections describe a CPU-based implementation (MIT’s Meep simulator), and our GPU-based GoLightly simulator. We

verify the GPU solution numerically, and compare performance between CPU- and GPU-based implementations. Finally, we consider future applications and enhancements.

## 1.1. FDTD Overview

FDTD expresses Maxwell's equations as a set of discretized time-domain equations[13]. These equations describe each electric field component in terms of its orthogonal, coupled magnetic fields, and each magnetic field component as a function of its coupled electric fields.

### 1.1.1. Wave equation

In a time-domain simulation, the wave equation for  $E_z$  is of the form:

$$\frac{\partial E_z}{\partial t} = K * \left( \frac{\partial H_x}{\partial y} + \frac{\partial H_y}{\partial x} \right) \quad (1.1)$$

Equation 1.1 states that the temporal derivative (change in amplitude) of  $E_z$  is a function of the  $Y$ -axis spatial derivative of the  $H_x$  field and the  $X$ -axis spatial derivative of the  $H_y$  field.

In order to apply this equation to a computational domain, FDTD defines a cell-based discretization strategy.

### 1.1.2. Yee Cell

Yee [13] defines a computational unit known as a "cell." The cell describes how each field component within a domain is related to its coupled fields. For instance, each  $E_Z$  field component depends on adjacent  $H_y$  and  $H_x$  components. The cell format used in such a simulation is of the form shown in Figure 1.1.

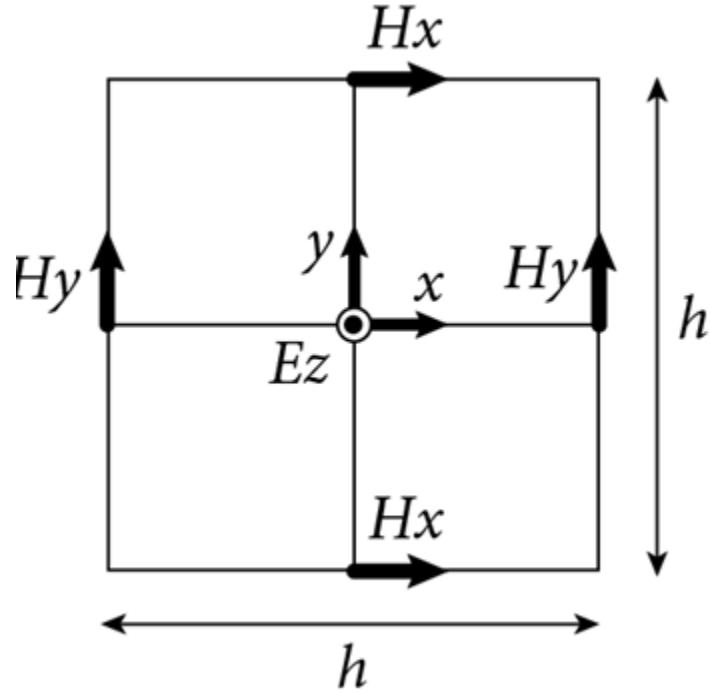


Figure 1.1. 2D  $TM_Z$  Yee Cell

More formally, we may expand the  $E_Z$  wave equation, arriving at:

$$E_{z,i,j}^t = C_a * E_{z,i,j}^{t-1} + C_b * (H_{x,i,j+\frac{1}{2}}^{t-\frac{1}{2}} - H_{x,i,j-\frac{1}{2}}^{t-\frac{1}{2}}) + C_b * (H_{y,i+\frac{1}{2},j}^{t-\frac{1}{2}} - H_{xy,i-\frac{1}{2},j}^{t-\frac{1}{2}}) \quad (1.2)$$

Similarly, the equations for the coupled fields  $H_x$  and  $H_y$  may be expressed as:

$$H_{x,i,j}^t = D_a * H_{x,i,j}^{t-1} + D_b * (E_{z,i,j+\frac{1}{2}}^{t-\frac{1}{2}} - E_{z,i,j-\frac{1}{2}}^{t-\frac{1}{2}}) \quad (1.3)$$

$$H_{y,i,j}^t = D_a * H_{y,i,j}^{t-1} + D_b * (E_{z,i+\frac{1}{2},j}^{t-\frac{1}{2}} - E_{z,i-\frac{1}{2},j}^{t-\frac{1}{2}}) \quad (1.4)$$

## Chapter 2

### Device Architecture

CPUs and GPUs each offer advantages for different computational tasks. Multi-core CPUs offer independent cores which are effectively discrete processors. GPUs, however, offer large-scale parallelization, but require strong data and code coherence in order to achieve acceptable performance.

#### 2.1. CPU

Users typically run many different applications in parallel: a web browser, music player, word processor and email client are a common combination.

In a modern multi-core CPU, each core provides a dedicated ALU and register set. This allows each core to operate as an independent device. This architecture is advantageous when the device is required to perform disparate operations.

However, this approach also introduces some performance limitations. If we wish to perform identical operations on large datasets, we are limited to 4-8 threads. The same processor that excels at executing many different tasks at the same time performs sub-optimally.

When running FDTD on a CPU, each core's ALU executes essentially the same operations on a dedicated register set. Each ALU performs the same operation at the same time, indicating that the additional ALUs are redundant. Thus, the flexible, general-purpose nature of a CPU becomes a liability.

#### 2.2. GPU

Unlike CPUs, GPUs are capable of running thousands of threads in parallel using a SIMD (single-instruction, multiple-data) model. This architecture allows rapid processing of large datasets wherein each datum exhibits little or no interdependence.

This independence is crucial to performant GPU-based applications, as detailed in subsection 2.2.1.

### 2.2.1. SIMD

In a SIMD[4] architecture, a core may consist of a single ALU, with multiple register banks. Separate "threads" load different data into dedicated register banks. The ALU executes identical operations on all registers simultaneously. In this way, each register bank is analogous to a CPU "thread." However, since the same instruction must be executed on all register sets at every step, this design is not as flexible as a system where each core is truly independent.

That said, the approach provides some benefits and limitations:

- Fewer components are required per core since fewer ALUs are required. Leads to reduced die space requirements.
- Better code caching. A single ALU, and corresponding cache, are used for many threads, eliminating the need to load or monitor code cache behavior per thread.
- Smaller footprint. Use of fewer components allows more cores to be placed within a given physical space.

### 2.2.2. FDTD in SIMD

FDTD's leap-frog update method, whereby E fields and H fields are successively calculated, is well-suited to a GPU implementation. E field values depend on adjacent H field values, and visa-versa. Since the E-field update equation requires knowledge

only of the H field state and previous E field state, each field component can be calculated independently with no opportunity for a pipeline stall or race condition.

## Chapter 3

### Meep

Meep[7] is a full-featured, open-source simulator produced by the Massachusetts Institute of Technology. In addition to its core FDTD-based simulation engine, it provides a scripting interface for defining models and simulation parameters, recording results, and other tasks.

#### 3.1. Modeling

One limitation of Meep is its use of an obscure, uncommon scripting language known as SCHEME. Models are defined in terms of constructive solid geometry (CSG) commands, whereby the user describes their model in terms of boolean operations and regular polyhedra.

While adequate for simple models, constructing an arbitrarily-shaped or dynamic structure in this way may be difficult. In practice, proprietary software may be used to convert more complex model definitions into a format usable by Meep.

It is worth noting that Meep provides a "material function" capability. This allows the user to dynamically determine the material properties of a point in space using their own algorithm rather than defining their model using CSG. However, to take advantage of this capability, the user must employ additional software or custom programming. This further increases the complexity of an already-complicated system.

#### 3.2. Performance

Meep is a mature, highly-optimized suite of tools. While complex to configure and use, it scales well across multiple machines, relying on the MPI protocol to keep nodes within a simulation cluster in sync.

While performant when compared to other FDTD software, Meep suffers from the same architecture-imposed limitations of all CPU-based implementations. The limited number of processing cores available on a general-purpose CPU restricts the number of data points that can be processed within a given time frame. This problem can be solved by provisioning additional computers which would run in parallel, distributing the computational load evenly across the resulting cluster.

This sort of cluster configuration incurs its own overhead. Although a domain may be divided into chunks and distributed across cluster nodes, FDTD boundary conditions require that, at some point, parts of the divided domain must be exchanged between nodes to maintain continuity. This necessitates installation of a high-speed local network and supporting hardware.

### 3.3. Limitations

The manner in which models are defined, poor performance on single machines and laptops, and incompatibility with Microsoft Windows all indicate that this system is not suitable for rapid design iteration. To that end, we present an alternative, GPU-based approach. (See chapter 4)

## Chapter 4

### GoLightly

GoLightly is the GPU-based FDTD simulator application that is the focus and product of this thesis. Written using a combination of C++, CUDA and OpenGL, it provides a lightweight yet complete FDTD solution.

#### 4.1. Goals

GoLightly is intended to address deficiencies common to CPU-based solutions. In particular, it is designed to be fast, friendly and portable.

- **Fast.** An iterative design process requires rapid feedback from the simulator. Long simulation times necessitated by existing solutions inhibit this process.
- **Friendly.** Definition of models and other simulation parameters should not require expertise in software development or quasi-proprietary scripting languages.
- **Portable.** Ideally, the simulator should run on a high-end consumer grade laptop and support the most common desktop operating systems (Microsoft Windows and Apple OS X).

To meet those goals, GoLightly takes advantage of the oft-underutilized programmable GPU available in common desktop and laptop computers, resulting in a dramatic speedup. Rather than relying on a proprietary model definition language or obscure, limited scripting system, we use industry-standard image and geometry file formats so that models may be defined using robust, familiar, readily-available

tools. By building the software specifically for Microsoft Windows, we ensure that it is compatible with the most common desktop operating system.

## 4.2. Architecture

GoLightly comprises three primary application blocks:

- Model Processor 4.2.1
- Simulator 4.2.2
- Visualizer 4.2.3

### 4.2.1. Model Processor

The model processor (MP) is responsible for initialization of the simulator. When launching the simulator, a domain size and image file, containing a coded image of the desired dielectric, as well as a max  $\epsilon$  are specified.

Table 4.1. Model processor inputs

Symbol	Data Type	Meaning	Typical value
Width	int	Domain size in X	1024
Height	int	Domain size in Y	1024
Media	float	$\epsilon_{max}$	9
Model	string	Model definition stored as a bitmap filename	

The MP allocates arrays to hold the dielectric properties for each Yee cell. These arrays are of the same dimensions as the domain, which may be different than the dimensions of the model.

Once the model image is loaded, the MP iterates through each element in the dielectric array. (See 4.2.1)

For each element:

1. Determine the normalized texel coordinate that corresponds to the current cell position
2. Read the red (R), green (G) and blue (B) color components from the image
3. If  $R > 128$ , this texel is part of a source. Add the cell to the list of sources
4. If  $G > 0$ , this texel has non-unity dielectric. Set  $C_{bi,j} = \epsilon_{max} * \frac{G}{255.0}$
5. If  $B > 0$ , this texel is part of a monitor. Add its position to the monitor definition with  $ID = B$

```
1 for( int j = 0; j < media.Size.y; j++)
2 {
3     int sourceY = j * height / media.Size.y;
4     for( int i = 0; i < media.Size.x; i++)
5     {
6         int sourceX = i * width / media.Size.x;
7         unsigned int sourceOffset = channels * (sourceY * width + sourceX);
8         unsigned int mediaOffset = j * media.Size.x + i;
9         unsigned char red = bytes[sourceOffset + 0];
10        unsigned char green = bytes[sourceOffset + 1];
11        unsigned char blue = bytes[sourceOffset + 2];
12        // is this pixel part of a source?
13        if (red > 128)
14            sourceOffsets.push_back(mediaOffset);
15
16        // fill default waveguide material (parameter n)
17        if (green > 0)
18        {
19            // interpolate n based on green value.
20            media.HostArray[mediaOffset] = epsilonMax * green * 1.f / 255;
```

```

21     }
22 }
23 }
```

Listing 4.1. Generating a model from an image

Once the dielectric, sources and monitors are derived from the model image, the model processor transfers control to the simulator.

#### 4.2.2. Simulator

The simulator block implements the FDTD algorithm. Given the dielectric, source and monitor configurations from the model processor, the simulator initializes the GPU, transfers required data from host memory to the GPU, and begins the simulation loop.

In addition to the dielectric and field arrays, the simulator generates a descriptor (4.2.2) for each field that will be updated. This structure is used by the kernels to assist in handling boundary conditions (PML) and other housekeeping duties. A similar, more compact descriptor (4.2.2) is generated from the host descriptor and passed to the kernels.

```

1 struct FieldDescriptor
2 {
3     /// <summary>
4     /// describes a split-field boundary region for PML
5     /// </summary>
6     struct BoundaryDescriptor
7     {
8         FieldType Name;
9         FieldDirection Direction;
10
11     /// CPU-resident fields
```

```

12     float *Amp;
13     float *Psi ;
14     float *Decay;
15
16     BoundaryDescriptor *DeviceDescriptor ;
17
18     unsigned int AmpDecayLength;
19
20     private:
21     CudaHelper *m_cuda;
22 };
23
24
25     float DefaultValue;
26     FieldType Name;
27
28     dim3 Size;
29     dim3 UpdateRangeStart;
30     dim3 UpdateRangeEnd;
31
32     vector<float> HostArray;
33     float *DeviceArray;
34
35     DeviceFieldDescriptor *DeviceDescriptor;
36
37     vector<GridBlock> GridBlocks;
38     map<FieldType , shared_ptr<BoundaryDescriptor>> Boundaries;
39 };

```

Listing 4.2. Host Field Descriptor structure

```

1 enum class FieldDirection { X,Y,Z };

```

```

2
3 struct DeviceFieldDescriptor
4 {
5     FieldType Name;
6     dim3 Size;
7     dim3 UpdateRangeStart;
8     dim3 UpdateRangeEnd;
9
10    float *Data;
11};

```

Listing 4.3. Device Field Descriptor

For each loop iteration, the simulator launches a CUDA kernel to update all  $E$  fields. Once the  $E$  update is complete, the simulator launches kernels to update all  $H$  fields.

The three kernels required for a  $TM_Z$  simulation are detailed below:

```

1 --global-- void UpdateEz(
2     dim3 threadOffset
3 )
4 {
5     unsigned int x = threadOffset.x + blockIdx.x * blockDim.x + threadIdx.
6         x;
7     unsigned int y = threadOffset.y + blockIdx.y * blockDim.y + threadIdx.
8         y;
9
10    if (y < 1 || x < 1)
11        return;
12
13    float cb = Cb->Data[y * Cb->Size.x + x];
14
15}

```

```

13 unsigned int center = y * Ez->Size.x + x;
14 float hxBottom = Hx->Data[y * Hx->Size.x + x];
15 float hxDTop = Hx->Data[(y - 1) * Hx->Size.x + x];
16 float dhx = (hxBottom - hxDTop);
17
18 float hyRight = Hy->Data[y * Hy->Size.x + x];
19 float hyLeft = Hy->Data[y * Hy->Size.x + x - 1];
20 float dhy = (hyLeft - hyRight);
21
22 float ezxPsi = 0.f;
23 float ezyPsi = 0.f;
24
25 // PML
26 if (x < 10 || x > Ez->UpdateRangeEnd.x - 10 || y < 10 || y > Ez->
27 UpdateRangeEnd.y - 10)
28 {
29     ezyPsi = Ezy->Decay[y] * Ezy->Psi[center] + Ezy->Amp[y] * dhx;
30     Ezy->Psi[center] = ezyPsi;
31     ezxPsi = Ezx->Decay[x] * Ezx->Psi[center] + Ezx->Amp[x] * dhy;
32     Ezx->Psi[center] = ezxPsi;
33 }
34
35 Ez->Data[center] = CA * Ez->Data[center] + cb * (dhy - dhx) + cb * (
36     ezxPsi - ezyPsi);
}

```

Listing 4.4. CUDA kernel for updating  $E_Z$

The majority of each kernel's source performs setup and bounds checking tasks. In each kernel, the FDTD equation implementation can be isolated to one or two lines of code.

For example, the line (from the  $E_Z$  update kernel),

```
1 Ez->Data[center] = CA * Ez->Data[center] + cb * (dhy - dhw) + cb * (
2   ezxPsi - ezyPsi);
```

corresponds to the FDTD  $E_Z$  equation. (See Equation 1.2)

```
1 __global__ void UpdateHx(dim3 threadOffset)
2 {
3   unsigned int x = threadOffset.x + blockIdx.x * blockDim.x + threadIdx.x;
4   unsigned int y = threadOffset.y + blockIdx.y * blockDim.y + threadIdx.y;
5
6   if (y >= Ez->Size.y - 1)
7     return;
8
9   unsigned int hxOffset = y * Hx->Size.x + x;
10 #ifdef USE_MAGNETIC_MATERIALS
11   float db = Db->Data[y * Db->Size.x + x];
12 #else
13   const float db = DbDefault;
14 #endif
15   //float ezTop = Ez->Data[y * Ez->Size.x + x];
16   //float ezBottom = Ez->Data[(y+1) * Ez->Size.x + x];
17
18   float dEz = Ez->Data[(y + 1) * Ez->Size.x + x] - Ez->Data[y * Ez->Size.x + x];
19
20   float hx = DA * Hx->Data[hxOffset] - db * dEz;
21
22   if (y < 10 || y > Hx->UpdateRangeEnd.y - 10 || x < 10 || x > Hx->
23     UpdateRangeEnd.x - 10)
24   {
```

```

24 // update boundaries
25 float decay = Hxy->Decay[y];
26 float amp = Hxy->Amp[y];
27 float psi = Hxy->Psi[hxOffset];
28
29 psi = decay * psi + amp * dEz / Configuration->Dx;
30
31 Hxy->Psi[hxOffset] = psi;
32 hx = hx - db * Configuration->Dx * psi;
33 }
34
35 Hx->Data[hxOffset] = hx;
36 }
37
38 -

```

Listing 4.5. CUDA kernel for updating  $H_X$

```

1 __global__ void UpdateHy(dim3 threadOffset)
2 {
3     unsigned int x = threadOffset.x + blockIdx.x * blockDim.x + threadIdx.x;
4     unsigned int y = threadOffset.y + blockIdx.y * blockDim.y + threadIdx.y;
5
6     if (x >= Ez->Size.x - 1)
7         return;
8
9     unsigned int hyOffset = y * Hy->Size.x + x;
10
11 #ifdef USE_MAGNETIC_MATERIALS
12     float db = Db->Data[y * Db->Size.x + x];

```

```

13 #else
14     const float db = DbDefault;
15 #endif
16
17     float ezLeft = Ez->Data[y * Ez->Size.x + x];
18     float ezRight = Ez->Data[y * Ez->Size.x + x + 1];
19
20     float dEz = ezRight - ezLeft;
21     float hy = DA * Hy->Data[hyOffset] - db * (ezRight - ezLeft);
22
23     if (x < 10 || y < 10 || x > Hy->UpdateRangeEnd.x - 10 || y > Hy->
24         UpdateRangeEnd.y - 10)
25     {
26         float psi = Hyx->Psi[hyOffset];
27         float decay = Hyx->Decay[x];
28         float amp = Hyx->Amp[x];
29
30         psi = decay * psi + amp * dEz / Configuration->Dx;
31
32         hy = hy - db * Configuration->Dx * psi;
33
34         Hyx->Psi[hyOffset] = psi;
35     }
36
37     Hy->Data[hyOffset] = hy;
38 }
```

Listing 4.6. CUDA kernel for updating  $H_Y$

Note that all  $E$  updates occur simultaneously, as do all  $H$  fields. However, given the dependence between the  $E$  and  $H$  fields, the  $E$  field update kernels must complete

before the  $H$  fields are updated.

The simulator repeats this operation until the application is closed, or the desired number of frames are completed.

Finally, the completed field arrays are copied to the host from the GPU, and saved to disk in bitmap and CSV format for later analysis.

#### 4.2.3. Visualizer

If enabled<sup>1</sup>, the visualizer application block provides interactive display of the simulation.

When running a simulation, the user may optionally specify a visualizer update frequency, indicating the number of simulation frames that should complete between visualizer updates. This aids in reducing the visualizer's performance impact.

A window and OpenGL context are created using GLFW and the glLoadGen OpenGL extension loader. An OpenGL pixel buffer object (PBO) is allocated to contain a copy<sup>2</sup> of the field that the user wishes to see<sup>3</sup>.

The visualizer also creates a screen-aligned quad on which the field texture will be rendered, and allocates an texture object which is then bound to the PBO.

After the required number of frames have been completed, the visualizer launches a CUDA kernel which samples the selected field and populates the PBO.

```
1 --global__ void visualizerUpdatePreviewTexture(
2     cudaSurfaceObject_t image
3     , int imageWidth
4     , int imageHeight
5     , float *fieldData
```

---

<sup>1</sup>The visualizer requires some GPU overhead. As such, its use may affect simulator performance.

<sup>2</sup>The PBO may be of different dimensions than the simulation domain. Since the PBO is used only for visualization, it is not necessary for it to contain the full-resolution field.

<sup>3</sup>The visualizer provides the ability to dynamically select which field(s) should be displayed

```

6     , int fieldWidth
7     , int fieldHeight
8     , float *materials
9 )
10 {
11     unsigned int x = blockIdx.x * blockDim.x + threadIdx.x;
12     unsigned int y = blockIdx.y * blockDim.y + threadIdx.y;
13     int readX = (int)(x * fieldWidth * 1.f / imageWidth);
14     int readY = (int)(y * fieldHeight * 1.f / imageHeight);
15     float fieldValue = fieldData[readY * fieldWidth + readX];
16     float cb = materials[readY * fieldWidth + readX];
17     float4 color = make_float4(fieldValue, cb, 0, 1); color.w = threadIdx.
18         x == 0 || threadIdx.y == 0;
19     surf2Dwrite(color, image, x * sizeof(float4), y, cudaBoundaryModeClamp
    );
}

```

Listing 4.7. CUDA kernel for updating visualizer pixel buffer object

A completed PBO is bound to a uniform input for a shader (4.2.3,4.2.3) which renders the texture to the visualizer window.

```

1 #version 430 core
2 uniform sampler2D CudaTexture;
3 uniform vec2 TextureSize;
4 uniform vec2 WindowSize;
5 in vec3 Position;
6 out vec4 vColor;
7 out vec3 vWorldPosition;
8 out vec2 vTexCoord;
9 void main()
10 {
11     vec4 worldPos = vec4(Position, 1);

```

```

12 // scale down the quad coordinates to match the texture aspect ratio
13 if (TextureSize.x > TextureSize.y)
14     worldPos.y *= TextureSize.y / TextureSize.x;
15 else
16     worldPos.x *= TextureSize.x / TextureSize.y;
17 float windowAspect = WindowSize.x / WindowSize.y;
18 if (windowAspect > 1)
19     worldPos.y *= windowAspect;
20 else if (windowAspect < 1)
21     worldPos.x /= windowAspect;
22 if (abs(worldPos.x) > 1) worldPos.xy /= abs(worldPos.x);
23 if (abs(worldPos.y) > 1) worldPos.xy /= abs(worldPos.y);
24 vWorldPosition = worldPos.xyz;
25 vTexCoord = vec2(Position.x,-Position.y);
26 gl_Position = worldPos;
27
28 vColor = vec4(1);
29 }
```

Listing 4.8. GLSL Vertex Shaderr

```

1 #version 430 core
2 uniform sampler2D CudaTexture;
3 uniform float ColorScale = 100.f;
4 in vec2 vTexCoord;
5 in vec3 vWorldPosition;
6 in vec4 vColor;
7 out vec4 fragmentColor;
8 vec4 saturate(vec4 val) { return clamp(val,vec4(0),vec4(1)); }
9 void main() {
10     vec2 texCoord = vTexCoord / 2 + vec2(0.5);
11     vec4 texel = texture2D(CudaTexture,texCoord);
```

```

12 float r = texel.r * ColorScale;
13 float b = -r;
14 float dielectric = texel.g;
15 float g = (dielectric >= 0.5) ? 0 : dielectric * 2;
16 vec4 t = saturate(vec4(r,g,b,1));
17 fragmentColor = saturate(vec4(r,g,b,1));
18 }
```

Listing 4.9. GLSL Fragment Shader

This process runs continuously. However, the texture is only updated at the frequency requested by the user when the simulation was launched.

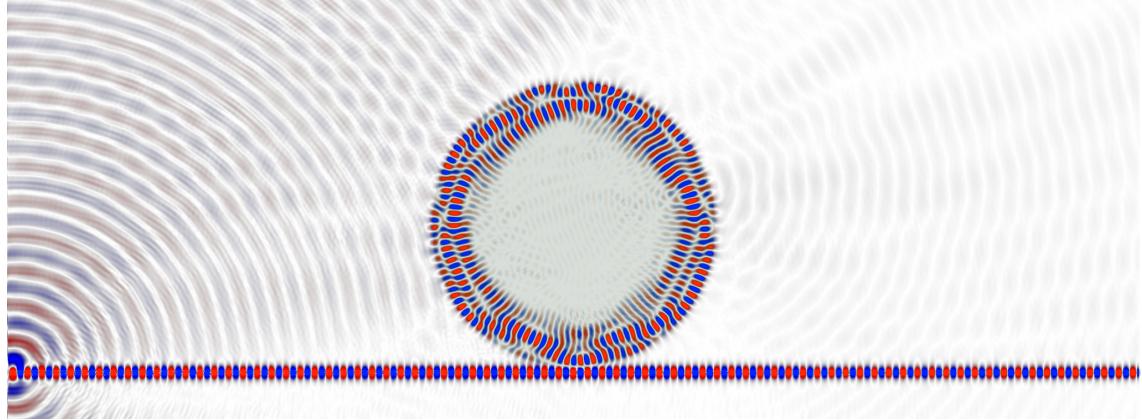


Figure 4.1. 2D Whispering Gallery Mode Sensor

### 4.3. Modeling approach

For simplicity, models are defined in any of a number of standard 32-bit color image formats.

In a 32-bit per pixel image, each color has 8-bit red, green, blue and alpha values. As mentioned in the model processor (subsection 4.2.1) section, each component is used to indicate some data about a given point in the simulation domain:

Table 4.2. Color component usage

Component	Meaning	Interpretation
Red	Source	normalized wavelength of the source
Green	Dielectric	$\epsilon_r = green * \frac{\epsilon_{max}}{255.0}$ )
Blue	Monitor	ID of the monitor to which this texel belongs
Alpha	Reserved	Reserved for future use.

Whenever a non-zero blue (monitor) value is encountered, it is used as an "ID" of a monitor. If the given ID has not yet been encountered, a new monitor is created. Texels with the same blue are added to the corresponding monitor. This allows monitors to be of any shape or size. Care must be taken during the design process to avoid nonsensical monitor configurations, such as disjoint or outlying pixels, or inconsistent thickness due to shape aliasing.

Using a tool such as Adobe Photoshop or Microsoft Paint, the user can specify all necessary data - sources, monitors and dielectric - in an intuitive fashion. Alternatively, these bitmaps could be generated by a custom tool which would voxelize a CAD model, assigning color components based on the model's metadata or object properties.

A significant advantage of this approach over the CSG method used in Meep is that all constructs - sources, monitors and dielectric - can be of any shape that can be drawn in a bitmap. In theory, any 3D voxel-based painting application could be used to build models in higher dimensions.

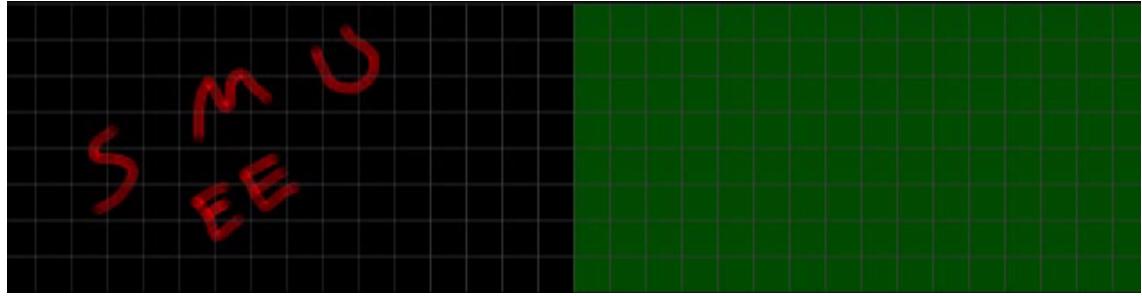


Figure 4.2. Arbitrarily-shaped source (Red pixels)

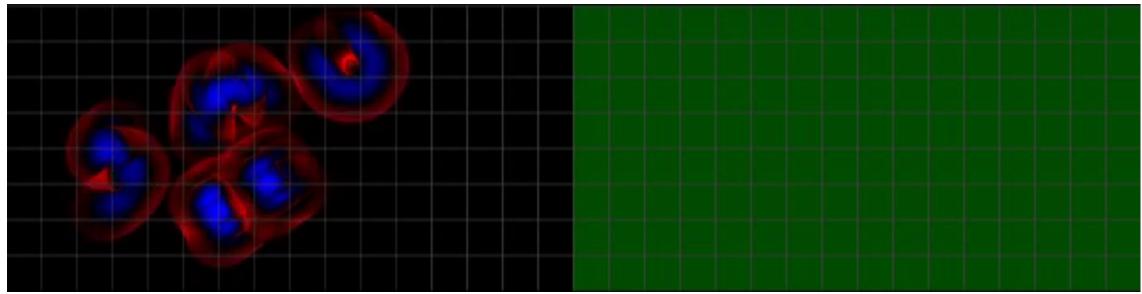


Figure 4.3. Arbitrarily-shaped source after 20 frames

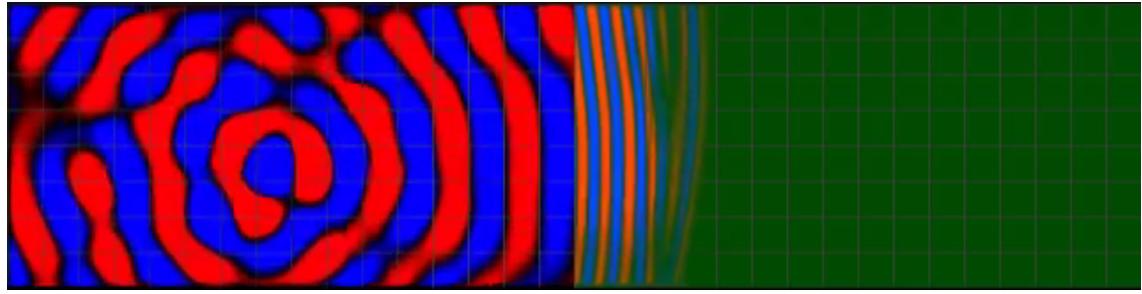


Figure 4.4. Arbitrarily-shaped source after 100 frames

#### 4.4. Testing and Validation Methodology

In order to validate the functionality of the simulator, a 2D  $TM_Z$  simulation was executed.

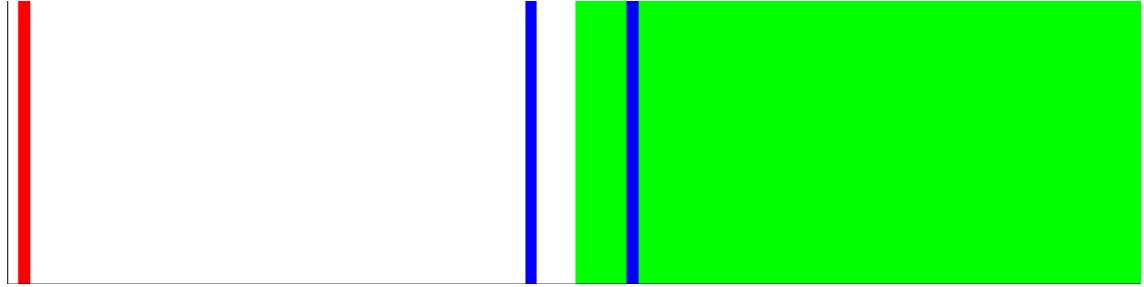


Figure 4.5.  $TM_Z$  Test Model with Plane Wave.

As can be seen in Figure 4.5, the simulation includes a plane wave source (red line), monitors for incident and transmitted waves, and a slab of dielectric with  $\epsilon_R = 9$ . An additional monitor at the source location is not clearly visible due to it's relatively low ID ( $B = 5$ ).

The simulation was first run with  $\epsilon_R = 1$  to record the incident magnitude in absence of any reflective interfaces.

The simulation was then run with  $\epsilon_R = 9$  in the dielectric slab area to record combined incidence and reflection, as well as transmittance within the dielectric.

In a post-processing step, the reflective magnitude was found by subtracting the incident wave magnitude obtained from the first simulation from the combined incidence and reflection magnitudes obtained from the second simulation.

$$|R| = |I + R| - |I| \quad (4.1)$$

Validation was performed by comparing the theoretical Fresnel coefficients for the test model with the time-averaged power (RMS) recorded during the simulation.

#### 4.4.1. Analytical Result

In the test configuration, a normalized ( $\lambda = 1$ )  $TM_Z$  plane wave is normally incident upon a dielectric interface with  $\epsilon_R = 9$ .

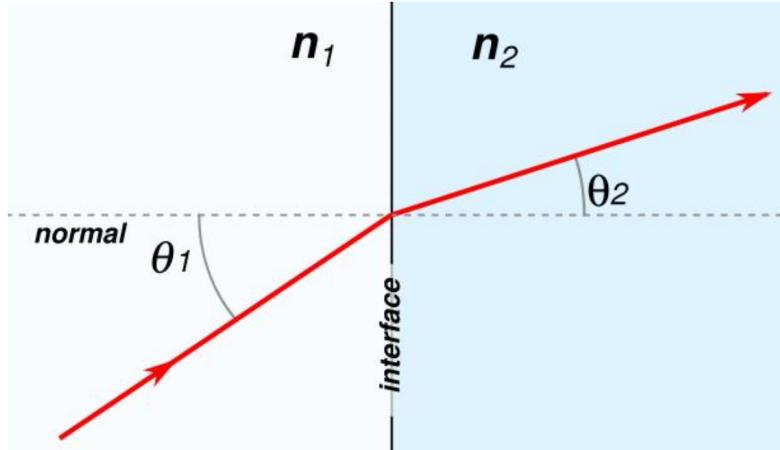


Figure 4.6. Snell's Law

Snell's Law states that the ratio of the refractive indices of the media at an interface, along with the angle of incidence, determine the angle of transmittance. (See Figure 4.6.)

Mathematically, this relationship may be expressed as

$$n_1 \sin \theta_1 = n_2 \sin \theta_2 \quad (4.2)$$

In the testmodel, the index  $n_2$  is calculated using the formula:

$$n = \sqrt{\epsilon_0 \epsilon_R \mu_0 \mu_R} \quad (4.3)$$

In our simulator,  $\epsilon_0$  and  $\mu_0$  are normalized to 1. Similarly, in the anisotropic media used in this case,

$$\mu_R = 1 \quad (4.4)$$

Using our test value of  $\epsilon_R = 9$  gives

$$n_2 = \sqrt{9} = 3 \quad (4.5)$$

For a normally-incident plane wave, the incident and refraction angles are

$$\theta_I = 0 \quad (4.6)$$

and

$$\theta_T = 0 \quad (4.7)$$

Evaluating the Fresnel equations for the reflection and transmission of a p-polarized wave,

$$r_p = \frac{n_2 \cos \theta_I - n_1 \cos \theta_T}{n_1 \cos \theta_T + n_2 \cos \theta_I} \quad (4.8)$$

$$t_p = \frac{2n_1 \cos \theta_I}{n_1 \cos \theta_T + n_2 \cos \theta_I} \quad (4.9)$$

gives the coefficients:

$$r_p = \frac{3 * 1 - 1 * 1}{1 * 1 + 3 * 1} = \frac{1}{2} \quad (4.10)$$

$$t_p = \frac{2 * 1 * 1}{1 * 1 + 3 * 1} = \frac{1}{2} \quad (4.11)$$

In this case, given a dielectric constant  $\epsilon_R = 9$ , the reflection and transmission coefficients are equal.

#### 4.4.2. Numerical Result

The observed steady-state output of the simulation for the baseline case ( $\epsilon_R = 1$ ) is shown in Figure 4.7.

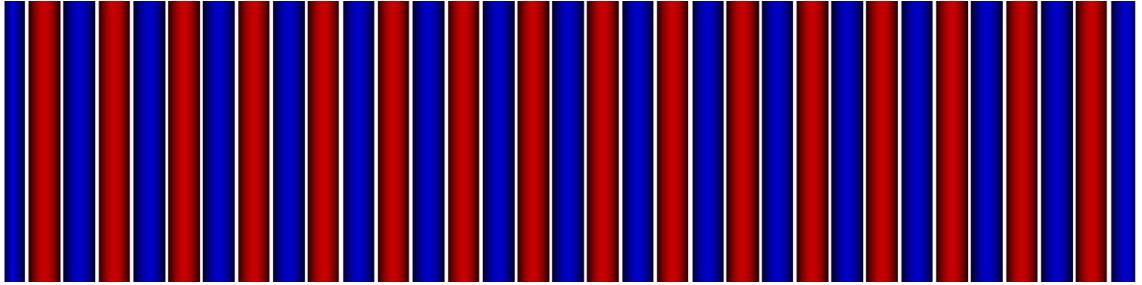


Figure 4.7. Plane wave with  $\epsilon_R = 1$

The time-averaged (RMS) numerical output for the incident and transmitted monitors is

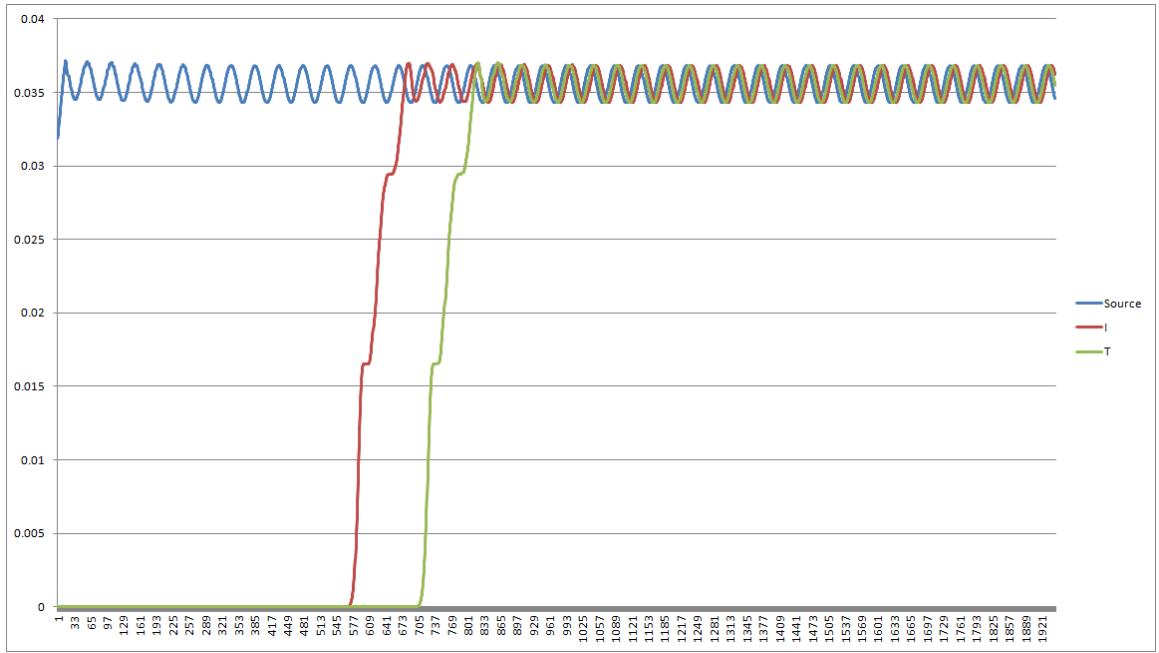


Figure 4.8. Time-averaged (RMS) output in free space  $\epsilon_R = 1$

As expected, the incident and transmitted magnitudes equal the source magnitude, offset by the time it takes for the wave to propagate from the source to the monitors.

The observed output of the same simulation, executed with dielectric constant

$$\epsilon_R = 9,$$

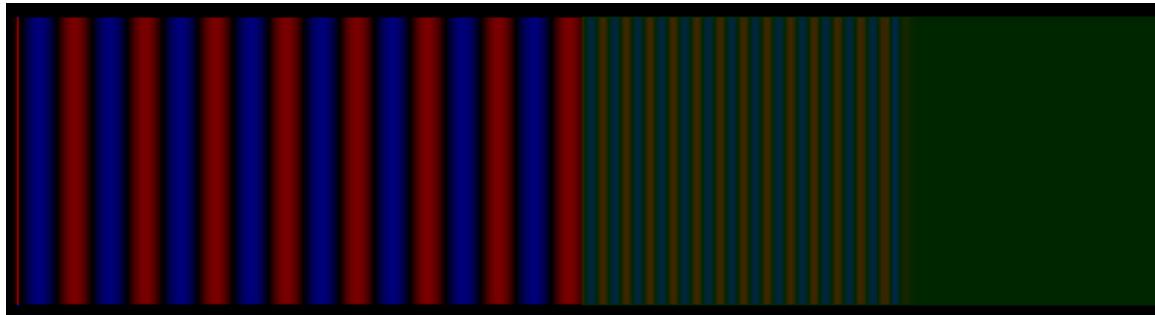


Figure 4.9. Steady state with  $\epsilon_R = 9$

Note the difference in  $\lambda$  once the wave enters the dielectric (green area). This simulation produces the following raw output:

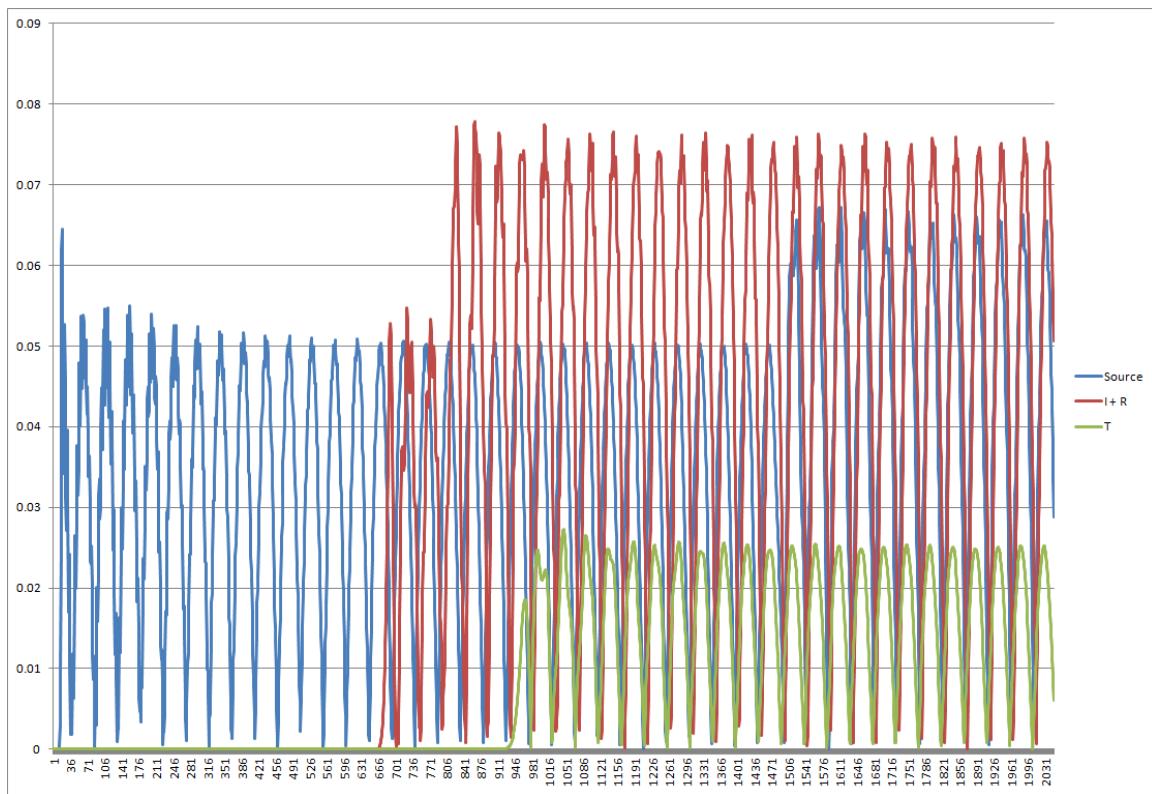


Figure 4.10. Raw monitor output with  $\epsilon_R = 9$

The time-averaged power for the monitors:

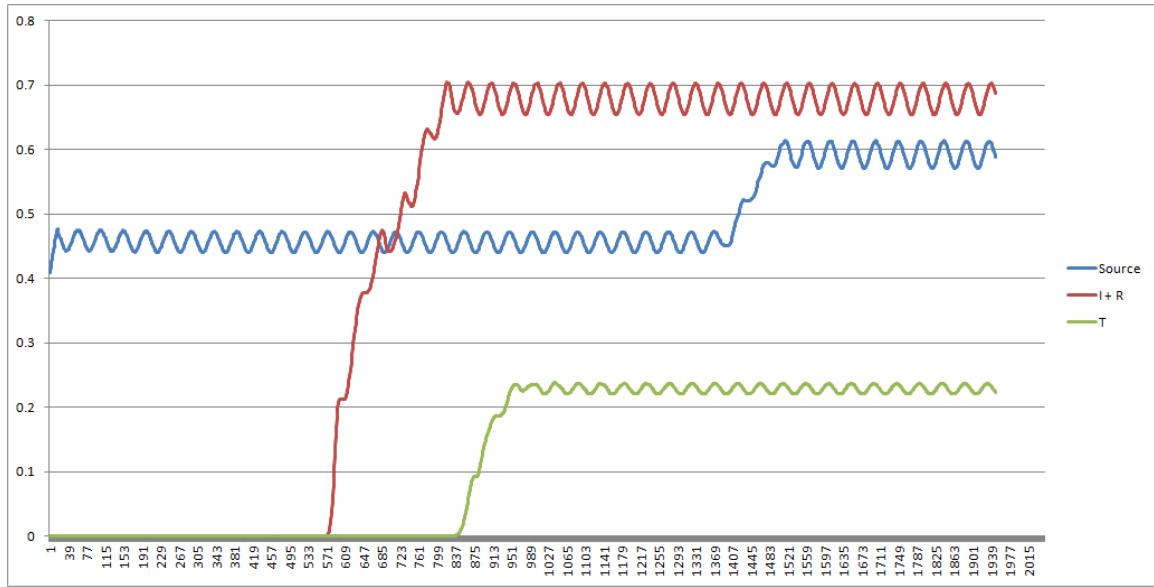


Figure 4.11. RMS output with  $\epsilon_R = 9$

Subtracting the incident magnitude from the simulation with  $\epsilon_R = 1$  from the second simulation with  $\epsilon_R = 9$  gives the following normalized result:

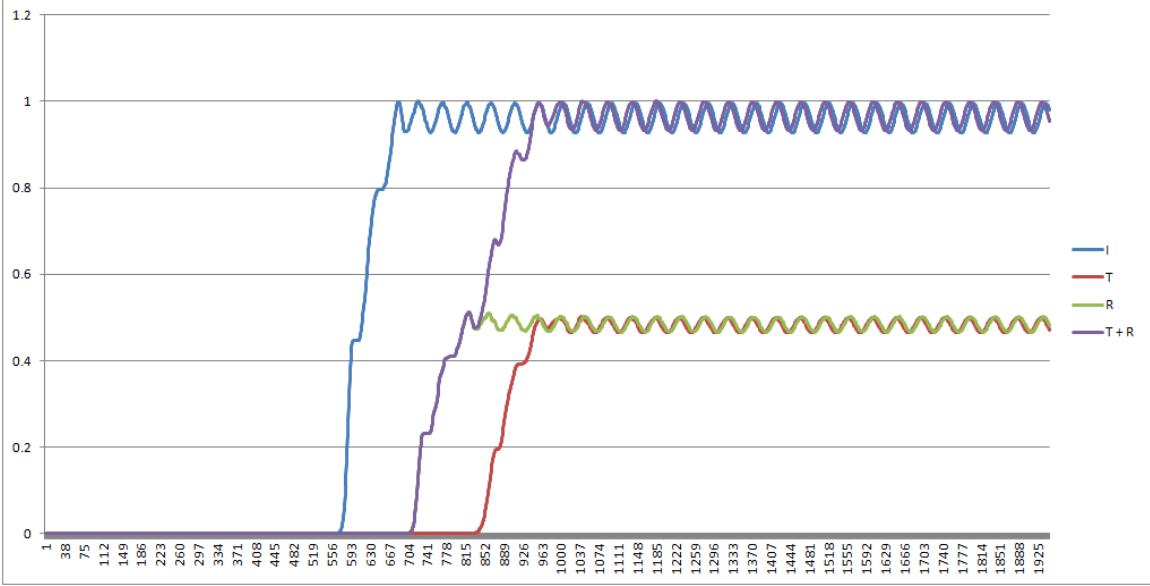


Figure 4.12. Normalized output with  $\epsilon_R = 9$

In this case, the average of each of the relevant values are determined to be:

$$I_{avg} = 0.96214\dots$$

$$T_{avg} = 0.4882\dots$$

$$R_{avg} = 0.484\dots$$

Conservation of energy requires that the sum of the transmitted (T) and reflected (R) magnitudes equals the incident magnitude (I).

$$I = T + R \quad (4.12)$$

The computed error of the computational result is

$$I_{error} = \left| \frac{I_{avg} - (T_{avg} + R_{avg})}{I - avg} \right| \quad (4.13)$$

Gives

$$\left| \frac{0.9621 - (0.4882 + 0.4843)}{0.9621} \right| = 1.011\%$$

Comparing to the analytically-derived Fresnel coefficients,

$$T = \frac{T_{avg}}{I_{avg}} = \frac{0.4882219...}{0.962135...} = 0.507 \quad (4.14)$$

$$R = \frac{R_{avg}}{I_{avg}} = \frac{0.484332.....}{0.962135...} = 0.503 \quad (4.15)$$

The error contribution of each component may be expressed as:

$$T_{error} = \left| 1 - \frac{T_{computed}}{T_{analytical}} \right| = \left| 1 - \frac{0.507}{0.5} \right| = 1.4\% \quad (4.16)$$

$$R_{error} = \left| 1 - \frac{R_{computed}}{R_{analytical}} \right| = \left| 1 - \frac{0.503}{0.5} \right| = 0.6\% \quad (4.17)$$

Careful analysis of different simulation results indicates that this error is due to floating-point precision errors. As every frame output depends on the state of the previous frame, small errors are compounded over time.

## Chapter 5

### Results

Metrics were calculated using domain sizes ranging from 128x64 to 8192x4096 over 5000 frames. (In this case, a frame represents a complete time step wherein both E and H fields are updated.) For benchmarking purposes, the visualizer was disabled.

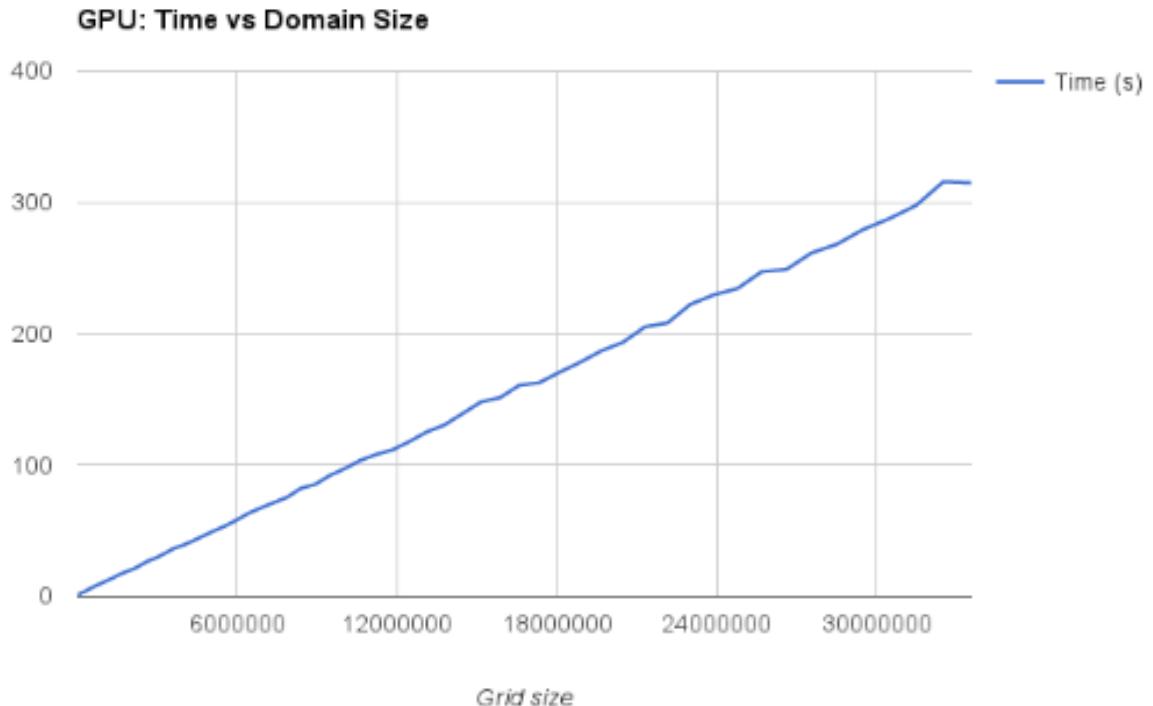


Figure 5.1. GoLightly: seconds for 5000 frames with the given domain size

As expected, As the number of computational cells in the simulation domain increases, the required computation time increases linearly.

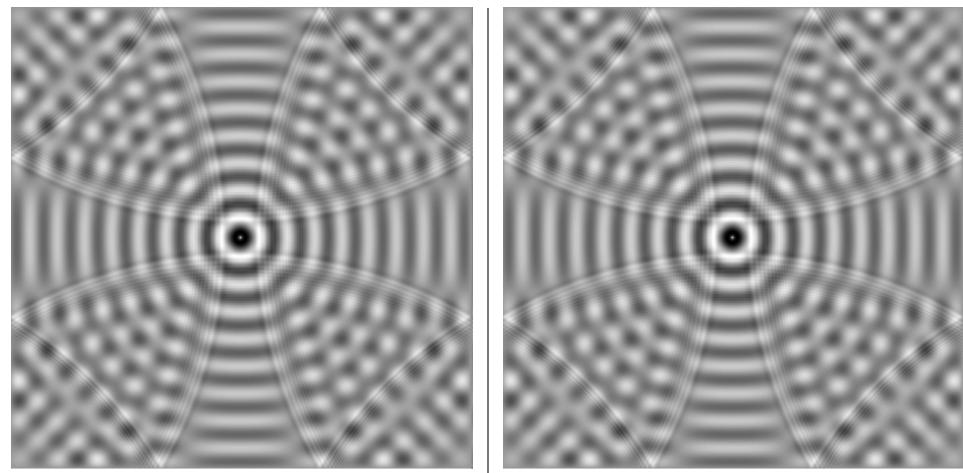


Table 5.1. blah

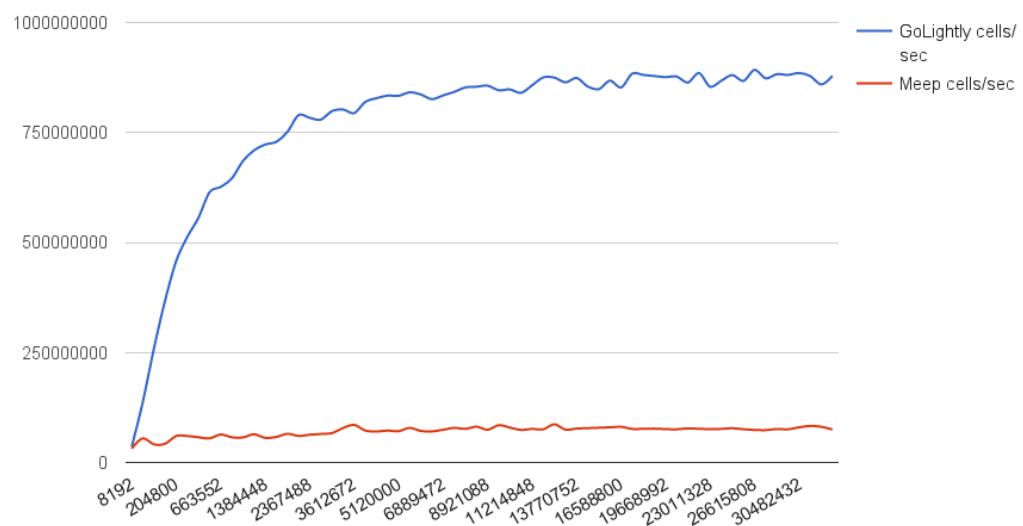


Figure 5.2. GoLightly: Completed E and H updates per second

Note that the total throughput (represented in the following figure as cell operations per second) increases dramatically as the domain size increases, until the point

where GPU initialization overhead is outweighed by computation time.

Comparing to Meep:

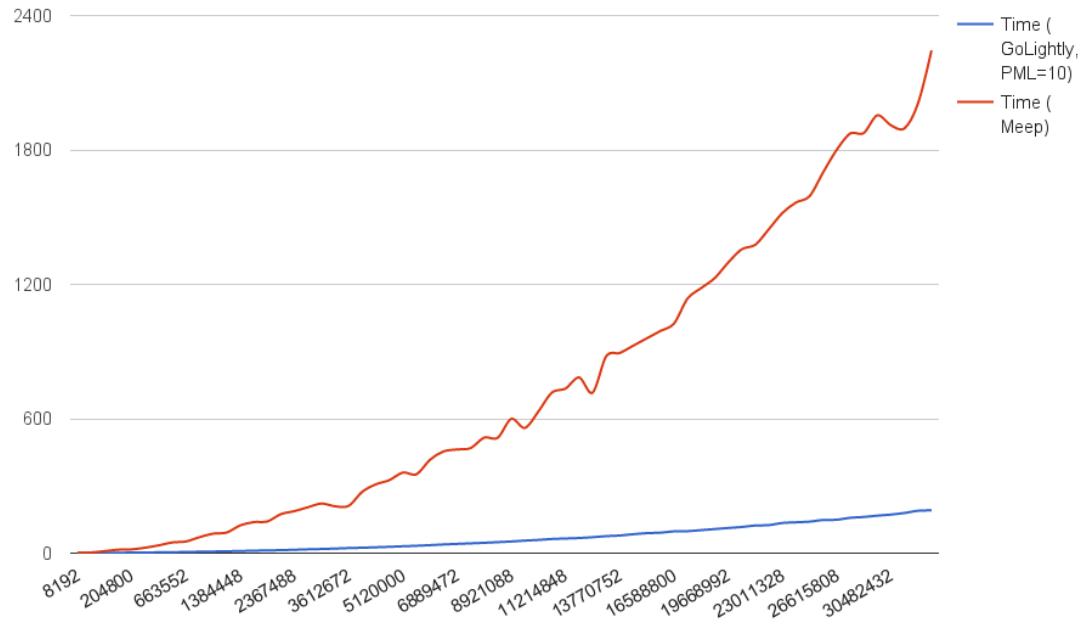


Figure 5.3. GoLightly vs Meep

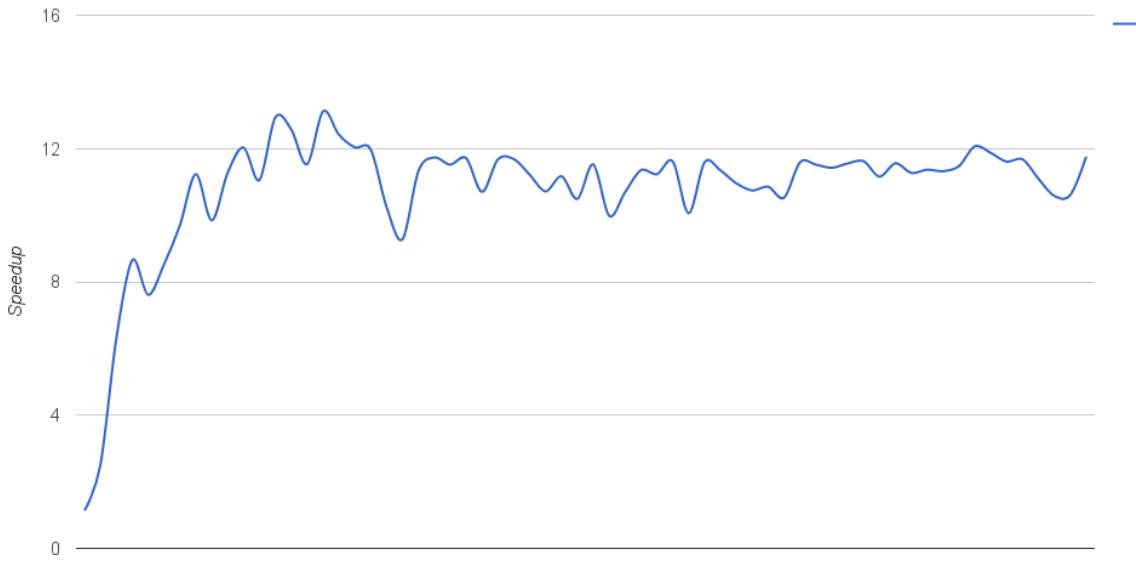


Figure 5.4. Speedup - Meep Time / GoLightly Time

This gives us a speedup ranging from 1.2 to 12, depending on the domain size. At lower resolutions, the overhead of initializing assets on the GPU can take more time than the simulation. In that case, a CPU-based simulation may outperform the GPU solution.

### 5.1. Optimization and Enhancements

Although a 1200% speed increase is significant, there is much room for improvement.

GPUs provide different memory spaces that vary in size and access speed. In addition to global device memory, each warp (group of 32 threads sharing an ALU)

has shared memory<sup>1</sup> and local memory<sup>2</sup>.

While global memory is the most flexible and plentiful - typically on the order of gigabytes on current generation hardware - it is also the slowest.

Shared memory is roughly equivalent to a CPU's L2 cache. It can be used for intra-thread synchronization and resource sharing. It is significantly faster than global memory.

Finally, local memory - similar to a CPU L1 cache - provides thread-local storage. Local memory provides the lowest latency of all of the memory spaces.

In its current form, GoLightly makes little use of shared or local memory. Modifying the application to take advantage of those may improve performance.

---

<sup>1</sup>Shared memory is physically local to the ALU and accessible by all threads within the warp.

<sup>2</sup>Each thread has local memory which is not accessible to any other thread

## Chapter 6

### Conclusions

Using the GoLightly GPU simulator, we have shown a speedup potential of 8 to 16 times that achieved by comparable software currently enjoying wide adoption.

#### 6.1. Usability

Although the software, in its current state, may not be suitable for production, it is clear that a GPU-based approach promises substantially improved performance. This results in increased productivity and, ultimately, better products. Eliminating the script-only modeling approach enforced by other packages makes utilization of this application significantly easier, and therefore accessible to a wider audience.

#### 6.2. Future Work

The vastly-improved performance afforded by this system presents some interesting opportunities:

##### 6.2.1. Genetic Algorithms

Genetic algorithms are intended to let software take over a part of the design process. By defining a problem in terms of a number of inputs and creating a fitness function, an application can potentially tests many different designs and use a feedback loop to suggest new permutations. This approach has been shown to be successful in such diverse fields as turbine design and pharmaceutical research.

A fast FDTD implementation will facilitate application of this technique to optical circuit design. By defining a circuit in terms of desired package size, available inputs, allowed waveguide shape, available dielectric properties, and others - and designing an appropriate fitness function, software will be able to rapidly evaluate different designs, and suggest new permutations. This may also dramatically reduce time-to-market and suggest new research avenues.

#### 6.2.2. Arbitrary Domain Shape and PML Sinks

Given the more flexible voxel-based model definition used in GoLightly, it becomes possible to create completely arbitrary domain shapes. Non-rectangular domains with PML "sinks" at any desired position within the domain would allow the designer to more tightly fit the computational domain to the circuit in question, and ignore areas that are not relevant to the task at hand.

### 6.3. Final Words

While not new, the field of general-purpose GPU computing offers potential that has to be realized in many of the areas where it shows the most promise. We pray that our investigation and demonstration detailed in this thesis encourages others to explore this area, and help to further reduce or eliminate boundaries that cripple or inhibit advances in the design process.

## REFERENCES

- [1]
- [2] CORPORATION, N. U.s. to build two flagship supercomputers for national labs, 2014.
- [3] MA, K.-L., XIE, J., AND YU, H. Scalable visualizing large 3d geodesic grid data with massively distributed gpus. In *In Proceedings of 2014 IEEE Large Data Analysis and Visualization Symposium (LDAV)* (Nov 2014), pp. 3–10.
- [4] MASSINGILL, B. L., MATTSON, T. G., AND SANDERS, B. A. Simd: An additional pattern for plpp (pattern language for parallel programming). In *Proceedings of the 14th Conference on Pattern Languages of Programs* (New York, NY, USA, 2007), PLOP ’07, ACM, pp. 6:1–6:15.
- [5] NICKOLLS, J., BUCK, I., GARLAND, M., AND SKADRON, K. Scalable parallel programming with cuda. *Queue* 6, 2 (Mar. 2008), 40–53.
- [6] NVIDIA CORPORATION. NVIDIA CUDA C programming guide, 2010. Version 3.2.
- [7] OSKOOI, A. F., ROUNDY, D., IBANESCU, M., BERMEL, P., JOANNOPOULOS, J. D., AND JOHNSON, S. G. MEEP: A flexible free-software package for electromagnetic simulations by the FDTD method. *Computer Physics Communications* 181 (January 2010), 687–702.
- [8] RAINA, R., MADHAVAN, A., AND NG, A. Y. Largescale deep unsupervised learning using graphics processors. In *International Conf. on Machine Learning* (2009).
- [9] RODRIGUEZ, P., PATTICHIS, M., AND JORDAN, R. Parallel single instruction multiple data (simd) fft: Algorithm and implementation. Tech. Rep. HPCERC2003-002, January 2003.
- [10] SHI, L., LIU, W., ZHANG, H., XIE, Y., AND WANG, D. A survey of gpu-based medical image computing techniques. *Quantitative Imaging in Medicine and Surgery* 2, 3 (2012).
- [11] STRZODKA, R., COHEN, J., AND POSEY, S. Gpu-accelerated algebraic multi-grid for applied {CFD}. *Procedia Engineering* 61 (2013), 381 – 387. 25th International Conference on Parallel Computational Fluid Dynamics.

- [12] WU, R., ZHANG, B., AND HSU, M. Clustering billions of data points using gpus. In *Proceedings of the combined workshops on UnConventional high performance computing workshop plus memory access workshop* (2009), ACM, pp. 1–6.
- [13] YEE, K. Numerical solution of initial boundary value problems involving maxwell’s equations in isotropic media. *IEEE Transactions on Antennas and Propagation* 14, 3 (1966), 302–307.