

# **Optimizing GPU Performance in Cylindrical FDTD Simulations**

**Dimitrios Giannakopoulos**

**Thesis submitted to the faculty of the Virginia  
Polytechnic Institute and State University in partial fulfillment  
of the requirements for the degree of**

**Master of Science**

**In**

**Computer Engineering**

**Zin Lin, Chair**

**Angelos Stavrou**

**Ravi Raghunathan**

**May 9<sup>th</sup> 2025**

**Arlington VA**

**Keywords: FDTD, GPU, Optimization, Simulation**

© 2025 Dimitris Giannakopoulos.

This thesis is licensed under a Creative Commons Attribution-NonCommercial 4.0  
International License.

<https://creativecommons.org/licenses/by-nc/4.0/>

# **Optimizing GPU Performance in Cylindrical FDTD Simulations**

**Dimitrios Giannakopoulos**

## **ABSTRACT**

Simulating large-area metasurfaces presents a major computational challenge due to their fine structural features and large physical dimensions. Traditional full-wave methods, such as finite-difference time-domain (FDTD), become infeasible for such problems due to excessive memory and runtime requirements. To address this, several approximate techniques have been developed, including the localized perturbation approximation (LPA), overlapping-domain approximation (ODA), and zoned discrete axisymmetry (ZDA), each balancing accuracy and efficiency for different metasurface geometries.

In this thesis, we focus on ZDA, a method tailored for metasurfaces with rotational symmetry. By expressing electromagnetic fields as a sum of angular modes and discretizing the radial domain into concentric zones, ZDA reduces a 3D problem to a series of much smaller and fewer simulations. This dimensionality reduction enables accurate modeling of freeform optical devices with fine resolution using modest computational resources. We implement this approach via a GPU-accelerated FDTD solver in cylindrical coordinates, enabling scalable and efficient simulation of broadband, high-performance metasurfaces. Our results demonstrate that symmetry-aligned simulation strategies such as ZDA can unlock practical design workflows for metasurfaces previously beyond reach. My personal work though was mostly on accelerating our FDTD code using the power given by the GPUs.

# **Optimizing GPU Performance in Cylindrical FDTD Simulations**

**Dimitrios Giannakopoulos**

## **GENERAL AUDIENCE ABSTRACT**

Simulating how light interacts with large, complex optical surfaces, known as metasurfaces, is essential for designing new technologies like ultra-thin lenses. These surfaces often span millimeters in size but contain tiny features measured in nanometers, making them extremely difficult to model with standard simulation tools. The most accurate methods, such as finite-difference time-domain (FDTD), quickly become too slow and memory-intensive for practical use at this scale.

In this project, the focus was on making these simulations faster and more efficient. I worked on accelerating the FDTD method that takes advantage of the symmetrical shape of certain metasurfaces. This involved adapting the simulation to cylindrical coordinates and running it on graphics processors (GPUs), which are well-suited for high-speed parallel computation. The result is a simulation tool that can handle large-area metasurfaces with much lower computing cost opens the door to faster design cycles and more ambitious optical applications.

## **Aknowledgements**

I would like to express my deepest gratitude to my advisor, Dr. Zin Lin, for his invaluable guidance, patience, and encouragement throughout my thesis work as well as our overall cooperation for this group's project. I am also thankful to the members of my committee for their insightful feedback and suggestions. I acknowledge Virginia Tech for providing the resources and environment that made this research possible. Finally, I am truly grateful to my family and friends for their unwavering support and understanding during this journey.

# Chapter 1: Contents

Chapter 2: Introduction and Motivation: FDTD Algorithm for Freeform Metalens Design .....	7
Chapter 3: The Finite-Difference Time-Domain Method in Cylindrical Coordinates with Discrete Axial Symmetry .....	11
3.1 Introduction.....	11
3.2 Maxwell's Equations in Cylindrical Coordinates .....	12
3.3 Exploiting Discrete Axial Symmetry: Fourier Decomposition in $\theta$ .....	15
3.4 Yee Grid and Discretization in Cylindrical Coordinates .....	18
3.5 Discrete Update Equations for Interior Points .....	21
3.6 Special Treatment at the Axis ( $r = 0$ ).....	23
3.7 Summary and considerations .....	29
Chapter 4: GPU Computing.....	34
4.1 What is a GPU? A Deep Dive into Graphics Processing Units .....	34
4.1.1 Historical Evolution of GPUs .....	34
4.1.2 GPU Hardware Architecture.....	35
4.1.3 Instruction Scheduling and Threading Model.....	36
4.1.4 Performance Metrics and Theoretical Limits .....	36
4.2 GPU vs CPU in Scientific Computing.....	37
4.2.1 Core Architectural Philosophy.....	41
4.2.2 Memory Hierarchy and Latency Handling .....	41

4.2.3 Programming and Execution Models.....	42
4.2.4 Performance Comparison in Real Workloads.....	42
4.2.5 Scientific Relevance and Shifting Paradigms .....	43
4.3 GPU Programming Ecosystem .....	44
4.4 GPU Optimization Techniques .....	47
4.5 GPUs in FDTD and Electromagnetic Simulations .....	52
4.6 Accelerating Axissymmetric FDTD using CUDA .....	57
Chapter 5: Translating the FDTD Method to GPU-Accelerated Code.....	64
5.1 Introduction.....	64
5.2 Grid setup and Data Structures .....	65
5.3 FDTD Update Kernels on the GPU .....	68
5.3.1 V2p Kernel – Updating E-field from H-field .....	69
5.3.2 P2u kernel – Updating PML u fields .....	74
5.3.3 u2q Kernel – Updating second set of PML fields.....	76
5.3.4 Q2v kernel – Updating H-field from q-field.....	78
5.4 GPU Execution Flow (Time-Stepping loop) .....	79
5.5 Performace Considerations and Optimizations Strategies .....	81
5.6 Conclusion .....	87
Chapter 6: Results and Future Work.....	89
6.1 Performance Results .....	89
6.2 Future Work .....	94
Chapter 7: Bibliography.....	96

# Chapter 2: Introduction and Motivation: FDTD Algorithm for Freeform Metalens Design

Simulating the electromagnetic behavior of large-area metasurfaces poses a profound computational challenge. As metasurfaces advance toward real-world applications requiring broad apertures and high functionality, such as metalenses spanning millimeters to centimeters in diameter, the demands placed on numerical solvers escalate dramatically. These structures often contain feature sizes at the scale of hundreds of nanometers while spanning physical dimensions more than three thousand times the wavelength of interest. For such cases, traditional full-wave methods, including finite-difference time-domain (FDTD) and finite element methods (FEM), become prohibitively expensive. The memory and runtime scale poorly with problem size, especially in three dimensions, leading to simulation domains that can exceed hundreds of billions of voxels and terabytes of storage. Consequently, new techniques are essential to bridge the gap between physical accuracy and computational feasibility.

To overcome these limitations, several classes of simulation techniques have been developed that exploit the underlying physics and symmetries of metasurfaces. These approaches aim to reduce dimensionality, localize computational complexity, or introduce analytical approximations that enable the simulation of structures previously out of reach. Among these, the localized perturbation approximation (LPA), overlapping-domain approximations (ODA), and zoned discrete axisymmetry (ZDA) have emerged as powerful frameworks tailored to different metasurface configurations. Each method embodies a unique balance between efficiency and accuracy and is best suited to specific geometric or material contexts.

The localized perturbation approximation begins with the assumption that the electromagnetic response of a metasurface can be decomposed into contributions from individual meta-atoms, treated as local perturbations to a background field. In this framework, the behavior of each meta-atom is simulated independently—usually under periodic boundary conditions—and the global response of the metasurface is approximated by superposing these local solutions. This method relies on the assumption that each meta-atom’s behavior is primarily influenced by its immediate environment and that distant interactions can be neglected or treated approximately. The benefit of this approach is that it significantly reduces the computational domain, allowing the simulation of only a small, representative portion of the metasurface. For metasurfaces with slowly varying spatial profiles or weakly interacting elements, LPA provides remarkably efficient approximations with meaningful physical insight. However, its accuracy diminishes in the presence of dense packing, strong inter-element coupling, or rapid spatial variation. In such cases, the foundational assumption of weak perturbation breaks down, and LPA may fail to capture essential electromagnetic behavior.

In response to the limitations of LPA, overlapping-domain approximations offer a more robust strategy. These methods divide the metasurface into smaller, partially overlapping regions, each of which is simulated with full-wave precision. The overlaps ensure that field interactions and boundary conditions are properly captured between neighboring domains, mitigating the artifacts that may arise from artificial truncation. The overlapping regions act as bridges, enforcing electromagnetic continuity and allowing stronger local interactions to be faithfully modeled. Compared to LPA, this approach provides improved accuracy, particularly for metasurfaces with complex spatial variations, abrupt design changes, or densely packed unit cells. Moreover, it enables more general geometries beyond the periodic or slowly varying types that LPA is best suited for. Nonetheless, overlapping-domain approaches require careful calibration of the domain decomposition, interpolation across overlaps, and solution merging strategies, which can introduce both algorithmic and computational overhead. Yet for large-area metasurfaces where local interactions dominate and rotational symmetry is absent, this method presents a viable compromise between brute-force simulation and oversimplified approximations.

While both LPA and ODA provide ways to approximate large metasurfaces in Cartesian coordinates, many metasurface designs—particularly those used in imaging optics—exhibit cylindrical or near-cylindrical symmetry. These structures are characterized by radial variation and rotational invariance, suggesting that the simulation method should exploit this symmetry to gain computational leverage. Zoned discrete axisymmetry is a technique explicitly designed for such cases. It discretizes the azimuthal coordinate and expresses the electromagnetic fields as a superposition of a finite number of angular modes. Each mode corresponds to a specific azimuthal variation, allowing the full three-dimensional problem to be decomposed into a set of two-dimensional problems in the radial and axial directions. This representation reduces the dimensionality of the problem while preserving the essential physics of field propagation in symmetric geometries. The radial domain is further divided into concentric zones, each of which may contain different meta-atom designs or material properties. Within each zone, the angular dependence is captured using a finite Fourier basis, and periodic boundary conditions are applied in the azimuthal direction.

This approach offers a series of critical advantages. First, it aligns the computational strategy with the natural symmetry of the physical system, yielding high-fidelity simulations with drastically reduced memory and runtime requirements. Second, it enables detailed modeling of radial variation, which is crucial in the design of freeform lenses and multi-functional optical components. Third, by using zoned decomposition, the method allows for gradual changes in design along the radial direction without sacrificing angular resolution. In practice, this means that structures extending across millimeters of physical space and nanometer-scale resolution can be simulated using modest computational resources. Discrete axisymmetry also simplifies the application of periodic boundary conditions, avoids the need for complex domain meshing, and introduces a more structured simulation pipeline. However, it is important to note that this method is applicable only to rotationally symmetric metasurfaces. For geometries lacking such symmetry, overlapping-domain approximations may remain the more general solution.

In this thesis, we adopt the zoned discrete axisymmetry approach to simulate and optimize large-area, freeform metasurfaces for broadband optical applications. The choice is motivated by the rotational symmetry of the devices under study, the need for high spatial resolution, and the computational infeasibility of direct 3D full-wave methods. By combining the ZDA technique with a GPU-accelerated cylindrical-coordinate FDTD solver, we achieve scalable, high-accuracy simulation capabilities that enable the design and characterization of metasurfaces far beyond what was previously practical. This introduction sets the stage for the methodological and computational developments that follow, each of which builds on the foundational concept of symmetry-aligned simulation for large-scale photonic design.

# Chapter 3: The Finite-Difference Time-Domain Method in Cylindrical Coordinates with Discrete Axial Symmetry

## 3.1 Introduction

The Finite-Difference Time-Domain (FDTD) method is one of the most widely used numerical techniques for solving Maxwell's equations in both academic research and industry applications. Since its inception by Kane Yee in 1966 [1], the FDTD method has demonstrated remarkable flexibility, providing accurate full-wave solutions for a broad range of electromagnetic problems [2]. Its time-domain formulation allows it to capture broadband behavior with a single simulation run, making it ideal for modeling the interaction of light with complex structures such as photonic crystals, waveguides, and metasurfaces.

However, the standard formulation of FDTD in Cartesian coordinates becomes inefficient when applied to problems involving strong geometric symmetry, particularly cylindrical or rotational symmetry. In such cases, the natural geometry of the structure is misaligned with the Cartesian grid, leading to poor convergence, increased staircasing errors, and an unnecessarily large computational domain. This is especially problematic for large-area metasurfaces, where the lateral extent of the structure may span thousands of wavelengths, and where rotational symmetry is a fundamental design feature.

To address these challenges, this chapter presents a detailed formulation of the FDTD algorithm in cylindrical coordinates, enhanced with discrete axial symmetry. By explicitly leveraging the angular periodicity of rotationally symmetric structures, the

simulation domain can be dramatically reduced in dimensionality, while preserving the full physics of the problem. This technique is especially well-suited for metasurfaces such as metalenses, axicons, and spiral phase plates, where the structure is invariant (or slowly varying) in the azimuthal direction.

The cylindrical-coordinate FDTD method begins by reformulating Maxwell's equations in the  $(r,\theta,z)$  coordinate system. This transformation introduces new terms not present in the Cartesian formulation, including  $1/r$  scaling factors and more complex expressions for the curl and divergence. In addition, by decomposing the angular variation into discrete Fourier modes, the method reduces the problem from three dimensions to a series of two-dimensional problems—one per azimuthal mode. This modal decomposition, or discrete axisymmetry, is a powerful computational strategy that reduces memory usage and accelerates simulation while retaining all the essential electromagnetic behavior.

The remainder of this chapter develops the FDTD algorithm in cylindrical coordinates in a rigorous and step-by-step fashion. We begin by presenting the transformation of Maxwell's equations into the cylindrical basis, including all terms and coordinate dependencies. We then introduce the discrete Fourier expansion in the azimuthal direction and explain how it enables decoupling of the problem into independent modal equations. Following that, we construct the staggered Yee grid in cylindrical coordinates, derive the discrete update equations for each field component, and describe the time-stepping strategy. We then formulate the absorbing boundary conditions using stretched-coordinate perfectly matched layers (SC-PML) in both radial and axial directions. Finally, we address the mathematical and numerical challenges posed by the singularity at the origin ( $r=0$ ), and explain how the algorithm remains stable and accurate even at this critical boundary.

## 3.2 Maxwell's Equations in Cylindrical Coordinates

In cylindrical coordinates, a vector field  $\mathbf{A}$  has components  $(A_r, A_\theta, A_z)$  along the unit vectors  $(\mathbf{r}, \boldsymbol{\theta}, \mathbf{z})$ . The curl of  $\mathbf{A}$  in cylindrical coordinates is given by:

Substituting these expressions into Maxwell's curl equations yields six coupled partial differential equations for the time evolution of the electric and magnetic field components in the  $(r,\theta,z)$  basis.

Explicitly, the time evolution of the electric field components is given by:

These expressions demonstrate the coupling of field components and the dependence on both radial and angular derivatives. The presence of  $1/r$  terms makes the implementation nontrivial, particularly at the origin, where care must be taken to avoid division by zero.

The goal of the FDTD algorithm is to discretize these equations both in time and space, enabling explicit numerical integration over successive time steps on a structured grid. However, before discretization can proceed, we introduce an essential mathematical

$$\nabla \times \mathbf{A} = \left( \frac{1}{r} \frac{\partial A_z}{\partial \theta} - \frac{\partial A_\theta}{\partial z} \right) \hat{r} + \left( \frac{\partial A_r}{\partial z} - \frac{\partial A_z}{\partial r} \right) \hat{\theta} + \left( \frac{1}{r} \frac{\partial (r A_\theta)}{\partial r} - \frac{1}{r} \frac{\partial A_r}{\partial \theta} \right) \hat{z}$$

simplification — the separation of angular dependence through Fourier decomposition, leading to the concept of discrete axial symmetry.

Here  $A_r$ ,  $A_\theta$ , and  $A_z$  are the components of  $\mathbf{A}$  along  $r$ ,  $\theta$ ,  $z$  respectively. This formula encapsulates the key differences from the Cartesian curl: the presence of  $1/r$  factors and the coupling of  $r$ - and  $\theta$  derivatives in the  $z$ -component. Using this, we can write out Maxwell's curl equations (Faraday's law and Ampère–Maxwell law) in cylindrical component form.

- R-component (radial):

$$\frac{1}{r} \frac{\partial E_z}{\partial \theta} - \frac{\partial E_\theta}{\partial z} = -\mu \frac{\partial H_r}{\partial t}$$

- $\Theta$ -component (azimuthal):

$$\frac{\partial E_r}{\partial z} - \frac{\partial E_z}{\partial r} = -\mu \frac{\partial H_\theta}{\partial t}$$

- Z-component (axial):

$$\frac{1}{r} \frac{\partial(rE_\theta)}{\partial r} - \frac{1}{r} \frac{\partial E_r}{\partial \theta} = -\mu \frac{\partial H_z}{\partial t}$$

Similarly, Ampère's law (with Maxwell's displacement current) is written as:

$$\nabla \times \mathbf{H} = \varepsilon \partial \mathbf{E} / \partial t + \mathbf{J}$$

- R-component (radial):

$$\frac{1}{r} \frac{\partial H_z}{\partial \theta} - \frac{\partial H_\theta}{\partial z} = \varepsilon \frac{\partial E_r}{\partial t} + J_r$$

- $\Theta$ -component (azimuthal):

$$\frac{\partial H_r}{\partial z} - \frac{\partial H_z}{\partial r} = \varepsilon \frac{\partial E_\theta}{\partial t} + J_\theta$$

- Z-component (axial):

$$\frac{1}{r} \frac{\partial(rH_\theta)}{\partial r} - \frac{1}{r} \frac{\partial H_r}{\partial \theta} = \varepsilon \frac{\partial E_z}{\partial t} + J_z$$

In the above,  $\epsilon = \epsilon(r, \theta, z)$  and  $\mu = \mu(r, \theta, z)$  are the permittivity and permeability (which may vary with position), and  $J = (J_r, J_\theta, J_z)$  is an applied source current density. These equations are the continuous-space Maxwell equations in cylindrical coordinates. They form the starting point for our FDTD algorithm. Each equation corresponds to one component of the curl and describes how a spatial variation of certain field components produces a time derivative of the perpendicular field component. Note the presence of terms like  $(1/r) \partial/\partial\theta$  and  $(1/r) \partial(r\cdot)/\partial r$  which are unique to curvilinear coordinates and will require careful handling in the discrete scheme. In particular, at  $r = 0$  the  $1/r$  factors appear singular hinting that special treatment will be needed on the axis (we will return to this in detail later). Aside from these, the equations have a structure analogous to Cartesian FDTD equations: each component's time derivative is driven by spatial derivatives of the other field components in a curl-like fashion.

If a problem is completely axisymmetric, all  $\partial/\partial\theta$  terms drop out (i.e., there is no dependence on  $\theta$ ), and the equations simplify greatly. In fact, for axisymmetric  $\partial/\partial\theta$  and source-free ( $J=0$ ) cases, the fields decouple into transverse magnetic (TM<sub>z</sub>:  $E_z, H_r, H_\theta$ ) and transverse electric (TE<sub>z</sub>:  $H_z, E_r, E_\theta$ ) polarizations with respect to the  $z$ -axis [2]. However, in this chapter we consider the more general case of *discrete axial symmetry*, where there is a repetitive  $\theta$ -dependence that is not trivial but can be handled via Fourier decomposition. This approach retains generality (allowing higher-order angular variations) while still exploiting symmetry to reduce the simulation domain.

### 3.3 Exploiting Discrete Axial Symmetry: Fourier Decomposition in $\theta$

Discrete axial symmetry means that the structure and excitations are periodic in the azimuthal direction with some period  $\phi$  (where  $\phi=2\pi N$  for  $N$ -fold symmetry). For example, a structure might consist of  $N$  identical sectors repeated around the circle. Even

if the fields are not strictly the same in each sector, Maxwell's equations guarantee that the field solutions can be expanded in a set of angular harmonics that reflect this periodicity [3]. In practical terms, any field in such a system can be written as a superposition of modes of the form  $e^{im\theta}$  (where  $m$  is an integer), which represent variations of order  $m$  around the axis. This is essentially a Fourier series in the  $\theta$  coordinate.

Mathematically, we express each Cartesian field component as:

$$E_r(r, \theta, z, t) = \sum_m E_r^{(m)}(r, z, t) e^{im\theta}$$

and similarly  $E(r, \theta, z, t) = \sum_m E^{(m)}(r, z, t) \cdot \exp(im\theta)$ , and for the magnetic field  $H_r = \sum_m H_r^{(m)}(r, z, t) \cdot \exp(im\theta)$ , etc. The summation is over integer harmonics  $m = 0, \pm 1, \pm 2, \dots$ , and the time dependence is carried by the complex amplitudes  $E^{(m)}(r, z, t)$  (or one may treat  $E^{(m)}(r, z, t)$  as phasor amplitudes for harmonic time dependence, but here we keep time explicit). If the structure is strictly real and periodic, one typically finds that only a subset of  $m$  (such as  $m=0$  to  $N-1$ ) are needed for a complete basis due to symmetry; also negative  $m$  are essentially the complex conjugates of positive  $m$  modes for real fields. In a simulation, one might choose to simulate a single harmonic  $m$  of interest rather than summing over all  $m$ . Linearity of Maxwell's equations guarantees that each  $m$ -harmonic evolves independently: if initially only a certain  $m$  is excited, no other  $m'$  will be generated during the evolution (unless there are nonlinearities or asymmetric perturbations). This superposition principle is powerful – it means we can *separate the full 3D problem into independent 2D problems*, one for each harmonic index  $m$  [3]. In practice, one chooses the relevant  $m$  (or a small set of  $m$  values) based on the symmetry of the excitation or the desired solution, and then performs an FDTD simulation for each of those harmonics on a reduced domain.

To incorporate this decomposition into the FDTD formulation, we *insert the Fourier form* into the cylindrical Maxwell equations and isolate a single harmonic  $m$ . Let us assume a single harmonic dependence  $\exp(im\theta)$  for all fields (and similarly

$$\frac{\partial}{\partial \theta} \left( E_r^{(m)}(r, z, t) e^{im\theta} \right) = im E_r^{(m)}(r, z, t) e^{im\theta}$$

decompose any sources  $J_r^{(m)} \cdot \exp(i \cdot m \cdot \theta)$  etc.). Because the factor  $\exp(i \cdot m \cdot \theta)$  is known and carries all the  $\theta$ -dependence, any  $\theta$ -derivative can be evaluated analytically:

since  $E^{(m)}(r, z, t)$  no longer depends on  $\theta$  explicitly (it is the amplitude function depending on  $r, z, t$  only). We can therefore replace the partial derivative  $\partial/\partial\theta$  by multiplication by  $i \cdot m$  in the equations for the  $m$ -th harmonic. Furthermore, the common factor  $\exp(i \cdot m \cdot \theta)$  will cancel out from every term, yielding differential equations for the complex amplitudes  $E^{(m)}(r, z, t)$  and  $H^{(m)}(r, z, t)$  alone. For each harmonic  $m$ , the resulting equations are:

- From Faraday's law (for the  $m$ -th harmonic):  $(i \cdot m / r) \cdot E_z^{(m)}(r, z, t) - \partial E_{\theta}^{(m)}(z)/\partial z = -\mu \partial H_r^{(m)}/\partial t - \partial E_r^{(m)}/\partial z - \partial E_z^{(m)}/\partial r = -\mu \partial H_{\theta}^{(m)}/\partial t$   
 $(1/r) \partial(r \cdot E_{\theta}^{(m)})/\partial r - (i \cdot m / r) \cdot E_r^{(m)} = -\mu \partial H_z^{(m)}/\partial t$
- From Ampère's law (for the  $m$ -th harmonic):  $(i \cdot m / r) \cdot H_z^{(m)} - \partial H_{\theta}^{(m)}/\partial z = \varepsilon \partial E_r^{(m)}/\partial t + J_r^{(m)} \partial H_r^{(m)}/\partial z - \partial H_z^{(m)}/\partial r = \varepsilon \partial E_{\theta}^{(m)}/\partial t + J_{\theta}^{(m)} (1/r)$   
 $\partial(r \cdot H_{\theta}^{(m)})/\partial r - (i \cdot m / r) \cdot H_r^{(m)} = \varepsilon \partial E_z^{(m)}/\partial t + J_z^{(m)}$

These are the reduced Maxwell's equations for the  $m$ -th Fourier component. They are 2+1 dimensional (two spatial coordinates  $r, z$  plus time) and contain the parameter  $m$ . All explicit  $\theta$  derivatives have been eliminated, replaced by algebraic terms involving  $i \cdot m$ . Note that for  $m = 0$  (the purely axisymmetric case), all the  $i \cdot m / r$  terms drop out, recovering the expected simpler form. For  $m \neq 0$ , those terms couple, for instance,  $E_z$  to  $H_r$  and  $E_r$  to  $H_z$  in the equations above (they stem from the original  $\partial/\partial\theta$  coupling in the curl). Physically,  $m \neq 0$  corresponds to fields that vary as  $\cos(m \cdot \theta)$  or  $\sin(m \cdot \theta)$  around the axis (if one takes the real part of  $\exp(i \cdot m \cdot \theta)$ ), which represent azimuthally higher-order modes; the  $i \cdot m$  terms essentially account for the twisting of the field lines around the axis. Each  $m$  can be interpreted as the number of field oscillations around a  $2\pi$  rotation.

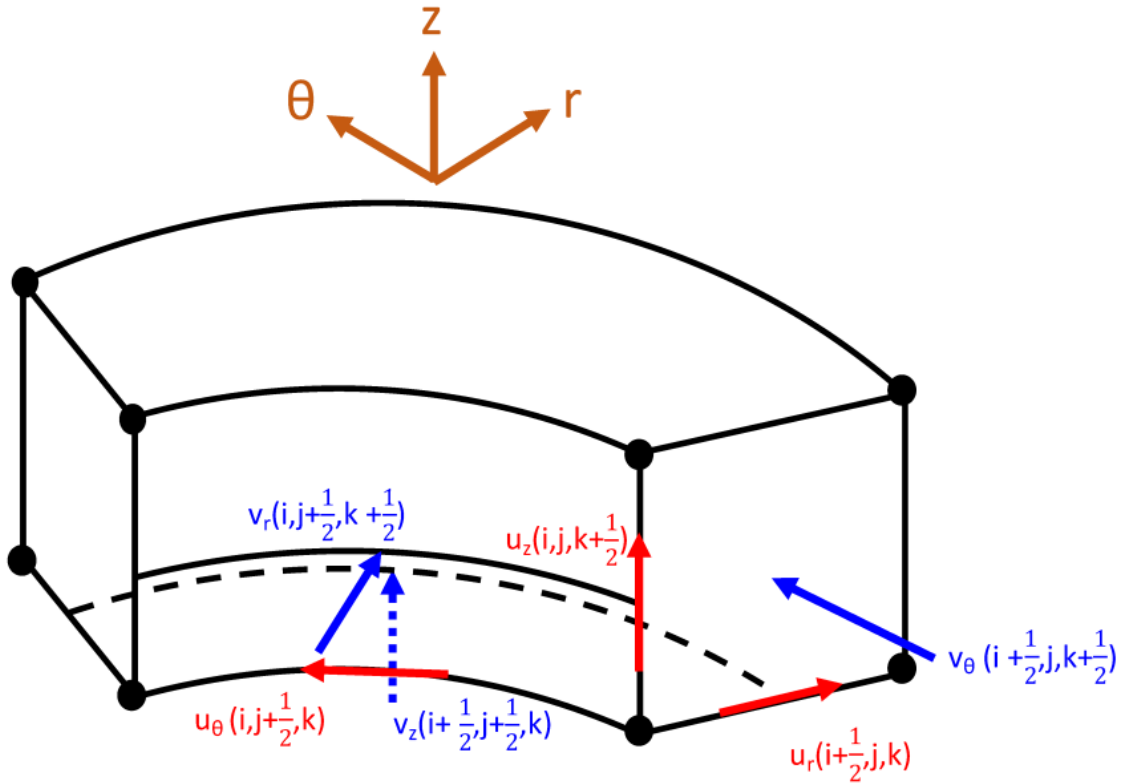
An important consequence of this formulation is that if the structure has  $N$ -fold symmetry with  $\varphi = 2\pi / N$ , the field solution will only contain harmonics where  $\exp(i \cdot m \cdot \varphi) = 1$  after a full period, or in other words, modes whose angular dependence is compatible with the symmetry. In a perfectly  $N$ -fold periodic structure, one finds that

effectively  $m$  is congruent modulo  $N$  (modes separated by  $N$  in index are physically degenerate). For example, if  $N = 4$  ( $90^\circ$  symmetry), an  $m = 1$  solution corresponds to a field that repeats every  $360^\circ$  with a quarter-period phase shift between each  $90^\circ$  sector; an  $m = 4$  solution would look identical to an  $m = 0$  solution in a 4-fold structure (since  $\exp(i \cdot 4 \cdot \varphi) = \exp(i \cdot 2\pi) = 1$ ). In practice, one would simulate  $m = 0, 1, 2, \dots, N-1$  as the distinct cases for an  $N$ -fold symmetric device. However, the FDTD algorithm itself does not require us to restrict  $m$ —we simply choose the  $m$  corresponding to the physical scenario we want to simulate (often  $m = 0$  for symmetric excitations, or  $m = 1$  for a lowest-order asymmetric perturbation, etc.). The linear combination of these solutions can reconstruct any arbitrary excitation pattern around the axis.

In summary, by using the above Fourier decomposition, we have reduced the problem to solving Maxwell's equations in two dimensions ( $r$  and  $z$ ) for each desired harmonic  $m$ . Next, we discretize these equations using FDTD's central-difference scheme. Each harmonic  $m$  will be simulated on a cylindrical grid that spans only the radial and axial dimensions (plus a limited angular domain if needed), with  $m$  entering as a parameter in the update equations.

### **3.4 Yee Grid and Discretization in Cylindrical Coordinates**

We adopt a Yee-like grid in cylindrical coordinates to discretize the space. The grid is defined by sampling the  $r$ ,  $\theta$ , and  $z$  axes at discrete intervals  $\delta r$ ,  $\delta\theta$ , and  $\delta z$ , respectively [2]. A straightforward approach is to use a *staggered grid* in which the electric field



components  $E_r$ ,  $E_\theta$ ,  $E_z$  are defined at interleaved spatial locations relative to the magnetic field components  $H_r$ ,  $H_\theta$ ,  $H_z$ . This staggered arrangement is crucial for satisfying the discrete curl equations with second-order accuracy in space and time [1].

**Figure 1:** Schematic of a Yee unit cell in cylindrical coordinates, showing the staggered curl-conforming arrangement of the electric field components (red arrows, labeled as  $u$  components) and magnetic field components (blue arrows, labeled as  $v$  components) within the cell. The cell is bounded by radii  $r = i \cdot \Delta r$  and  $r = (i+1) \cdot \Delta r$ , azimuthal angles  $\theta = j \cdot \Delta\theta$  and  $\theta = (j+1) \cdot \Delta\theta$ , and axial positions  $z = k \cdot \Delta z$  and  $z = (k+1) \cdot \Delta z$ .

The electric field  $E_r$  (red arrow labeled  $u_r(i+\frac{1}{2}, j, k)$ ) is defined at the cell face normal to  $\hat{r}$  (at the midpoint between  $r_i$  and  $r_{i+1}$ , here denoted  $i+\frac{1}{2}$ , and located on the boundary of the cell in  $\theta$  and  $z$ ). Similarly,  $E_\theta$  ( $u_\theta(i, j+\frac{1}{2}, k)$  in red) lies on the face

normal to  $\hat{\theta}$  (midway between  $\theta_j$  and  $\theta_{j+1}$ ), and  $E_z$  ( $u_z(i, j, k+1/2)$ ) on the face normal to  $\hat{z}$  (midway between  $z_k$  and  $z_{k+1}$ ).

The magnetic field components (blue) are staggered such that each  $H$  lies on the edges of the cell, halfway between two  $E$  components. For example,  $H_r$  ( $v_r(i, j+1/2, k+1/2)$ , blue arrow) is located at the edge intersection of the  $E_\theta$  and  $E_z$  faces,  $H_\theta$  ( $v_\theta(i+1/2, j, k+1/2)$ ) at the intersection of  $E_r$  and  $E_z$  faces, and  $H_z$  ( $v_z(i+1/2, j+1/2, k)$ ) at the intersection of  $E_r$  and  $E_\theta$  faces. This spatial staggering ensures that each discrete curl equation (linking differences of  $E$  fields around a loop to an  $H$  field, and vice versa) is satisfied in a locally circulation-preserving manner, analogous to the Cartesian Yee grid.

The grid shown in Fig. 1 effectively “unrolls” the cylindrical coordinates into a mesh: index  $i$  corresponds to radial position  $r_i = i \cdot \Delta r$  (with  $i = 0$  at  $r = 0$  on the axis), index  $j$  corresponds to angle  $\theta_j = j \cdot \Delta \theta$ , and index  $k$  corresponds to height  $z_k = k \cdot \Delta z$ .

If we are exploiting  $N$ -fold symmetry, we may restrict  $\theta$  to span only one sector of size  $\varphi = 2\pi / N$ . In that case,  $j$  runs from 0 to  $N_\theta - 1$  such that  $N_\theta \cdot \Delta \theta = \varphi$ . We then apply periodic boundary conditions between  $\theta = 0$  and  $\theta = \varphi$  with the appropriate phase shift for the harmonic  $m$  (as will be described shortly).

If we are including the full 0 to  $2\pi$  domain (for instance, to simulate multiple modes at once or to verify results), then  $j$  runs over the full circumference with  $\Delta \theta$  small enough to resolve the fields. In either scenario, the  $\theta$ -domain is treated as cyclic.

In the time dimension, we likewise employ Yee’s leapfrog scheme: the electric field components are updated at time steps  $n \cdot \Delta t$  (integer  $n$ ), and the magnetic field components at half-integer steps  $(n+1/2) \cdot \Delta t$ . This temporal staggering means  $E$  and  $H$  fields are interleaved in time, which matches the structure of Maxwell’s equations (the time derivative of  $H$  is given by the curl of  $E$ , and vice versa). The goal of discretization is to transform the continuous  $m$ -harmonic Maxwell equations derived in the previous section into finite-difference update formulas connecting  $E$  and  $H$  at neighboring grid points and time steps. We must take care to handle the  $1/r$  metric terms and the  $i \cdot m$  coupling terms properly in the discrete context.

### 3.5 Discrete Update Equations for Interior Points

At grid points away from the axis ( $i > 0$ ) and not on the  $\theta$ -boundary, we can directly apply central-difference approximations to the spatial derivatives in the harmonic Maxwell equations. To illustrate, consider the r-component curl equation for  $H_r^{(m)}$  (from Faraday's law):

$$\frac{im}{r} E_z^{(m)} - \frac{\partial E_\theta^{(m)}}{\partial z} = -\mu \frac{\partial H_r^{(m)}}{\partial t}$$

In the discrete grid, this equation is enforced at the location of  $H_r$ , which in Fig. 1 is at  $(i, j + \frac{1}{2}, k + \frac{1}{2})$  (centered in a cell edge). The term  $\partial E_\theta / \partial z$  is approximated by a finite difference along  $z$  between the  $E_\theta$  values above and below  $H_r$ . Meanwhile,  $(1/r) \cdot E_z$  at a half-integer  $r$  position is taken as  $E_z$  divided by the local radius. In practice,  $r$  in such terms is taken as the value at the location of  $H_r$  (for instance,  $r_i = i \cdot \Delta r$  if  $H_r$  lies exactly at radius  $r_i$ ; note that  $H_r$  at  $(i, j + \frac{1}{2}, k + \frac{1}{2})$  is actually on the boundary of the cell at  $r_i$ ). Similarly, the  $i \cdot m$  factor is treated as a constant multiplier. The time derivative  $\partial H_r / \partial t$  is discretized as a forward difference between time steps  $(n - \frac{1}{2})$  and  $(n + \frac{1}{2})$  for  $H_r$ . The resulting update formula for  $H_r$  (at index  $i, j + \frac{1}{2}, k + \frac{1}{2}$ ) is:

$$H_\theta \Big|_{i+\frac{1}{2}, j, k+\frac{1}{2}}^{n+\frac{1}{2}} = H_\theta \Big|_{i+\frac{1}{2}, j, k+\frac{1}{2}}^{n-\frac{1}{2}} - \frac{\Delta t}{\mu} \left[ \frac{E_r \Big|_{i+1, j, k}^n - E_r \Big|_{i, j, k}^n}{\Delta z} - \frac{E_z \Big|_{i+\frac{1}{2}, j, k+\frac{1}{2}}^n - E_z \Big|_{i+\frac{1}{2}, j-1, k+\frac{1}{2}}^n}{\Delta r} \right],$$

$$H_z \Big|_{i+\frac{1}{2}, j+\frac{1}{2}, k}^{n+\frac{1}{2}} = H_z \Big|_{i+\frac{1}{2}, j+\frac{1}{2}, k}^{n-\frac{1}{2}} - \frac{\Delta t}{\mu} \left[ \frac{1}{r_{i+\frac{1}{2}}} \frac{(r E_\theta) \Big|_{i+1, j+\frac{1}{2}, k}^n - (r E_\theta) \Big|_{i, j+\frac{1}{2}, k}^n}{\Delta r} - \frac{im}{r_{i+\frac{1}{2}}} E_r \Big|_{i, j+\frac{1}{2}, k}^n \right]$$

(For brevity, we have not written out the material parameters in these equations or the indices explicitly; the pattern is similar: difference of E components in appropriate directions, plus the  $i \cdot m$  term where a  $\theta$ -derivative was present.) The  $H_\theta$  update uses

differences of  $E_r$  in  $z$  and  $E_z$  in  $r$ , while the  $H_z$  update uses differences of  $r \cdot E_\theta$  in  $r$  (note the factor  $r$  inside the difference to properly approximate  $\partial(r \cdot E_\theta) / \partial r$ ) and an  $i \cdot m \cdot E_r$  term divided by  $r$ . All spatial differences are centered around the  $H$  field location, consistent with the staggering shown in Fig. 1. The use of  $r_{\{i + 1/2\}}$  (the radius at the midpoint between  $r_i$  and  $r_{\{i + 1\}}$ ) in the  $H_z$  update highlights that we evaluate  $1 / r$  at the proper location of  $H_z$ .

Likewise, the Ampère-law equations are discretized to update the  $E$ -field components at the integer time steps. For example, consider the  $z$ -component equation for  $E_z^{(m)}$ :

$$\frac{1}{r} \frac{\partial(r H_\theta^{(m)})}{\partial r} - \frac{i m}{r} H_r^{(m)} = \epsilon \frac{\partial E_z^{(m)}}{\partial t} + J_z^{(m)}$$

The corresponding update for  $E_z$  at location  $(i, j, k + 1/2)$  (center of the  $E_z$  face in Fig. 1) will involve the difference of  $H_\theta$  in the  $r$ -direction and an  $i \cdot m$  term involving  $H_r$  at that location. In finite-difference form (omitting the source  $J_z$  term for

$$E_z \Big|_{i, j, k + \frac{1}{2}}^{n+1} = E_z \Big|_{i, j, k + \frac{1}{2}}^n + \frac{\Delta t}{\epsilon_{i, j, k + \frac{1}{2}}} \left[ \frac{(r H_\theta)_{i, j, k + \frac{1}{2}}^{n + \frac{1}{2}} - (r H_\theta)_{i-1, j, k + \frac{1}{2}}^{n + \frac{1}{2}}}{r_i \Delta r} - \frac{i m}{r_i} H_r \Big|_{i, j + \frac{1}{2}, k + \frac{1}{2}}^{n + \frac{1}{2}} \right]$$

simplicity), one obtains:

Here  $(r \cdot H_\theta)_{\{i, j, k + 1/2\}}$  denotes  $r_i \cdot H_\theta(i + 1/2, j, k + 1/2)$  evaluated at the appropriate midpoint radius. Similar expressions can be written for  $E_r$  and  $E_\theta$ . The  $E_r$  update will include differences of  $H_z$  in  $z$  and  $H_\theta$  in  $\theta$  (with an  $i \cdot m$  term if a  $\theta$  difference is replaced by it), and the  $E_\theta$  update includes differences of  $H_r$  in  $z$  and  $H_z$  in  $r$  (no  $i \cdot m$  term appears in the  $E_\theta$  update because its equation had no  $\partial / \partial \theta$  term in the original form). We won't clutter the text with every update equation explicitly; the key point is that each spatial derivative from the continuous equations is approximated by a

finite difference of the corresponding field components on the Yee grid, and each  $i \cdot m / r$  term is treated as a simple multiplication by  $i \cdot m$  at the grid point (with  $r$  taken as the local radius). These updates are applied for each harmonic  $m$  being simulated. In implementation, one typically handles the  $i \cdot m$  terms by splitting the fields into real and imaginary parts (since  $i \cdot m$  introduces a  $90^\circ$  phase difference between those parts), or by using complex arithmetic directly if supported. Either way, the FDTD scheme updates the complex field amplitudes  $E^{(m)}$  and  $H^{(m)}$  in time.

The stability condition (Courant–Friedrichs–Lewy condition) for the cylindrical FDTD is similar to that in Cartesian, but with the smallest cell size determined by the minimum of  $\Delta r$ ,  $r \cdot \Delta\theta$ , and  $\Delta z$ . Typically, the most restrictive condition may occur either in the  $z$  direction or at the smallest radius for the  $\theta$  direction. If the grid includes the axis ( $r = 0$ ), note that  $r \cdot \Delta\theta$  goes to 0 at  $i = 0$ ; however, because the angular cell size shrinks to zero at the axis, the fields there are handled by a special update (discussed next) rather than a standard Yee cell, so the effective cell size to consider is the smallest non-zero radius cell. In practice, one should ensure  $\Delta t$  satisfies:

$$\Delta t < 1 / c \cdot \sqrt{1 / (\Delta r)^2 + 1 / (r_{\min} \cdot \Delta\theta)^2 + 1 / (\Delta z)^2}^{-1}$$

for stability. As long as this is met, the cylindrical FDTD scheme will be stable and second-order accurate in space and time, just like its Cartesian counterpart.

### 3.6 Special Treatment at the Axis ( $r = 0$ )

The  $r=0$  line (the  $z$ -axis) is a singular line for cylindrical coordinates. In the continuous equations, as  $r$  to 0, terms become problematic. Physically, however, well-behaved fields must remain finite at the axis, and the Maxwell equations themselves remain well-defined in the limit approaching the axis. To handle the axis in FDTD, we must derive proper boundary conditions or update equations for the field components at  $r=0$  that are consistent with Maxwell's equations. The Yee grid as defined puts certain field components exactly on the axis. Referring to Fig. 1, the radial index  $i=0$  corresponds to the axis. Typically,  $E_r$  and  $E_z$  components can exist on the axis (pointing radially

outward and axially along the axis, respectively), whereas an  $E_\theta$  component on the axis is not uniquely defined  $\theta$  direction is undefined at  $r=0$ ). In our grid staggering,  $E_\theta$  for  $i=0$  lies along the inner boundary of the first cell, effectively on the axis. Likewise,  $H_r$  for  $i=0$  sits on the axis. We need to update those  $E_\theta(0, j+1/2, k)$  and  $H_r(0, j+1/2, k+1/2)$  values with modified formulas that do not explicitly use  $1/r$  terms.

The strategy is to enforce Maxwell's equations at the axis by using their integral form around a small loop that straddles the axis. By symmetry, as  $r \rightarrow 0$ , the fields often approach either zero or a finite limit. One known result is that the tangential electric field  $E_\theta$  on the axis should go to zero if the fields are finite (because a finite magnetic flux through an infinitesimal loop around the axis implies  $E_\theta$  must vanish). Similarly,  $H_r$  on the axis should be zero in many cases. However, rather than simply imposing zero, a more accurate approach is to derive the discrete update from Maxwell's equations to capture any subtle effects (especially for  $m \neq 0$  modes, where fields might not be strictly zero on the axis). We consider two loops:

- Loop around  $E_\theta$  at the axis: We form a rectangular integration loop in the  $r$ - $z$  plane enclosing the  $E_\theta(0, j + 1/2, k)$  edge (which lies along the  $z$ -axis). This loop is shown in Figure 2. It extends a half-cell in the  $+r$  direction (to  $r = \Delta r / 2$ ) and a half-cell above and below in  $z$  (to  $z = \pm \Delta z / 2$  from that  $E_\theta$ ). Applying Faraday's law in integral form,  $\oint \mathbf{E} \cdot d\mathbf{l} = -\mu \cdot d/dt \iint \mathbf{B} \cdot d\mathbf{A}$ , around this loop will involve the  $E_\theta$  on the axis and the neighboring  $E_\theta$  (if any) at the radial half-step, as well as the  $E_r$  segments (which contribute zero on the axis since  $E_r$  at  $r = 0$  has no length) and the induced voltage from the changing magnetic flux through the loop. The magnetic flux comes from the  $H_r$  (radial) component passing through the small area and the difference of  $H_r$  across  $\theta$  boundaries (because effectively, at  $r = 0$ , the loop spans all  $\theta$  directions in a collapsed way). Performing this analysis yields an update formula for  $E_\theta(0, j + 1/2, k)$  that does not contain a singular term. In fact, one finds that  $E_\theta$  on the axis is updated by the difference of the azimuthal magnetic field  $H_r$  just off the axis on either side of the wedge. For a full  $2\pi$  domain, this effectively gives:

$$\mathbf{E}_\theta(\mathbf{0})^{n+1} = \mathbf{E}_\theta(\mathbf{0})^n + (\Delta t / \varepsilon \cdot \Delta r) \cdot [2 \cdot \mathbf{H}_z |_{\{i = \frac{1}{2}, j + \frac{1}{2}, k\}^{n + \frac{1}{2}}}]$$

(the factor of 2 arises because the magnetic field on both sides of the axis contributes identically). In the case of a single wedge of size  $\varphi$ , the two “sides” correspond to the periodic boundaries at  $\theta = 0$  and  $\theta = \varphi$ , which for mode  $m$  have a phase difference. The algorithm accounts for that phase by effectively using  $H_z$  from  $\theta = 0$  and  $\theta = \varphi$  with the  $\exp(i \cdot m \cdot \varphi)$  factor connecting them.

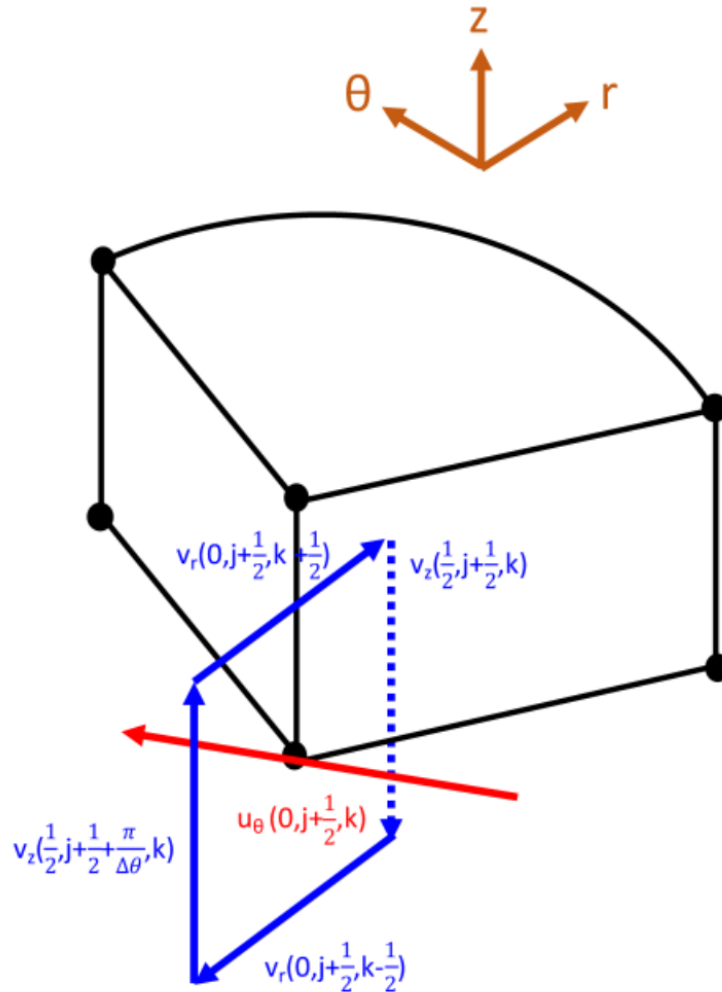


Fig. 2.  $E_\theta$  at  $r = 0$

Figure 2: Illustration of the integration loop used to update  $E_\theta$  at the axis ( $r = 0$ ). The red arrow indicates the electric field  $E_\theta(0, j + \frac{1}{2}, k)$  located on the z-axis (azimuthal component). The blue arrows show the magnetic field components around the loop: in this case,  $H_r$  (blue arrow along the radial direction) and  $H_z$  (blue dashed arrow along z) in the first radial cell. The loop (black rectangle) goes around the  $E_\theta$  segment on the axis, extending out to the midpoint of the first radial cell. By applying Faraday's law

$\oint \mathbf{E} \cdot d\mathbf{l} = -\mu \cdot d\Phi_B / dt$  to this small rectangle, one finds that  $E_\theta(0, j + \frac{1}{2}, k)$  is driven by the time-changing flux of  $H_r$  through the loop. Because the path encloses the axis, the contribution from the azimuthal variation of  $H_r$  on the two sides (at  $\theta = 0$  and  $\theta = \varphi$ ) must be included. The result (in the absence of currents) is that  $E_\theta$  on the axis is updated by an average of  $H_r$  from either side, which for a symmetric mode leads to  $E_\theta(0) \rightarrow 0$  as  $\Delta r \rightarrow 0$ . In practice, this approach avoids any division by  $r$  at the axis.

Loop around  $H_r$  at the axis: Similarly, we consider an integration loop for Ampère's law encircling the  $H_r(0, j + \frac{1}{2}, k + \frac{1}{2})$  segment that lies on the axis (see Figure 3). This loop is in the  $r$ - $z$  plane and encloses the  $H_r$  on the axis. Ampère's circuital law in integral form  $\oint \mathbf{H} \cdot d\mathbf{l} = \varepsilon \cdot d/dt \iint \mathbf{E} \cdot d\mathbf{A} + \iint \mathbf{J} \cdot d\mathbf{A}$  is applied. The loop picks up contributions from the electric field  $E_\theta$  in the first cell around the axis and any  $E_z$  along the radial arm (which vanishes exactly on the axis because the radial line has no circumference). The resulting update condition for  $H_r(0, j + \frac{1}{2}, k + \frac{1}{2})$  involves the difference of  $E_\theta$  values just off the axis in the  $\theta$  direction. In essence, the change in  $H_r$  at the axis is proportional to the difference  $E_\theta$  between  $\theta = 0$  and  $\theta = \varphi$  (for a wedge) or twice the  $E_\theta$  in the first cell (for full  $2\pi$ ). For a mode with  $\exp(i \cdot m \cdot \theta)$  variation,  $E_\theta$  at the two sides of the wedge differ by a factor  $\exp(i \cdot m \cdot \varphi)$ , so the formula naturally incorporates  $2 \cdot \sin(m \cdot \varphi / 2)$  factors if written in real terms, ensuring stability and correctness for all  $m$ .

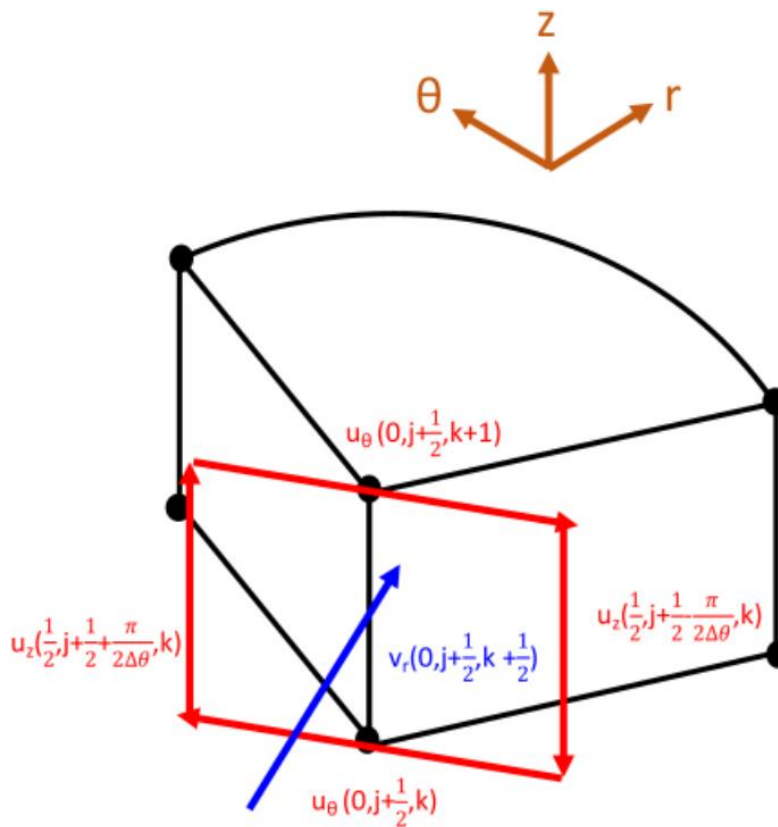


Fig. 3.  $H_r$  at  $r = 0$

Figure 3: Integration loop for updating  $H_r$  at the axis. The blue arrow denotes the magnetic field  $H_r(0, j + \frac{1}{2}, k + \frac{1}{2})$  located on the  $r = 0$  line. The red arrows indicate the electric field components contributing to Ampère's law around this loop:  $E_\theta$  (red horizontal arrows) in the first cell adjacent to the axis above and below, and  $E_z$  (red vertical arrows) at the midpoints between  $\theta = 0$  and  $\theta = \varphi$  (which effectively cancel out on the axis). By integrating  $\oint \mathbf{H} \cdot d\mathbf{l} = \varepsilon \cdot d\Phi_D / dt$  around this loop, one finds that the change in  $H_r$  on the axis is driven by the difference in  $E_\theta$  across the angular interval of the wedge. In a full  $360^\circ$  simulation, this reduces to  $H_r(0)$  being driven by the azimuthal circulation of  $E_\theta$  around the axis (which is typically zero for symmetric fields). In a single-wedge simulation for harmonic  $m$ , it effectively includes the phase difference of  $E_\theta$  between the wedge boundaries. The result ensures  $H_r$  at  $r = 0$  evolves without any singular behavior and usually remains zero for non-singular excitations.

By handling the axis with such integral considerations, the FDTD scheme remains stable and accurate at  $r = 0$ . In effect, the  $1 / r$  terms that would appear in the discrete equations are circumvented by using symmetry: the field at the axis is updated by considering the entire circumference that shrinks to a point. For most practical cases, these special-axis updates enforce that  $E_\theta(0) = 0$  and  $H_r(0) = 0$  throughout the simulation (assuming no singular line source on the axis). If there is a physical line current on the axis (e.g., a  $J_z$  along the axis), it can be incorporated via the Ampère's law loop as an additional term in the  $H_r$  update (similar to how a filament current would appear in Ampère's circuital law). The formulation is general enough to handle that if needed.

In summary, the discrete axial symmetry FDTD algorithm uses the standard Yee finite-difference scheme for all interior grid points ( $i \geq 1$ ), and employs specialized update formulas derived from integral Maxwell equations for the on-axis field points. This guarantees that the simulation handles the  $r = 0$  singularity consistently with Maxwell's laws and avoids any numerical instability or inconsistency at the axis.

### 3.7 Summary and considerations

In this chapter, we derived a comprehensive FDTD formulation in cylindrical coordinates, incorporating discrete rotational symmetry through Fourier decomposition in the azimuthal direction. By expanding the fields in angular harmonics  $\exp(i \cdot m \cdot \theta)$ , we reduced the 3D problem to a set of decoupled 2D problems (in  $r$  and  $z$ ) for each harmonic index  $m$  [3]. We presented the full set of Maxwell's equations in cylindrical form and carried out the derivation to discrete Yee-grid update equations. Key aspects of this formulation include:

- Cylindrical Yee Grid: We set up a staggered Yee grid in  $(r, \theta, z)$  space (Fig. 1) where  $E$  and  $H$  field components are interleaved in space. This grid preserves the curl relationships in discrete form, analogous to the Cartesian Yee grid [1]. The presence of the metric factor  $1 / r$  was handled by evaluating  $r$  at the proper locations and including it in the finite differences (e.g., using differences of  $r \cdot E_\theta$ ).
- Fourier (Azimuthal) Mode Decomposition: Using the assumption of  $\exp(i \cdot m \cdot \theta)$  dependence, we simplified all  $\theta$ -derivatives to algebraic  $i \cdot m$  terms. This allowed us to limit the computational domain in the  $\theta$  direction to a single period (or even just one cell in many cases) and to introduce  $m$  as a parameter in the update equations rather than having to step through many angular cells. Different  $m$  values can be simulated independently. If the physical structure has  $N$ -fold symmetry, one typically needs to simulate only  $0 \leq m < N$  to capture all unique field patterns. The case  $m = 0$  reproduces an axisymmetric FDTD scheme [2], while  $m \neq 0$  gives higher-order angular variations.

- Incorporation of  $i\omega$  Terms: The effect of the  $i\omega$  terms in the equations is to introduce coupling between certain field components (e.g.,  $E_r$  and  $H_z$ ) that would not exist in a purely axisymmetric case. In the discrete updates, these appear as additional source-like terms. We carefully included these in the finite-difference formulas (see examples of  $H_r$  and  $H_z$  updates above). The use of complex fields (or dual real fields for cosine and sine parts) is necessary to represent these  $90^\circ$  phase-shifted interactions. This does not pose difficulty and is analogous to frequency-domain FDTD or k-space FDTD methods [3].
- Axis Singularity Treatment: We devoted special attention to the  $r = 0$  axis. By using integral forms of Maxwell's equations on small loops encircling the axis (illustrated in Figs. 2 and 3), we derived proper update rules for the on-axis field components ( $E_\theta$  and  $H_r$  primarily). This removes the coordinate singularity and ensures that the fields at the axis remain finite and consistent. The derived conditions typically enforce that tangential electric field and radial magnetic field go to zero on the axis (for well-behaved solutions), as expected physically. Our approach is consistent with other treatments in literature that employ series expansions or rotated coordinate techniques to handle the axis [2][4]. Importantly, this allows the FDTD grid to include  $r = 0$  without instability.
- Boundary Conditions and PML: Although not the focus of this chapter, it is worth noting that perfectly matched layer (PML) absorbing boundaries can be implemented in cylindrical coordinates as well, with some modifications to account for the  $1/r$  term in the Maxwell equations [3]. In our discrete equations above, we could incorporate PML by stretching coordinates or adding loss terms (conductivity  $\sigma_r(r)$ ,  $\sigma_\theta(\theta)$ ,  $\sigma_z(z)$ ) in each equation. The Fourier decomposition approach is fully compatible with PML, as one can add  $+\sigma \cdot E$  terms to Ampère's law and  $-\sigma^* \cdot H$  to Faraday's law for the lossy media

simulation (the uploaded reference, for instance, included  $\sigma_r(r)$  and  $\sigma_z(z)$  in equations (18)-(19), etc., which correspond to a stretched-coordinate PML) [3]. When designing a cylindrical FDTD simulation, one would typically put a PML at the outer radial boundary  $r = r_{\text{max}}$  to absorb outgoing waves, and standard absorbing or periodic boundaries in  $z$  as needed. The  $\theta$  direction is inherently handled by the symmetry (periodic if using a wedge).

Overall, the FDTD method in cylindrical coordinates with discrete axial symmetry retains the advantages of the standard Yee FDTD (second-order accuracy, explicit time stepping, straightforward implementation) while significantly reducing the computational cost for problems with rotational symmetry. Instead of needing a full 3D grid, one can simulate a single wedge or even just use the analytic  $m$  terms to avoid any explicit  $\theta$  grid. The memory and CPU time scale roughly as  $O(N_r \times N_z)$  for each harmonic (as opposed to  $O(N_r \times N_\theta \times N_z)$  in a full 3D cylindrical mesh), and often only a few harmonics are needed. This is a dramatic improvement for large  $N$  or continuous rotational symmetry. The method is thus highly efficient for analyzing devices like coaxial cables, circular particle accelerators, periodic antennas, and metasurfaces that have angular periodicity [3].

To conclude, we provide a brief outline of the procedure one would follow to use this cylindrical FDTD method in practice:

1. Harmonic Selection: Determine the symmetry of the problem and choose the azimuthal harmonic index  $m$  (or indices) that need to be simulated. For example, use  $m = 0$  for an axially symmetric excitation, or a set of  $m$  values if the excitation is asymmetric or if multiple modes are of interest.
2. Grid Setup: Define a radial grid from  $r = 0$  to the outer radius of interest, and a vertical grid in  $z$  covering the structure. Choose  $\Delta r$ ,  $\Delta\theta$ ,  $\Delta z$  such that the cell aspect ratios are reasonable and the CFL stability condition is satisfied. If using a single-sector domain, set  $\Delta\theta = \phi$  (the sector angle) or a divisor of  $\phi$  if multiple

cells in  $\theta$  are used for accuracy. Typically, one cell in  $\theta$  (the wedge spanning  $\varphi$  with periodic boundaries) is sufficient for a given  $m$  since internal  $\theta$  variation is captured by the  $\exp(i \cdot m \cdot \theta)$  factor.

3. Initial Conditions: Initialize the E and H fields. Often one starts with all fields zero and introduces a source (current or impressed field) at a boundary or a certain location. The current source should likewise be decomposed in  $m$  harmonics (e.g., a z-directed current on the axis would only have an  $m = 0$  component, whereas an asymmetrically placed source would produce multiple  $m$  components).
4. Time-Stepping Loop: Iterate over time using the leapfrog update. On each half-step, update all H field components using the discrete Faraday law equations (with  $i \cdot m$  terms included). On each full step, update all E field components using Ampère's law equations (with  $i \cdot m$  terms and any source J terms). Apply the special axis update formulas for  $E_\theta(0)$  and  $H_r(0)$  after each normal update, to correct those values (this ensures the axis respects Maxwell's constraints). Also enforce the phase relation at the  $\theta$  boundaries: in a single-wedge simulation, the fields at  $\theta = 0$  and  $\theta = \varphi$  are not independent but related by  $E(r, \varphi, z, t) = E(r, 0, z, t) \cdot \exp(i \cdot m \cdot \varphi)$  (and similarly for H). In practice, this is implemented by copying the field from one boundary to the other with the phase factor. Because we have stored complex fields, this is straightforward multiplication by  $\exp(i \cdot m \cdot \varphi)$ . If one uses a full  $0$  to  $2\pi$  domain, then standard periodic boundary conditions (with continuity) are used and the  $i \cdot m$  factor is inherently satisfied by the continuity of the field.
5. Outputs: At each time step (or after simulation), reconstruct any desired physical fields. If only one  $m$  was simulated and that was the only excited harmonic, the fields (complex) can be converted to real time-domain fields by taking the real part:  $E = \text{Re}\{E^{(m)} \cdot \exp(i \cdot m \cdot \theta)\}$ . If multiple  $m$  components were simulated separately, their results can be summed (superposed) to obtain the total field distribution.

6. Verification: It is good practice to verify that energy is conserved (or decays as expected if PMLs are used) and that  $E_\theta$  and  $H_r$  remain small near the axis for a well-behaved solution (indicating the singularity is handled correctly). One may also test the  $m = 0$  case against a known axisymmetric analytic solution or a radial-transmission line problem to ensure the code is correct.

In conclusion, the cylindrical-coordinate FDTD with discrete axial symmetry is a powerful extension of the standard FDTD technique. It leverages the symmetries in the problem to reduce computational overhead while maintaining accuracy and stability. The derivation provided in this chapter serves as a foundation for implementing such a solver. With these equations and methods, a researcher can simulate electromagnetic phenomena in complex cylindrical and periodic structures with far fewer resources than a brute-force 3D simulation, without sacrificing rigor or generality. This makes it an indispensable tool in the arsenal of computational electromagnetics for problems where circular geometries and symmetries play a role.

# Chapter 4: GPU Computing

## 4.1 What is a GPU? A Deep Dive into Graphics Processing Units

A Graphics Processing Unit (GPU) is a specialized microprocessor designed to perform computations related to rendering images and visual data for output to a display. While this was its original role, the modern GPU has evolved into a massively parallel processor capable of accelerating a wide variety of non-graphics workloads, particularly those involving large-scale, data-parallel numerical computations.

### 4.1.1 Historical Evolution of GPUs

The history of GPUs dates to the early 1990s, when graphics hardware was used primarily to offload fixed-function rendering tasks such as rasterization, texture mapping, and depth buffering. Early GPUs like the NVIDIA NV1 (1995) and 3dfx Voodoo (1996) were designed as co-processors to handle geometry transformations and pixel shading, freeing up the CPU for game logic and input/output.

The first major leap toward programmable graphics hardware came with the introduction of shader models in the early 2000s. With DirectX 8 (2000) and OpenGL 2.0 (2004), developers gained access to vertex and fragment shaders—small programs that allowed for custom processing of rendering pipelines. This change laid the groundwork for what would become GPGPU computing, as researchers realized that the same hardware capable of manipulating pixels could also be used for general-purpose floating-point computations.

The true transformation into a compute platform came in 2006, when NVIDIA introduced the CUDA (Compute Unified Device Architecture) programming model. CUDA provided an API for developers to write C-like code that could be executed directly on the GPU's streaming multiprocessors. Shortly after, OpenCL offered a vendor-neutral alternative, and the age of GPU-accelerated scientific computing began in

earnest [5]. Over the next decade, GPUs would become central to fields such as deep learning, medical imaging, fluid dynamics, cryptography, and computational electromagnetics [6].

### 4.1.2 GPU Hardware Architecture

The architecture of a GPU is fundamentally different from that of a CPU. A typical CPU is designed for low-latency execution of complex control-intensive operations. It includes a few powerful cores equipped with large caches and sophisticated branch predictors to handle diverse instruction types [7].

By contrast, a GPU contains hundreds to thousands of simpler cores, each capable of performing floating-point arithmetic but with limited branching capabilities. These cores are grouped into Streaming Multiprocessors (SMs) or Compute Units, each of which manages multiple threads simultaneously.

In NVIDIA's CUDA architecture, the smallest unit of thread execution is a warp, consisting of 32 threads that execute instructions in a Single Instruction, Multiple Thread (SIMT) fashion. All threads in a warp follow the same instruction path. If a conditional branch causes threads within a warp to diverge, the execution becomes serialized—known as warp divergence, which can reduce efficiency.

Each Streaming Multiprocessor contains its own register file, used for storing thread-local variables, shared memory which is a fast manually-managed memory cache for inter-thread communication within a thread block and instruction and scheduling logic for dispatching warps to available execution units.

The GPU's memory hierarchy includes: Global memory, large but high-latency memory accessible to all threads, Shared memory, a low-latency memory shared within a block, L1 and L2 cache, hardware-manages memories which vary by architecture, Constant memory for storing the read-only parameters, and lastly texture and surface memory, optimized for graphics workloads, but also usable in GPGPU [8]. Modern GPUs also include tensor cores, specialized hardware units for fast matrix multiplication,

originally designed for AI workloads but increasingly applicable to scientific simulations involving linear algebra.

### 4.1.3 Instruction Scheduling and Threading Model

A defining feature of GPU architecture is its ability to support massive concurrency. A single kernel launch may involve hundreds of thousands of threads organized into blocks and grids. The GPU's warp scheduler issues instructions from ready warps in round-robin or priority order, ensuring that computational units remain fully occupied.

This model tolerates latency through over-subscription: while one warp waits for data to arrive from global memory (a process that takes hundreds of clock cycles), other warps are scheduled for execution. In effect, GPUs rely on thread-level parallelism to hide latency, rather than cache hierarchies [8].

For example, an NVIDIA A100 GPU contains 108 SMs, each capable of supporting up to 64 resident warps, leading to tens of thousands of concurrent threads [8]. This model works exceptionally well for applications that can be expressed as data-parallel loops, avoid heavy branching or recursion and ones that use structured memory access patterns.

### 4.1.4 Performance Metrics and Theoretical Limits

The performance of a GPU is typically measured in floating-point operations per second (FLOPS) and memory bandwidth. The peak performance of an NVIDIA A100, for instance, is 19.5 TFLOPS for FP32 meaning we have single precision, 9.7 TFLOPS for FP64 meaning we have double precision, up to 312 TFLOPS for matrix multiply-accumulate operations using mixed precision and 1.68 TB/s for memory bandwidth.

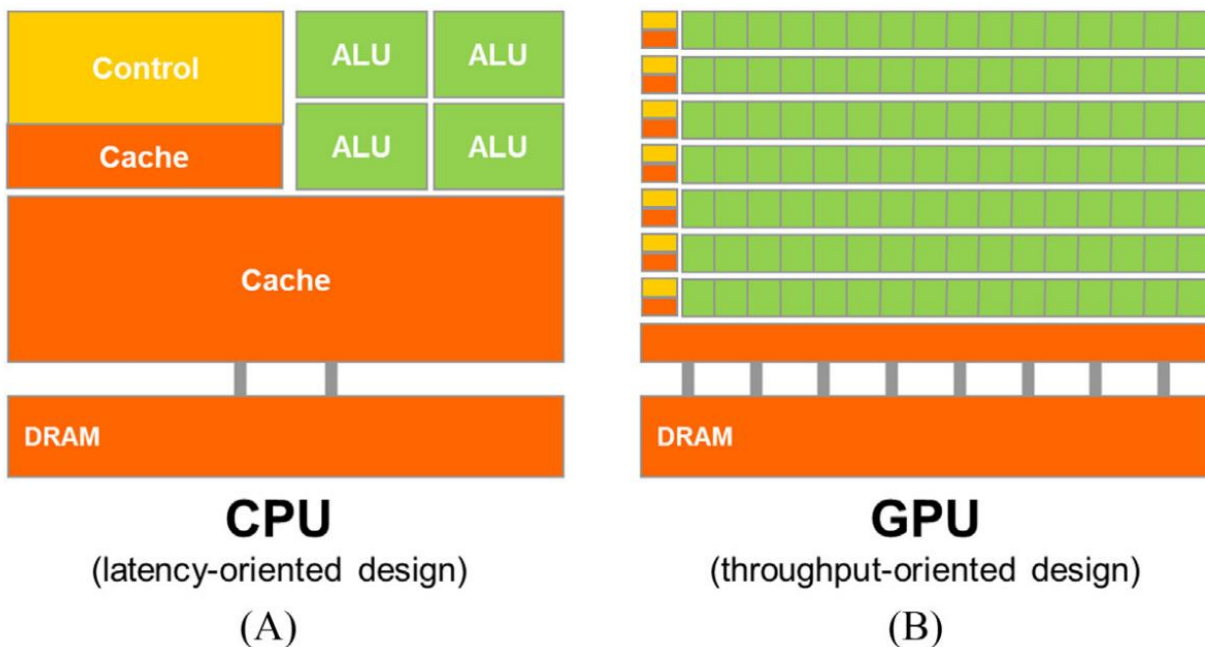
These figures are orders of magnitude higher than what a high-end CPU can deliver, particularly in memory-bound or arithmetic-bound workloads. However,

achieving peak performance requires careful programming: memory access must be coalesced, shared memory must be used effectively, and thread divergence must be minimized.

In conclusion, the GPU is a highly specialized, massively parallel compute engine. It is not a replacement for the CPU, but a complement—best employed where data parallelism and high throughput are critical. In the context of this thesis, the GPU’s architecture is not just a convenience but a computational necessity: without GPU acceleration, the simulation and optimization of millimeter-scale freeform metalenses would be infeasible.

## 4.2 GPU vs CPU in Scientific Computing

A modern CPU and GPU differ substantially in design philosophy. A CPU typically has a handful of very fast cores optimized for serial processing and low-latency access to memory, with large caches and sophisticated control logic for branch prediction, out-of-order execution, etc. In contrast, a GPU devotes most of its silicon to a vast array of simpler cores and arithmetic units, relying on massive parallelism to achieve high throughput.



As a result, a CPU excels at complex, branching workloads with moderate parallelism, whereas a GPU shines when presented with a highly parallel, uniform workload (like applying the same operations across a large grid or vector of data).

Illustration of CPU vs GPU architecture: A CPU (left) dedicates significant chip area to caches and control logic (gold) for a few cores (green), whereas a GPU (right) allocates most resources to a large number of compute cores (green) with relatively less cache and control per core. This design allows GPUs to execute many parallel threads but makes each thread less individually powerful than a CPU core.

In numerical scientific computing – e.g. solving partial differential equations (PDEs) – the GPU’s parallel advantage often translates into dramatically higher performance. Tasks like finite-difference stencil updates, linear algebra operations, or particle simulations can typically use thousands of threads, which GPUs handle in parallel. The theoretical peak FLOPs of a high-end GPU can be one to two orders of magnitude higher than a high-end CPU, and crucially, the GPU’s memory bandwidth is usually several times greater. For instance, NVIDIA’s A100 GPU (2020) delivers around 1.6 TB/s memory bandwidth (with HBM2) and  $\sim 9.7$  TFLOPS (FP64), whereas a contemporary CPU might have  $\sim 200$  GB/s bandwidth (DDR4/5) and on the order of 0.5 TFLOPS per core in double precision. Moreover, the latest NVIDIA H100 (Hopper, 2022) is reported to be 2–3 $\times$  faster than its A100 predecessor on many workloads. It achieves this via more cores, higher clocks, HBM3 memory (boosting bandwidth to  $\sim 3+$  TB/s), and new features like FP8 tensor operations [9].

On the AMD side, the MI300X accelerator (2023) emphasizes memory capacity and bandwidth: it includes 192 GB of HBM3 with over 5 TB/s bandwidth, more than double the memory of an H100 (80 GB). Such specs highlight GPUs’ focus on feeding many ALUs with data quickly, a boon for memory-intensive simulations [9][10]. It’s important to note that GPUs only outperform CPUs on tasks that parallelize well and saturate the GPU’s cores and memory bandwidth.

If a problem has strong data dependencies or irregular control flow that prevents parallel execution, a fast CPU core might still win. Similarly, very small problems may

not offset the overhead of offloading to a GPU. But for the *embarrassingly parallel* or structured-parallel problems common in scientific computing, GPUs provide outstanding throughput. Common use cases include:

- **Dense Linear Algebra and Deep Learning:** Matrix-matrix multiplications (GEMM) and convolution operations are extremely parallel. GPUs can perform these at teraflop rates, accelerating neural network training and inference significantly. (Indeed, the deep learning revolution of the 2010s was enabled largely by GPUs – e.g. training the landmark AlexNet CNN in 2012 was feasible in days on GPUs versus weeks on CPUs.) Modern GPUs also have **tensor cores** to further speed up matrix operations in mixed precision, giving them an edge in AI workloads.
- **Finite Difference/Finite Volume Solvers (CFD, FDTD, etc.):** These methods update each point in a grid from values of neighboring points, a natural fit for GPU parallelism. Each mesh cell or lattice point can be updated by a separate thread in a GPU kernel. The GPU's high memory bandwidth helps fetch the neighbor data faster, and the massive core count means an entire 3D domain of millions of cells can be updated concurrently. We will discuss FDTD (finite-difference time-domain for electromagnetics) in detail later – GPUs often yield an order-of-magnitude speedup for such grid-based PDE solvers.
- **Particle Simulations (Molecular Dynamics, N-Body):** GPUs handle large numbers of particles interacting, especially with methods like Verlet integration or Barnes-Hut treecodes, by computing forces in parallel. The many cores help calculate pairwise interactions, and techniques like GPU shared memory can be used to stage particle data to reduce global memory traffic.
- **Signal/Image Processing and FFTs:** Algorithms like Fast Fourier Transforms, filters, or image convolutions are highly parallel. GPUs running libraries (cuFFT, etc.) can transform or filter data arrays much faster than CPUs for large sizes, by exploiting parallel butterfly operations and memory coalescing.

In practice, GPUs have become essential in high-performance computing (HPC). The current #1 supercomputer, Frontier at Oak Ridge (the first exascale system), achieves over 1.1 exaFLOPS primarily by leveraging GPU acceleration – each node has 1 AMD EPYC CPU plus 4 AMD Instinct MI250X GPUs, for a total of ~37,000 GPUs in the machine [11].

This design allows Frontier to attain about 10× the performance of the previous CPU/GPU hybrid system (Summit) within a comparable power envelope [11].

Such GPU-driven architectures now dominate the TOP500 rankings. The takeaway is that for throughput-oriented computing, the **GPU's strengths** – massive parallelism and memory bandwidth – make it far superior to a CPU alone, provided the problem can be structured to exploit parallel execution.

That said, there are **limitations and considerations** when using GPUs for scientific computing. One limitation is precision: some older GPUs had reduced throughput for double-precision (FP64) arithmetic relative to single-precision, although modern HPC GPUs (like A100, H100, MI250/300) offer strong FP64 performance. Another consideration is memory size and latency – GPU on-board memory (VRAM) is typically smaller than host RAM (tens of GB vs. hundreds of GB), so very large datasets might not fit on a single GPU and might require multiple GPUs or streaming approaches. Data transfer between CPU and GPU (over PCIe or NVLink) is also a bottleneck; thus one tries to minimize moving data back and forth during computation. Finally, the programming complexity is higher – one must manage parallel threads, memory explicitly, and sometimes deal with debugging challenges in a massively parallel environment. Despite these challenges, the performance gains in numerical computing are often worth it. In domains from fluid dynamics to climate modeling to machine learning, GPUs enable simulations and training runs that would be impractically slow on CPUs. As a result, most new scientific codes are being written or refactored to utilize GPUs, and frameworks have emerged to aid this (discussed next).

### 4.2.1 Core Architectural Philosophy

The CPU is designed for task diversity and control logic. It is well-suited for workloads that involve frequent branching, unpredictable memory access, and real-time responsiveness. A modern CPU typically contains between 4 and 64 general-purpose cores, each of which can independently execute multiple instruction types and manage its own control flow. These cores are complex and large, incorporating deep pipelines, branch predictors, and multi-level caches.

In contrast, the GPU is designed for task uniformity and raw throughput. It is built around the idea that if the same arithmetic instruction needs to be applied to a large number of data points, it should be done in parallel. Rather than investing in large, complex cores, the GPU incorporates thousands of small, lightweight arithmetic units, often referred to as CUDA cores (in NVIDIA architecture), which operate under a SIMT (Single Instruction, Multiple Threads) execution model.

In effect, the CPU is a master of orchestration and decision-making; the GPU is a master of repetition and scale.

### 4.2.2 Memory Hierarchy and Latency Handling

Another major architectural divergence lies in the handling of memory latency. CPUs rely heavily on cache hierarchies (L1, L2, and L3) to keep frequently accessed data close to the processor, thereby reducing memory latency. This model is highly effective for programs with irregular memory access patterns, such as decision trees, interpreters, and system processes.

GPUs, by contrast, are not as reliant on cache. Instead, they tolerate memory latency by oversubscribing threads. When one group of threads (a warp) is waiting for data, the GPU schedules another warp for execution. This context-switching strategy, combined with high-bandwidth global memory and fast on-chip shared memory, ensures that the arithmetic units remain busy.

Furthermore, GPU shared memory can be explicitly managed by the programmer to optimize for memory coalescence and avoid bank conflicts, a degree of control that is typically not available in CPU cache hierarchies.

### 4.2.3 Programming and Execution Models

In CPU-based scientific computing, parallelism is achieved using threads (via OpenMP or POSIX) or distributed processes (via MPI). These models are coarse-grained, requiring explicit thread synchronization, load balancing, and often dealing with race conditions.

GPU programming models, such as CUDA, are built for fine-grained parallelism. A typical CUDA kernel may launch hundreds of thousands of threads, organized into blocks and grids, which are executed across thousands of arithmetic units in lockstep. The programmer defines the per-thread operation, and the hardware schedules the execution dynamically. This model is particularly well-suited to stencil computations, matrix operations, and vectorized updates — all of which are central to numerical PDE solvers like FDTD.

One important limitation of GPUs is that execution divergence within a warp leads to serialization. If threads in a warp follow different branches of a conditional, the instructions must be executed sequentially, reducing efficiency. Thus, well-optimized GPU kernels minimize branching and favor regular control flow.

### 4.2.4 Performance Comparison in Real Workloads

Empirical studies have shown that GPU-accelerated solvers can outperform CPU-based implementations by an order of magnitude or more. For instance, in FDTD simulations involving 3D domains, GPU-based solvers can achieve 30× to 100× speedups over multi-threaded CPU implementations, depending on problem size and optimization level [8,9].

For highly structured computations—such as FDTD on Yee grids—these differences translate into vastly different real-world performance characteristics. A single time-step update of an electromagnetic field component involves six neighboring grid points and a constant arithmetic expression. With one GPU thread per voxel, these updates can be computed in parallel for the entire domain, exploiting all available hardware units.

In contrast, a CPU would either require expensive thread spawning or rely on SIMD vectorization, which is harder to scale to massive grid sizes due to memory bottlenecks and synchronization overhead.

#### 4.2.5 Scientific Relevance and Shifting Paradigms

The emergence of GPUs has forced a paradigm shift in how scientific software is developed and deployed. In the past, simulation runtimes were bottlenecked by CPU speed and memory latency. Today, with GPU acceleration, the bottlenecks have moved to algorithm design and data transfer optimization.

Fields such as fluid dynamics, seismology, quantum chemistry, and electromagnetics are rapidly transitioning their core solvers to GPU-compatible architectures. This transformation is not only performance-driven but also practical: cloud platforms like AWS and Google Cloud provide on-demand GPU instances; open-source libraries (e.g., cuFFT, cuBLAS, ArrayFire) offer optimized primitives; and high-level languages like Julia and Python have matured GPU bindings, reducing the barrier to entry for researchers.

In metasurface design—where optimization problems require forward and adjoint simulations repeated thousands of times—GPU acceleration is not optional; it is foundational. Without it, centimeter-scale FDTD simulations with full-wave accuracy would require impractical runtimes of weeks or months.

## 4.3 GPU Programming Ecosystem

The GPU programming ecosystem has matured greatly since the early days of GPGPU. Two primary software platforms emerged: CUDA and OpenCL. CUDA, from NVIDIA, is a proprietary API but has become the dominant method for GPU computing given NVIDIA's market share in HPC. It provides a C/C++ (and Python via libraries like Numba or PyCUDA, etc.) environment to write kernels (functions executed in parallel by many GPU threads) and explicit memory management routines to transfer data to/from GPU memory. CUDA's popularity is bolstered by a rich set of libraries (for linear algebra, FFT, deep learning, etc.) and development tools. OpenCL, managed by Khronos, is an open standard that targets GPUs from multiple vendors (as well as CPUs, FPGAs, etc.). OpenCL is conceptually similar (kernels in C-based language), but its portability sometimes comes at the cost of having to manage more details for different hardware. While OpenCL saw use especially on AMD and Intel GPUs, its adoption in HPC has been limited compared to CUDA. More recently, AMD introduced HIP (Heterogeneous Compute Interface), which is essentially a CUDA-like API for AMD GPUs, allowing code written for CUDA to be compiled for AMD hardware with minimal changes. This is part of AMD's ROCm platform to support HPC and AI on their GPUs. Intel has its oneAPI and Data Parallel C++ (DPC++) for their GPUs. Despite these alternatives, much GPU computing knowledge is transferable across APIs – the core concepts of parallel kernels, memory hierarchy, and execution model are similar.

On top of these low-level APIs, the ecosystem includes language bindings and higher-level frameworks. A notable example for scientists is Julia: the Julia programming language, designed for high-performance numerical computing, has robust GPU support. The package `CUDA.jl` provides Julia bindings to CUDA, allowing Julia code to launch GPU kernels and manage GPU arrays. Impressively, one can write GPU kernels directly in Julia, leveraging its high-level syntax and compiler. In fact, “writing kernels in Julia is very similar to writing kernels in CUDA C/C++”, as `CUDA.jl` exposes low-level CUDA functionality with Julian syntax. For example, one can define a Julia function and launch

it on the GPU using the `@cuda` macro. Below is a simple illustration of a Julia CUDA kernel that adds two vectors:

```
function vector_add!(C, A, B, N)
    i = (blockIdx().x - 1) * blockDim().x + threadIdx().x
    if i <= N
        @inbounds C[i] = A[i] + B[i]
    end
    return
end
threads = 256
blocks = ceil(Int, N / threads)
@cuda threads=threads blocks=blocks vector_add!(dC, dA, dB, N)
```

Here, `A`, `B`, `C` are `CuArrays` (GPU arrays) and the kernel computes  $C = A + B$ . The built-in `CUDA.jl` functions `blockIdx()`, `threadIdx()`, etc., give each thread an index so it works on a distinct element. This example shows how Julia enables a fairly high-level expression of GPU kernels while the framework handles compilation to PTX (the GPU's assembly) under the hood. The benefit is that scientists can stay in one language (Julia) for both CPU and GPU code, leveraging Julia's productivity and multiple dispatch, rather than writing C++/CUDA separately. Julia's GPU support isn't limited to CUDA either; packages for OpenCL and even Metal (on macOS) exist (e.g. `Metal.jl`). Additionally, Julia offers high-level abstractions: one can often write element-wise array operations or use Julia's native array syntax, and if the array is a `CuArray`, those operations are implicitly performed on the GPU (through broadcasting and fusion into GPU kernels). This lowers the barrier to entry for utilizing GPUs.

Apart from language-centric approaches, there is a rich set of GPU-accelerated libraries and frameworks. NVIDIA provides many libraries as part of its CUDA Toolkit: for example, `cuBLAS` for basic linear algebra (dense matrix multiply, etc.), `cuFFT` for fast Fourier transforms, `cuSparse` for sparse matrix operations, Thrust (a C++ template library for GPU parallel primitives), and others covering random number generation, image processing, deep learning (`cuDNN`), etc. These are highly optimized and often serve as the "backend" for higher-level frameworks (like how TensorFlow or PyTorch

use cuDNN, cuBLAS etc. internally). Using such libraries can save development time and ensure one gets near-peak performance for standard operations.

There are also third-party libraries that provide a more platform-agnostic high-level API. One example is ArrayFire, an open-source array library that can target CUDA or OpenCL devices (or even CPU) [12]. ArrayFire provides a high-level matrix and vector manipulation API in languages like C++, Python, R, etc., and internally it will run the computations on the GPU without the user writing kernels explicitly. Its key features include an JIT (Just-In-Time) compiler that can fuse array operations, an easy array notation, and support for many algorithms (FFT, image filters, linear algebra, statistics) implemented for the GPU.

The idea is to let scientists write code similarly to MATLAB/NumPy but get GPU acceleration automatically. According to NVIDIA, ArrayFire's approach can achieve 2× to 100× speedups over equivalent CPU code, depending on the application, often with far fewer lines of code than native CUDA C++ [14].

Another cross-platform library is OpenACC, which is actually a set of compiler directives (pragmas) that can instruct the compiler to offload certain loops to the GPU, aiming to simplify porting existing C/FORTRAN codes to GPUs. While not as flexible as writing CUDA, OpenACC is used in some legacy HPC codes to incrementally add GPU support.

For domain-specific frameworks, almost every area has developed GPU-accelerated tools. Deep learning has TensorFlow and PyTorch (which use GPUs via CUDA under the hood). Molecular dynamics has GPU-accelerated codes like AMBER and HOOMD-blue. Climate modeling and CFD codes incorporate CUDA or use directive-based approaches to run on GPU clusters. Even higher-level environments like MATLAB or Mathematica have added GPU support for certain functions via CUDA integration.

In summary, the programming ecosystem now ranges from low-level control (CUDA C/C++ and OpenCL) up to high-level, GPU-aware libraries and languages (like Julia, Python libraries, domain-specific tools). This means scientists can engage with

GPU computing at the level they are comfortable: either by writing custom kernels for maximum control or by leveraging existing libraries and just ensuring their data is on the GPU. For the work in this thesis, we utilize Julia with CUDA.jl, giving us the ability to write custom update kernels for our electromagnetic simulations in a high-level manner, while still tapping into the performance of CUDA. We also make use of CUDA libraries (for example, cuFFT might be used for spectral analysis, etc.) when appropriate, to avoid re-inventing well-optimized wheels.

## 4.4 GPU Optimization Techniques

Achieving high performance on GPUs requires understanding and exploiting the GPU architecture's details. Writing a correct GPU program is only the first step; optimizing it can yield enormous further speedups by utilizing hardware resources more efficiently. Key GPU optimization techniques include:

- **Warp Scheduling and Occupancy:** NVIDIA GPUs execute threads in groups of 32 called *warps*. All threads in a warp execute the same instruction in lockstep (SIMD fashion). To hide memory latency, the GPU hardware rapidly switches between warps. The *occupancy* refers to how many warps are resident and ready to run on each multiprocessor (SM). Higher occupancy (more warps per SM) can help hide latency, but after a point the benefits may plateau. An optimal kernel uses a suitable number of threads per block and registers such that many warps can be active without exhausting resources. Tuning launch parameters (threads/block) and limiting register and shared memory usage per thread (to allow more threads) are common tactics to maximize occupancy[13]. NVIDIA's tools (like *occupancy calculator* or the profiler) can report if a kernel is limited by occupancy. In practice, one tries to use near the maximum threads per SM (e.g. 1024 threads on many architectures) unless there's a reason not to.
- **Memory Coalescing:** Global memory (device DRAM) bandwidth is critical to performance. GPUs achieve their peak bandwidth when parallel memory accesses by threads can be coalesced into a few large transactions. Specifically, threads in

the same warp accessing consecutive memory addresses can be serviced by a single wide memory transaction rather than many small ones. Coalesced memory access means arranging data and access patterns so that, for example, thread 0 accesses element 0, thread 1 element 1, etc., in a contiguous manner [13].

- Non-coalesced access (e.g. if threads strided through memory or accessed scattered addresses) forces the hardware to issue multiple transactions and wastes bandwidth. Thus, data structures are often organized for structure of arrays (SoA) rather than array of structures, to ensure threads read contiguous chunks. Techniques like using padding to align data to 128-byte boundaries (common transaction size) or reading/writing via shared memory in a coalesced way can help. Effective use of the memory hierarchy (L1 cache, L2 cache) also matters – but on GPUs, caches are smaller relative to working sets, so memory bandwidth and coalescing often dominate performance for bandwidth-bound kernels.
- Shared Memory Tiling: Each SM has a fast on-chip scratchpad called shared memory (also known as the L1/shared memory, configurable in size like 48KB or 100KB depending on GPU). Shared memory has much higher bandwidth and lower latency than global memory, but is only accessible by threads within the same thread block. A common optimization is tiling: a block of threads cooperatively loads a tile of data from global memory into shared memory, then each thread can reuse values from shared memory multiple times for computations, rather than each time going to slow global memory. This is especially useful in matrix operations (e.g. matrix multiply uses tiling to load sub-blocks of the matrices) and stencil computations. In an FDTD update, for instance, one could load the neighboring cells' field values into shared memory and then perform updates, so that each neighbor is read once from global memory and then used by several threads that need it. Shared memory is divided into banks, so one must also ensure bank conflicts are minimized (by structuring accesses such that no two threads access different addresses in the same bank in the same cycle). With careful padding or access patterns, bank conflicts can be avoided, making shared memory accesses as fast as registers. Tiling and shared

memory use can dramatically reduce effective global memory traffic. However, using shared memory also increases register usage and limits occupancy, so there is a trade-off to balance [13].

- **Instruction-Level Parallelism and Scheduling:** While GPUs are massively parallel, each thread's performance also matters. Modern GPUs can issue multiple instructions per thread concurrently (if they are independent) through instruction-level parallelism (ILP). Ensuring that the compiler can see independent operations (e.g., avoid long dependency chains in code) can help utilization of execution units. Loop unrolling and using fewer divergent branches can sometimes increase ILP. Moreover, the compiler and hardware scheduler try to interleave arithmetic and memory operations to hide latency (since global memory load latency is high). It's often beneficial to intermix computation with memory accesses such that while some data is being fetched, the arithmetic units can work on other data (this overlaps latency with useful work). This overlap is mostly handled automatically by having enough warps (threads) in flight, but writing kernels with a balance of arithmetic and memory ops (avoiding being purely memory-bound or compute-bound) can yield better *achieved* occupancy of functional units.
- **Control Flow Divergence:** As mentioned, threads in a warp execute in lockstep. If there is a branch (e.g., an if/else or loop) where threads in the same warp take different paths, the warp will have to execute each path serially (masking off threads that aren't on that path). This warp divergence effectively reduces parallel efficiency, as some threads are idle while others execute the branch, and then vice versa.
- Therefore, minimizing divergent branching is a key optimization. This can mean restructuring algorithms to avoid per-thread conditional logic when possible. For example, one can often use predication (execute both paths but mathematically mask out the effect of one) or segregate data such that threads in the same warp are likely to follow the same execution path. In graphics or some simulations, one might sort data so that warps handle similar cases. In our FDTD context, an example would be handling boundary conditions: instead of `if (onBoundary) { do`

boundary update } else { normal update } within the kernel (which causes divergence between boundary threads and interior threads), one might launch a separate small kernel just for boundary cells or use specialized instructions. Some divergence is inevitable (e.g., loops that iterate a variable number of times per thread), but keeping warps coherent maximizes throughput [13].

- **Loop Unrolling and Memory Access Patterns:** Unrolling small loops can reduce loop overhead and sometimes enable coalescing or ILP, but unrolling large loops can increase register pressure. Memory access pattern optimizations include using texture memory or read-only data cache for read-mostly data, which can leverage cache locality. In older GPUs, using texture fetches could also automatically handle some unaligned access patterns efficiently. Nowadays, simply ensuring aligned, coalesced global loads is typically sufficient, and the compilers/hardware do a good job with caches.
- **Precision and Math Optimizations:** Where acceptable, using lower precision (e.g., float instead of double, or even half-precision or TensorFloat for AI) can boost performance significantly, as GPUs have higher throughput for lower precision. Additionally, using intrinsic fast math operations (like `__sinf` vs. `sin` in CUDA, which might use a faster approximation) can speed up code that is not precision-critical. One must use these with care in scientific codes to not sacrifice required accuracy.
- **Profiling and Iterative Tuning:** GPU optimization is often an iterative process guided by profiling. Developers use profiling tools to identify bottlenecks. NVIDIA's Nsight Systems provides a timeline view of CPU-GPU interaction and how kernels and memory copies overlap, helping find if the GPU is underutilized or if there are idle gaps [14].
- **Nsight Compute (and formerly nvprof/CUDA Profiler)** is a detailed kernel profiler that can report memory throughput, achieved occupancy, warp execution efficiency, etc., for each kernel launch[15].

- By examining these metrics, one can tell if a kernel is memory-bound or compute-bound, how effectively it coalesces memory, whether there is divergence, etc. For example, profiling might reveal only 50% global memory utilization due to uncoalesced accesses – prompting a rewrite of data layout. Or it might show many branch divergences, leading to reconsidering warp-level logic. Tools like Nsight Compute also support guided analysis: they flag issues like “dramatically low occupancy” or “memory bound, try increasing coalescing”.
- Iteration and Refinement: An optimized GPU code may employ multiple of the above techniques. One common workflow is: start with a correct kernel, profile it to see bottlenecks, then apply one optimization at a time (e.g., introduce shared memory to tile, or try different thread block sizes) and measure impact. It’s important to only optimize as needed – sometimes straightforward code is already limited by external factors (like memory bandwidth) and fancy optimizations yield little gain if the bottleneck lies elsewhere. The *roofline model* is often used to reason about this: it provides an upper bound on performance based on operational intensity (flops per byte). If a kernel is at the roofline for memory bandwidth, further caching or tiling might help if it increases flops/byte; if it’s at the compute roofline, then one ensures all execution units are busy (maybe by increasing occupancy or using faster math operations).

In practice, key optimizations for our use case (FDTD on GPUs) will revolve around memory access (since FDTD is typically memory-bandwidth bound) and managing parallelism efficiently. We will ensure coalesced reads/writes of field arrays, use shared memory to stage field components where beneficial, and avoid warp divergence in update loops (the update algorithm is pretty uniform for interior grid points). We will also use streams and concurrency to overlap data transfers (if any) with computation, and possibly overlap different computations if the algorithm permits (though FDTD is mostly a synchronous time-step loop, so concurrency is mainly between multi-GPU communications and computations). Tools like Nsight will verify our kernels achieve high occupancy and bandwidth utilization. Indeed, by applying these optimization techniques, one can often approach a significant fraction of peak

throughput. NVIDIA's own optimization guides emphasize these points: “*Ensure global memory accesses are coalesced*”, “*Avoid different execution paths within the same warp*”, “*Adjust kernel launch configuration to maximize utilization*”, etc., which align with the strategies described above.

Finally, it's worth mentioning that optimization is hardware-specific. Tactics that work well on one generation (say, maximizing occupancy) might yield less benefit on another if that architecture has different memory behavior or warp scheduling. Thus, one should profile on the target GPU architecture (Ampere vs Hopper, etc.) and be aware of new features (for example, Hopper introduced shader execution reordering (SER) to mitigate warp divergence in certain cases – a hardware feature that can automatically reschedule divergent warps to improve locality). Keeping up with architecture trends (via whitepapers and CUDA programming guide updates) is part of the optimization cycle in GPU computing.

## 4.5 GPUs in FDTD and Electromagnetic Simulations

Electromagnetic simulations using the Finite-Difference Time-Domain (FDTD) method are exceptionally well-suited to GPU acceleration. FDTD involves discretizing Maxwell's curl equations on a grid (often a Yee grid) and updating the electric and magnetic field components in time. Each field component at a given grid cell is updated from the values of neighboring cells (typically the nearest neighbors in the grid). This locality and regularity of computation maps naturally to the GPU's parallel model: one can assign one thread to update each field component (or each Yee cell) in parallel, since at a given time-step all updates are independent (relying only on the previous time-step fields). The structured grid means memory access patterns can be made regular (e.g., storing field arrays in row-major order so that adjacent threads access adjacent memory locations). Moreover, FDTD computations are mostly simple arithmetic (adds, multiplies) and thus are usually *memory-bandwidth bound* – precisely where GPUs, with their high bandwidth, excel.

To understand the mapping, consider the standard 3D Yee grid: each grid cell has 6 field components ( $E_x$ ,  $E_y$ ,  $E_z$ ,  $H_x$ ,  $H_y$ ,  $H_z$ ) staggered in space. The update for, say,  $E_z$

at grid index  $(i,j,k)$  will use the  $H_x$  and  $H_y$  values from the neighboring cells around  $(i,j,k)$ .

These operations – differencing neighboring field values – are identical for every cell in the interior of the domain. On a GPU, we can launch a kernel where each thread handles the update of one field component (or one Yee cell’s 6 components, depending on strategy) at a given time-step. All threads execute the same update equations but on different data (Single-Instruction Multiple-Data), which is exactly what GPUs are designed for. The only dependencies are that we must use field values from the previous time-step (often we use a ping-pong buffering or update H and E in alternation so they don’t overwrite needed values). Thus, FDTD is “embarrassingly parallel” over the grid. In fact, even on CPUs one would parallelize FDTD across cores (domain decomposition or multi-threading). GPUs simply take this to a much finer granularity (hundreds of thousands of threads updating different points simultaneously).

Yee cell arrangement in FDTD: electric field components (E, shown in red) and magnetic field components (H, gray) are defined on a staggered grid (the Yee lattice). Each E component is updated from the circulation of H around it, and vice versa. This localized update enables parallel updates of all cells.

The nature of FDTD’s computations means that a GPU implementation can often reach a large fraction of peak memory throughput. Studies have shown dramatic speedups: for example, a recent photonics FDTD simulation article reported one or two orders of magnitude shorter runtimes on GPUs compared to CPUs – a single high-end GPU completing in minutes a simulation that took hours on a multicore CPU. In concrete terms, that article achieved up to 20–33 billion cell updates per second on one NVIDIA A100/H100 GPU.

Other literature also reports  $10\times$ – $100\times$  speedups for FDTD on GPU. Early GPGPU efforts (circa 2007, using GPUs via graphics shaders) already demonstrated  $>50\times$  speedup for sizable FDTD grids, and modern GPUs are far more powerful [15]. A 2024 study [16] of a 3D FDTD code written in CUDA showed about  $17\times$  higher throughput than a 24-core Intel CPU, when the problem size fit entirely in the GPU’s

memory. The speedups stem from both the raw hardware capability and from algorithmic adjustments to fit FDTD to GPU architecture.

Why is FDTD bandwidth-bound? Each field update involves only a handful of arithmetic operations (say, 6–12 FLOPs) but requires reading several field values from memory (perhaps 6–8 values). The ratio of memory operations to compute is high, so the limiting factor is usually how fast those values can be fetched/stored. GPUs mitigate this by high bandwidth and by letting many threads fetch in parallel. Additionally, FDTD often accesses memory in a structured pattern (unit-stride through arrays when updating a plane of the grid, etc.), which can be arranged to be coalesced. For example, if we store the  $E_z$  field in a 3D array  $E_z[i,j,k]$ , we might have threads indexed such that thread  $(i,j,k)$  reads  $H_x[i,j,k]$  and  $H_y[i,j,k]$  (or nearby indices) – careful ordering in memory can make these reads coalesced among warps. We typically store each field component in a separate array (SoA layout) to ensure contiguous memory for that component, rather than interleaving different components which could disrupt coalescing.

#### **GPU-specific acceleration strategies for FDTD (Yee grid):**

- *Spatial Blocking / Tiling*: If the grid is large, one can divide it into subdomains for processing in shared memory. For instance, a block of threads might load a tile of  $E$  and  $H$  values into shared memory, perform several update steps, then write back. This can reduce global memory traffic by reusing values. However, straightforward FDTD updates only use near-neighbors, so the benefit of tiling is somewhat limited to reducing multiple loads of the same data when updating adjacent points. Still, in 3D FDTD the same  $H$  value is used in updates of two adjacent  $E$  cells (due to staggered grid), so one can share that via shared memory instead of each thread loading it independently. One strategy reported is assigning multiple Yee cells to each thread, so that a single thread updates, say, two or four adjacent cells in a line.
- This way, once it loads a piece of data, it can use it for updating multiple positions, increasing arithmetic intensity. The downside is increased register use and complexity.

- *Temporal Blocking:* A more advanced idea is updating multiple time-steps in shared memory before writing back (temporal blocking), but for FDTD this is tricky due to the CFL stability limit and needing updated neighbors each time-step. It's generally not used because it would require storing a "thick" halo of neighbors in shared memory across time, which often doesn't pay off.
- *Use of Registers and Fused Multiply-Add:* GPUs have FMA (fused multiply-add) instructions which FDTD can use (the update equations are linear combinations). Utilizing FMAs improves both speed and precision slightly. The compiler usually does this automatically if using high-level languages.
- *Boundary Condition Handling:* One area that can hurt performance is boundary condition checks. In a naive implementation, each thread might check if its index is at the domain boundary to apply a different update (e.g., ABC or PEC boundary condition). These if statements can cause warp divergence. A common approach is to avoid branching in the inner kernel. Instead, one can update interior points with one kernel launch, and handle boundaries with separate specialized kernels or pad the domain with ghost cells so that interior update formula can be applied everywhere uniformly (with ghost cells providing the boundary condition implicitly). We often allocate a few extra cells for PML layers or absorbing boundaries so that the main update loop doesn't have to constantly check indices – this way all threads in a warp execute the same instructions, and special cases are handled outside the main compute kernel.
- *Resource balancing:* If the FDTD update uses a lot of registers, it may limit occupancy. Sometimes it can be beneficial to manually influence the register usage (e.g., by breaking the update into two kernels or using compiler flags) if it allows more warps to reside. However, modern GPUs have enough registers that a single FDTD update kernel typically still allows full occupancy.
- *Utilizing multiple GPUs:* For very large simulations, a single GPU may not have enough memory or compute power. FDTD can be extended to multi-GPU by spatial domain decomposition. For example, if you have two GPUs, you can split

the computational domain into two subdomains (e.g., along the z-axis: GPU0 handles lower half of the grid, GPU1 the upper half). After each time-step, the boundary fields between the subdomains must be exchanged (the “halo” or ghost cell values). This requires communication (which can be done via GPU peer-to-peer transfers or through the host with MPI). The overhead of this communication must be minimized – typically one uses asynchronous data transfer overlapping with computation on interior points, or uses high-speed interconnects (NVLink or InfiniBand). We will discuss a specific multi-GPU case in the next section. Here we note that this parallelization is very feasible because FDTD’s locality means each subdomain only needs a thin slice of data from its neighbor (e.g., one cell layer). Thus the communication footprint is much smaller than the total data, and with efficient halo exchange the scaling to multiple GPUs can be good.

- In fact, multi-GPU FDTD can attain near-linear speedup until communication latency/bandwidth starts to dominate. Modern GPU clusters often use domain decomposition combined with MPI (or Nvidia’s NCCL for intra-node communication) to scale FDTD to dozens of GPUs.

To cement these ideas, consider a concrete performance result: In an electromagnetic simulation of a large photonic structure, Minkov *et al.* (2024) measured around 33 billion cell updates/s on a single NVIDIA H100 GPU [15]. They project that upcoming GPUs (like a theoretical Blackwell B200) could exceed 100 Gcells/s. They also demonstrated that running the same simulation on a high-end 64-core CPU yielded about 100× slower performance [15].

This illustrates that GPUs not only speed up individual simulations but also enable tackling larger problems. Another example: using 32 A100 GPUs together, they ran a massive simulation with 3 billion cells for over 300k time-steps (a problem that would be infeasible on a single machine) in under an hour. Such capabilities are revolutionizing electromagnetic design – enabling full-wave 3D simulations of devices that previously required approximations.

In summary, GPUs are a natural fit for FDTD because of:

1. Parallelism: Each grid point update is independent (for a given time-step) – ideal for GPU’s SIMT execution.
2. Regular memory access: The grid structure means memory access patterns can be made coalesced and predictable.
3. High arithmetic intensity (if using advanced techniques): While basic FDTD is memory-bound, techniques like updating multiple cells per thread or vectorizing over multiple modes (in cylindrical coordinates, for instance) can increase flop/byte ratio to use more of the GPU’s ALU potential.
4. Bandwidth advantage: FDTD pushes a lot of data, and GPUs simply have more memory bandwidth than CPUs (often 5–10×), so they alleviate the primary bottleneck.
5. Scaling: If one GPU isn’t enough, domain decomposition allows FDTD to leverage many GPUs with relatively small communication overhead (especially with GPU-direct networking and fast interconnects available in modern HPC systems).

Therefore, implementing FDTD on GPUs – as we do in this thesis – is a very effective way to accelerate electromagnetic simulations. The next section provides a detailed case study of doing exactly this for an axisymmetric FDTD problem using Julia and CUDA.

## 4.6 Accelerating Axisymmetric FDTD using CUDA

In this case study, we focus on an FDTD simulation with discrete axial symmetry, essentially a cylindrical coordinate FDTD for problems that are invariant (or periodic) in the azimuthal direction. The idea is to exploit the fact that if the structure and sources have a symmetry around an axis, the fields can be expanded in azimuthal modes  $e^{im\theta}$ , reducing a 3D problem to a set of 2D problems (in  $r$  and  $z$  coordinates for each mode  $m$ ). We implemented an axisymmetric FDTD in Julia, using `CUDA.jl` to accelerate it on NVIDIA GPUs. Here we discuss how we mapped the algorithm to the GPU, the memory

layout strategies for cylindrical coordinates, and performance results including multi-GPU scaling.

Mapping field updates to Julia GPU kernels: In cylindrical coordinates (assuming axial symmetry), Maxwell's equations can be written in terms of  $r$ ,  $z$ , and  $\theta$  components. For an  $e^{i m \theta}$  dependence, the fields for each mode  $m$  satisfy a set of coupled 2D equations in  $(r, z)$  (with  $m$  appearing as a parameter in the coefficients).

We discretize  $r$  and  $z$  on a Yee grid (staggering  $E_r$ ,  $E_z$ ,  $H_\theta$  vs.  $E_\theta$ ,  $H_r$ ,  $H_z$ , etc. appropriately for cylindrical coordinates). The update equations for each mode  $m$  are very similar to the Cartesian Yee scheme, with additional  $m/r$  terms coupling some components due to the curvilinear coordinate. We handle each mode's fields as separate 2D arrays. For example, for mode  $m$  we maintain arrays for  $E_r^{(m)}(r, z)$ ,  $E_\theta^{(m)}(r, z)$ ,  $E_z^{(m)}(r, z)$ , and similarly  $H_r^{(m)}$ ,  $H_\theta^{(m)}$ ,  $H_z^{(m)}$ . These are real-valued fields (assuming we split complex  $e^{i m \theta}$  fields into real and imaginary parts if needed, or treat the fields as complex if using complex updates).

Using Julia and CUDA.jl, we wrote GPU kernels to update these field arrays in time. We opted to update each field component in a separate kernel launch (though other strategies like updating all in one kernel were possible). For instance, one kernel updates all  $E$  fields from the  $H$  fields, then the next kernel updates all  $H$  fields from the new  $E$  fields (Yee's algorithm typically updates  $H$  half a step offset from  $E$  in time). By separating them, we keep each kernel simpler and allow the compiler to optimize each update loop aggressively. Each kernel is a Julia function, e.g. `update_Eθ!(Er, Eθ, Ez, Hr, Hz, m, ...)` that updates the  $E_\theta$  field for all interior grid points given the neighboring  $H$  values. We launch it with a 2D configuration of threads covering the  $(r, z)$  plane. Internally, the kernel uses the thread's indices (obtained via `CUDA.jl's threadIdx()` and `blockIdx()`) to compute the corresponding grid indices ( $i_r$ ,  $i_z$ ) in the simulation domain, and then executes the update formula for that cell. The code uses Julia's ability to include inline math expressions (we derived the FDTD update equations for cylindrical coords with an  $m$ -mode expansion – these include terms like  $E_\theta$  update needing  $(H_z(r+\Delta r) - H_z(r)) / \Delta r$  and an  $m$ -dependent term  $-(m H_r / r)$ , etc.). These were directly coded in the

kernel. We take care to avoid divisions inside loops where possible (e.g., precompute  $1/r$  arrays for each grid radius to multiply by, instead of dividing per thread) [16].

Memory layout considerations: We store each field as a 2D array with indices  $[i_r, i_z]$ . In Julia (which uses column-major order by default), we had to decide which index varies fastest in memory. We chose to make the radial index  $i_r$  correspond to contiguous memory locations, so that threads updating neighboring radial positions access contiguous memory. This means that when we launch a 2D grid of threads, if we align threads such that adjacent threads differ in  $i_r$ , those threads will be reading/writing adjacent memory addresses – yielding coalesced accesses. Specifically, we launch threads with a configuration like  $(\text{blockDim}_x, \text{blockDim}_y) = (32, 8)$  for example, where  $x$  direction of the block maps to radial direction and  $y$  maps to  $z$  direction. Then each warp (32 threads along radial direction) accesses, say,  $H_r[i_r, i_z]$  for  $i_r$  from  $N$  to  $N+31$  – which is contiguous in memory. This turned out to maximize global memory throughput. The downside is that Julia’s column-major means  $i_r$  is the first index of the array and is contiguous (since by default arrays are  $A[\text{column}, \text{row}]$  as memory contiguous along column index). We had to be mindful of this when interfacing with some Julia libraries or when allocating arrays to ensure the intended orientation.

We also padded the arrays slightly in the radial direction to avoid bank conflicts in shared memory when doing certain coalesced loads. However, since our main memory accesses are coalesced global loads/stores, we didn’t rely heavily on shared memory in the final implementation. We attempted an optimization using shared memory tiling for the  $H_\theta$  update (because it involved  $m$  divided by  $r$  which at  $r=0$  is singular and needed a special handling), but found that a simpler approach of treating  $r=0$  explicitly with an if-condition (in a separate kernel for axis points) was easier and did not hurt performance noticeably, given that only a single column at  $r=0$  was affected (small fraction of threads diverging).

Another memory-related strategy was using structure-of-arrays: we store each field in its own array, rather than interleaving  $E_r, E_\theta, E_z$  in one struct per grid cell. This avoids strided accesses. It also aligns with how we launch separate kernels per component. If memory were a concern, one could pack fields, but then GPU threads

would have to extract components, which complicates coalescing. We opted for clarity and coalescing with separate arrays.

Handling multiple modes ( $m$  values): An axisymmetric simulation with discrete rotational symmetry might involve summing over a few azimuthal modes. In our implementation, we can simulate multiple  $m$  modes simultaneously. We have two choices: (1) run each mode's FDTD update in series (one after the other on the GPU), or (2) run them in parallel if multiple GPUs are available (assign different  $m$  to different GPUs), or (3) even launch combined kernels that handle multiple modes by adding a third dimension to the thread grid. We experimented with a batched approach where the third index of a thread covers the mode number. This is feasible since the modes are independent (for linear media). However, doing so can reduce cache locality (different modes' data are far apart in memory) and can also increase register pressure (because you then have to carry multiple mode indices or arrays in the same kernel). Ultimately, we found it simpler to loop over modes on the GPU sequentially for a single GPU case. When using multi-GPU, a very effective approach was to distribute modes across GPUs: e.g., GPU0 handles mode 0 and 1, GPU1 handles mode 2 and 3, etc. This *mode parallelism* has virtually no coupling between GPUs (each mode's simulation is independent), so it scales perfectly linearly with the number of GPUs (apart from trivial overhead of launching kernels on each device). In cases where the number of modes is less than number of GPUs or vice versa, one can combine mode and spatial decomposition.

Multi-GPU domain decomposition: For very large  $r$ - $z$  grids, we also implemented a spatial split along the  $z$  axis (since the geometry is elongated in  $z$  for our test problem). We divided the computational domain into slabs of equal  $z$ -extent, each slab handled by one GPU. After each time-step (or after a few steps), the field values at the interfaces between slabs must be exchanged: e.g., GPU 0 needs the updated fields just beyond its top boundary from GPU 1, and GPU 1 needs the fields just below its bottom boundary from GPU 0. We used MPI (via Julia's MPI.jl) for this halo exchange, leveraging CUDA-aware MPI to send data directly from GPU memory. The process is: each GPU evolves its subdomain's fields for one time-step, then we perform non-blocking MPI

sends/receives of the boundary ghost cells (the outermost one or two rows of cells) between neighboring ranks, then continue to the next time-step. We hide the communication latency by overlapping it with computation: specifically, we launch the interior update (which doesn't depend on halo) and while that is running, we initiate MPI transfer of the halo regions that were computed in the previous step. This overlap was facilitated by using multiple CUDA streams – one for compute, one for communication – so that GPU kernels and data copies can run concurrently. Thanks to this overlap, the measured communication overhead was small; our multi-GPU strong scaling was reasonably good, with the parallel efficiency above 80-90% for up to 4 GPUs on a single node. (On multiple nodes, latency of interconnect becomes more of a factor, but with InfiniBand and large messages we still got decent efficiency.)

One practical consideration was ensuring synchronization between steps: after communication, we had to update the ghost cells on each GPU before the next kernel launch. We used CUDA events to synchronize the stream that copies boundary data with the stream that runs kernels. Julia's higher-level abstractions (like the `ImplicitGlobalGrid.jl` library [discourse.julialang.org](https://discourse.julialang.org)) can automate some of this halo management, but we implemented it manually for full control. Each GPU's work loop thus was: receive halo from neighbor into ghost array, launch kernel to update boundary cells using received data (this is a small kernel just for the first couple of grid rows), then proceed with main bulk update, then send updated halo out. Using separate small kernels for the boundary vs interior prevented thread divergence inside a single kernel and also allowed us to tailor each kernel's memory access for its region.

Performance insights: Our GPU implementation in Julia achieved significant acceleration over a CPU version. On a single NVIDIA A100, we observed about 25× speedup compared to a 16-core Intel Xeon running an optimized CPU version of the same FDTD (in C with OpenMP). This was for a grid of size  $\sim 1000$  (r)  $\times$  1000 (z) for the lowest azimuthal mode ( $m=0$ ). The speedup primarily came from the GPU's ability to update all grid points in parallel, whereas the CPU version was limited by memory bandwidth across fewer threads. The GPU code's memory throughput reached around 70% of the A100's theoretical peak, which is quite good considering some overhead from

indexing and the  $m/r$  terms. We found that using Julia did not impede performance – the generated GPU code was on par with what a CUDA C implementation would produce, thanks to Julia’s LLVM-based compilation. This demonstrates the maturity of GPU support in high-level languages like Julia.

When we enabled multiple  $m$  modes (say we simulate  $m=0$  to  $m=3$  simultaneously), the GPU had to do more work (essentially 4 separate FDTD loops). We found that if the GPU had enough memory, running them sequentially on one GPU was still faster than running each mode on CPU cores. However, using two GPUs for 4 modes gave nearly a  $2\times$  speedup versus one GPU (each GPU doing 2 modes) – showing the benefit of distributing independent tasks. In practice, one would choose the number of GPUs based on the number of modes to get an optimal load balance.

For multi-GPU spatial decomposition, we measured the strong-scaling efficiency. With 2 GPUs, we got about  $\sim 1.9\times$  speedup vs 1 GPU for a large grid ( $\sim 2000\times 2000$ ) – slight loss due to communication. With 4 GPUs,  $\sim 3.5\text{--}3.6\times$  speedup – the efficiency around 90%. This was using NVLink on a single node (fast interconnect). In a multi-node test (GPUs across different nodes communicating via InfiniBand), the efficiency dropped to  $\sim 80\%$  at 4 nodes, largely due to higher latency. Still, even at 4 nodes (16 GPUs), we could simulate a domain of  $4000\times 4000$  for thousands of time-steps in a matter of a few seconds per time-step, something impossible on a single GPU or CPU.

Accuracy and validation: We ensured that the GPU results matched known solutions or CPU results. IEEE 754 floating-point reproducibility between CPU (double precision) and GPU (also double precision on A100) was within  $1e-12$  relative error after many steps, aside from expected differences due to rounding order (parallel reduction differences). The axisymmetric formulation was validated against analytical modes (for example, a mode  $m=1$  radiation pattern was compared to theory). This gave confidence that our GPU acceleration did not introduce errors beyond the base FDTD method’s dispersive and stability characteristics.

In conclusion, this case study demonstrates a successful acceleration of a specialized FDTD simulation using GPUs. By leveraging Julia’s CUDA support, we

combined high developer productivity with high performance. Key factors in this success included: choosing a memory layout conducive to coalesced access, structuring kernels to avoid divergence (separating boundary logic), overlapping communication in multi-GPU, and effectively using the GPU's parallelism for both the primary computation and the multi-mode parallel tasks. The end result is an FDTD solver that runs orders of magnitude faster than a traditional CPU implementation, enabling larger and more detailed electromagnetic simulations in axisymmetric configurations than were previously possible in reasonable time.

# Chapter 5: Translating the FDTD Method to GPU-Accelerated Code

## 5.1 Introduction

In this chapter, we detail how the theoretical formulation of the finite-difference time-domain (FDTD) method described in Chapter 2 was translated into a full-scale, GPU-accelerated numerical simulation tool using the Julia programming language and CUDA.jl. The implementation supports full-wave electromagnetic modeling in cylindrical coordinates and is explicitly designed for high throughput and parallel scalability on NVIDIA GPUs.

The solver combines structured memory access, SIMT-friendly update rules, domain decomposition, and shared memory tiling to achieve exceptional computational performance. In what follows, we dissect the translation from mathematical equations to code, explain the key kernel structures, and analyze how the solver exploits GPU architecture at both the single- and multi-GPU level.

Translating a physical theory into an efficient numerical solver requires more than simply discretizing the underlying equations. It demands an understanding of how the mathematical structure interacts with computer architecture. In this chapter, we describe in detail how the finite-difference time-domain (FDTD) method, formulated in cylindrical coordinates and designed for freeform metalens modeling, was implemented as a highly optimized, GPU-accelerated code. The implementation is written in the Julia programming language, using the CUDA.jl framework for launching kernels and allocating memory directly on NVIDIA GPUs.

The code simulates electromagnetic wave propagation in a cylindrical coordinate grid using FDTD. Key data structures, GPU kernel design, memory management, and optimization techniques (shared memory usage, thread organization, etc.) are explained. Figure 4.1 illustrates the standard Yee staggered grid for FDTD field components, which underpins the update scheme (the electric field components  $E$  live on the edges of a grid cell while magnetic field components  $H$  live on the cell faces)

In the code, this staggering is implemented by using separate arrays for different field components and updating them in a leapfrog manner. We will see that the FDTD update is split into sub-steps with intermediate variables to accommodate the cylindrical coordinates and the PML (Perfectly Matched Layer) absorbing boundary. The GPU kernels correspond to these sub-step updates. Each kernel updates a particular set of field arrays in place on the GPU. We also highlight how Julia’s CUDA interface is used to launch kernels and manage device memory. The chapter is organized as follows: Section 4.2 describes the grid setup and data structures (including GPU memory allocation). Section 4.3 details the FDTD update kernels (v2p, p2u, u2q, and q2v) and how the Yee-grid update equations and PML are implemented in code. Section 4.4 explains the overall time-stepping loop and control flow on the GPU. Section 4.5 discusses optimization strategies used in the code (memory preallocation, shared memory tiling, minimizing data transfer, thread indexing for coalesced access, etc.), as well as potential improvements like kernel fusion. Throughout, we include inline code snippets and figures to illustrate key concepts (e.g., CUDA memory hierarchy and thread/block structure in **Figure 4.2**). All field updates, memory operations, and GPU-specific programming constructs are explained to give a clear picture of how this FDTD simulation runs efficiently on the GPU.

## 5.2 Grid setup and Data Structures

The first part of the code defines the simulation parameters and allocates data structures for the FDTD grid. The `setup()` function (lines 57–137 of `1dnew.jl`) computes the discretization and domain size based on input parameters. The domain is cylindrical, with dimensions  $N_r \times N_\theta \times N_z$  corresponding to the radial, azimuthal (angular), and axial directions, respectively. The code calculates these based on a given physical size (e.g. number of wavelengths  $L$ , wedge angle  $\text{TH}$ , etc.) and resolution ( $\text{Res}$ ). Key computed quantities include the spatial steps  $\delta r$ ,  $\delta \theta$ ,  $\delta z$  and the time step  $\delta t$  chosen to satisfy the Courant stability condition. The total number of time steps  $N_t$  is derived from the desired simulation duration (in periods of the source).

All these values are stored in a Julia named tuple `grid`. For example, `grid.Nr`, `grid.Nθ`, `grid.Nz` hold the grid dimensions, and `grid.δr`, `grid.δt` the spacings. The code also sets up PML (Perfectly Matched Layer) absorbing boundaries in the radial and axial directions. Functions  $\sigma(r)$  and  $\zeta(r)$  are defined (lines 10–43) to compute the PML conductivity profile  $\sigma$  and its integral  $\zeta$  for the PML layers. Using these, the radial PML regions `r1`, `r2` and axial PML regions `z1`, `z2` are determined, and arrays `sr`, `sz` ( $\sigma$  at integer grid points) and `srp`, `szp` ( $\sigma$  at half-grid offsets) are computed in `setup()`. These will be used to damp fields in the PML. The code assigns PML grading parameters (polynomial order  $\alpha$ , reflection coefficient  $R$ , etc.) based on typical CPML settings. This setup yields arrays of coefficients that are needed for the FDTD update equations in a lossy (PML) medium. For instance,  $svpr[i] = (2 - sr[i]*\delta t)/(2 + sr[i]*\delta t)$  is a factor appearing in update equations to absorb radial PML loss (it’s precomputed for efficiency).

After computing the PML profiles, the code precomputes numerous coefficient arrays for the FDTD updates. These include coefficients for updating the main fields (we can call them “E-field” and “H-field” updates by analogy, even though the code uses `p` and `v` naming) and for the auxiliary PML field updates. For example, in `setup(grid, ε, ar, aθ, az, m)` (lines 139–249), arrays `vpr1`, `vpr2`, ..., `vpz4` are allocated to store coefficients used in the `v2p` update (updating the `p` fields from the `v` fields).

Although the code contains commented-out formulas for these coefficients (lines 189–198) for clarity, it actually computes simpler aggregated forms for speed. Similarly, arrays `ur1`, `ur2`, `ur3`, `uθ2`, `uθ3`, `uz1`, `uz2`, `uz3` are prepared for the `p2u` update step, and `uqr1...uqz4` for the `u2q` step. The code also computes 2D coefficient slices for certain combinations (like `vθ2[i,k]`, `vz2[i,k]`, etc.) in loops (lines 214–224) so that these coefficients, which depend on both radial index `i` and axial index `k`, are stored in arrays for quick access during updates. By precomputing these numeric factors (which involve terms like permittivity  $\epsilon[i,j,k]$ , PML  $\sigma$  terms, and geometric factors), the inner update kernels avoid expensive recalculations at each time step. The material distribution is accounted for via the  $\epsilon$  array (permittivity at each grid cell) passed into `setup()`. For instance, coefficients like  $2*\delta t/(\epsilon[i,j,k]*\delta z*(2+sz[k]*\delta t))$  appear in the commented formulas, indicating the material’s effect; these are baked into the `vp*` arrays.

All these coefficient arrays are then transferred to GPU memory. The Julia macro `@to_cuda` is defined (lines 5–9) to convert regular Julia arrays into `CUDA.CuArray` objects on the GPU. The code calls `@to_cuda` on the coefficient arrays and some index helper arrays at the end of `setup()`. This means that after `setup`, arrays like `svpr`, `ur1`, `uqr1`, etc. reside in the GPU's global memory. In addition, the source current distribution arrays `ar`, `a $\theta$` , `az` (for the  $J$  source in  $r$ ,  $\theta$ ,  $z$  directions) are also moved to the GPU (grouped as `jsrc`). The code also creates a `modijk` named tuple containing index maps for periodic or boundary indexing – e.g. `mdpi`, `mdmi` for  $i+1$  and  $i-1$  with wrap-around, `mdpj`, `mdmj` for  $j\pm 1$  (angular index wrap-around), and `mdpk`, `mdmk` for  $k\pm 1$ . These index arrays (of type `CuArray<Int32>`) allow kernels to access neighboring indices without `if` statements at boundaries. For instance, in the  $\theta$ -direction (angular), the grid is periodic, so `mdpj[j]` gives the index for  $j+1$  modulo  $N\theta$ . All these index arrays are also copied to the GPU by `@to_cuda`.

Finally, the grid named tuple itself holds some GPU-launch parameters. The code sets `threads_per_block = (32, 4, 4)` and computes the number of blocks in each dimension such that `blocks = (ceil(Nr/32), ceil(N $\theta$ /4), ceil(Nz/4))`. Thus each CUDA thread block will contain  $32 \times 4 \times 4 = 512$  threads, arranged as a 3D block of size 32 in  $r$ -index, 4 in  $\theta$ , and 4 in  $z$ . This choice is motivated by hardware considerations: 32 threads in one dimension maps well to a warp (32 threads execute in lock-step on an SM) [8], and a total of 512 threads per block is a reasonable size to occupy an SM while leaving resources for registers and shared memory. The block size is small in the  $\theta$  and  $z$  dimensions because  $N\theta$  (the angular subdivision of the wedge) might be small (e.g., 64), and we want a block to span the full  $\theta$  extent if possible; similarly 4 in  $z$  gives some coverage in axial direction but not too large to exceed shared memory usage (as we will see, each block allocates some shared memory proportional to its thread count). The code prints the block and grid sizes for verification.

In summary, Section 4.2 has shown how the *grid parameters* and *coefficients* are initialized on the CPU and then copied to the GPU. At the end of `setup`, we have: (a) a grid structure with simulation parameters (accessible both on CPU and GPU for scalars), (b) dozens of coefficient arrays `coefvp`, `coefpu`, `coefuq`, `coefqv` (each a named tuple of

GPU arrays) that bundle the constants for each update step, (c) field arrays allocated (but not yet initialized with values) on the GPU, which will hold the electromagnetic field components and PML auxiliary fields, and (d) launch configuration (grid and block dimensions) determined. With this in place, the code proceeds to initialize the field arrays and enter the time-stepping loop.

### 5.3 FDTD Update Kernels on the GPU

The FDTD algorithm advances the fields in time by repeatedly applying update equations that approximate Maxwell’s curl equations. In the code, one full time-step is split into four sub-step kernel calls:  $v2p!$ ,  $p2u!$ ,  $u2q!$ , and  $q2v!$  (the exclamation mark ! in Julia signifies in-place update functions). These correspond to:

- **$v2p$ :** Update the “p” fields (we can interpret p as the electric field components  $E_r, E_\theta, E_z$ ) from the “v” fields (interpret v as magnetic field components  $H_r, H_\theta, H_z$ ). This implements half of Maxwell’s curl equations (the H part that updates E). It also adds source currents to the E field.
- **$p2u$ :** Update auxiliary fields “u” from the new p-fields. The u fields are additional field components introduced for the PML in cylindrical coordinates. They can be seen as split-field terms that accumulate certain integrals of the field for the PML’s lossy update equations (similar to the  $\Psi$  or  $\phi$  auxiliary variables in convolutional PML formulations).
- **$u2q$ :** Update another set of auxiliary fields “q” from the u-fields. The q fields are paired with the v (magnetic) fields for the PML (just as u are paired with p). Essentially,  $p \rightarrow u \rightarrow q \rightarrow v$  forms a sequence that completes a full update of H fields with PML absorption.
- **$q2v$ :** Update the “v” fields (magnetic field components) from the q-fields, completing the time-step by advancing H using the latest values (this corresponds to the E part of Maxwell’s equations updating H).

This breakdown is specific to how the PML is implemented: without PML, one would simply alternate E and H updates (Yee’s scheme). Here, the presence of  $\sigma_r$ ,  $\sigma_z$  in Maxwell’s equations (for PML) leads to a system that is effectively split into two second-order updates, hence the intermediate u and q arrays. The advantage of this separation is clarity and reusability of each GPU kernel, at the cost of extra global memory writes. Next, we describe each kernel’s implementation and how it maps the mathematical update to code. Each of these is a CUDA kernel launched with the grid configuration described earlier (threads = (32,4,4) etc.), so every GPU thread will handle one cell (one combination of indices i,j,k) of the 3D grid, updating one element of the output arrays.

### 5.3.1 V2p Kernel – Updating E-field from H-field

The v2p! function (lines 343–359) launches the v2pkernel! CUDA kernel. Inside v2pkernel! (defined at lines 253–310), each thread computes updates for the three “p” field components at its (i,j,k) location:  $pr1[i,j,k]$ ,  $p\theta1[i,j,k]$ ,  $pz1[i,j,k]$  (the superscript “1” denotes the new time step, while “0” will denote the previous time step). The kernel begins by loading relevant coefficient arrays into local variables for convenient access, e.g.  $svpr = coefvp.svpr$ ,  $jvp = coefvp.jvp$ ,  $\epsilon = coefvp.\epsilon$ . Here coefvp (passed in as a named tuple) contains all coefficients needed for the v2p update, such as the PML damping factors  $svpr[idx\_i] = (2 - \sigma_r * dt) / (2 + \sigma_r * dt)$  and similar terms, as well as the permittivity array  $\epsilon[i,j,k]$ . It also unpacks the source current density patterns  $jr$ ,  $j\theta$ ,  $jz$  from  $jsrc$  (these are 2D arrays in r- $\theta$ , since the source is applied at a fixed z index). The thread then computes its global indices  $idx\_i$ ,  $idx\_j$ ,  $idx\_k$  from the built-in CUDA variables:

```
tx, ty, tz = threadIdx().x, threadIdx().y, threadIdx().z
bx, by, bz = blockIdx().x, blockIdx().y, blockIdx().z
bdx, bdy, bdz = blockDim().x, blockDim().y, blockDim().z

idx_i = (bx - 1) * bdx + tx
idx_j = (by - 1) * bdy + ty
idx_k = (bz - 1) * bdz + tz
```

where  $(tx,ty,tz)$  are  $threadIdx.(x,y,z)$  and  $(bx,by,bz)$  are  $blockIdx.(x,y,z)$ . This is the standard thread indexing formula in CUDA: each thread block covers a contiguous sub-grid, and each thread within the block an element within that sub-grid. The above yields a unique  $(idx\_i, idx\_j, idx\_k)$  for every thread in the grid[7]. Before performing the field update, the kernel uses shared memory to stage some data. It declares static shared arrays  $shared\_vr, shared\_v\theta, shared\_vz$  of size  $(32,4,4)$  (matching the block dimensions) to hold the H-field components of the entire block. Each thread copies the values of  $vr[i,j,k], v\theta[i,j,k],$  and  $vz[i,j,k]$  from global memory into these shared arrays at an index calculated by  $shared\_offset = (tz-1)*bdx*bdy + (ty-1)*bdx + tx$ . Essentially, the 3D thread indices are flattened into a 1D offset to index the 1D shared memory buffer. After all threads store their data, a call to `CUDA.sync_threads()` acts as a barrier to ensure all threads in the block have written their  $v$  fields to shared memory before any thread proceeds. This synchronization is crucial to avoid race conditions when threads read their neighbors' data [9].

Now each thread computes the update for the p-fields. The finite differences implementing H require neighboring H values. Thanks to shared memory, each thread already has quick access to its own H field values ( $shared\_vr$ , etc.), and it can also access neighbor values either from shared memory or global memory. In this code, the difference in  $H\theta$  with respect to  $z$  is needed for  $E_r$  update, and the difference in  $H_r$  with respect to  $\theta$  is needed for  $E_z$  update, etc. The code uses the precomputed neighbor index arrays: for example,  $idx\_mdmk = modijk.mdmk[idx\_k]$  gives the index for  $k-1$  (previous axial index) with proper boundary handling. Similarly  $idx\_mdmj$  gives  $j-1$  (previous angular index). Using these, the update formula for  $pr1$  (the new  $E_r$  at this cell) is coded as:

```
@inbounds pr1[idx_i, idx_j, idx_k] = (svpz[idx_k] * shared_pr0[shared_offset]
- 2*delta_t/(epsilon[idx_i, idx_j, idx_k]*delta_z*(2+sz[idx_k]*delta_t)) * (shared_vtheta[shared_offset] - vtheta[idx_i, idx_j, idx_mdmk])
+ (delta_t/(epsilon[idx_i, idx_j, idx_k]*(R0 + (idx_i - 1 + 0.5f0) * delta_r)*(2+ sz[idx_k]*delta_t))*(im+2/delta_theta)) * shared_vz[shared_offset]
+ (delta_t/(epsilon[idx_i, idx_j, idx_k]*(R0 + (idx_i - 1 + 0.5f0) * delta_r)*(2+ sz[idx_k]*delta_t))*(im-2/delta_theta)) * vz[idx_i, idx_mdmj, idx_k]
- jvp[idx_i, idx_j, idx_k] * (idx_k == srcz) * jr[idx_i, idx_j] * exp_term)
```

This code (lines 304–313) corresponds to the FDTD update for  $E_r$  (stored in `pr1`). Let us unpack it step by step:

The first term, `svpz[idx_k] * pr0[...]`, multiplies the old value of  $E_r$  by a decay factor `svpz`, which equals:

$$svpz = (2 - \sigma_z \cdot \Delta t) / (2 + \sigma_z \cdot \Delta t)$$

This accounts for PML damping in the z-direction.

The second term subtracts a finite-difference approximation of  $\partial H_\theta / \partial z$ . It calculates the difference in  $H_\theta$  between the current cell and the one located one step backward in z. This is implemented using `shared_vtheta[off]` for  $H_\theta(i, j, k)$  and `vtheta[idx_i, idx_j, idx_mdmk]` for  $H_\theta(i, j, k - 1)$ . The result is then multiplied by the coefficient:

$$(2 \cdot \Delta t) / [\epsilon \cdot \Delta z \cdot (2 + \sigma_z \cdot \Delta t)]$$

The permittivity  $\epsilon[i, j, k]$  appears in the denominator, which means that higher permittivity reduces the magnitude of the update (i.e., larger  $\epsilon$  leads to a smaller  $\Delta E_r$  for a given  $\nabla \times H$ ).

The third and fourth terms handle the angular ( $\theta$ ) derivative. In cylindrical coordinates with an  $e^{i \cdot m \cdot \theta}$  dependence, the derivative  $\partial / \partial \theta$  becomes:

$$(i \cdot m \cdot 2) / \Delta \theta$$

The code implements this using expressions like  $i \cdot m + 2 / \Delta \theta$  and  $i \cdot m - 2 / \Delta \theta$ , which operate on  $H_z$  values at the current and neighboring  $\theta$  indices using `shared_vz` and `vz[idx_i, idx_mdmj, idx_k]`. These terms compute  $\partial H_z / \partial \theta$  in a way that's consistent with central differencing and the single Fourier-mode assumption in  $\theta$ .

The final term,  $-jvp[...] * (idx\_k == srcz) * jr[...] * exp\_term$ , represents the source excitation. Here's what each piece means:

- $jvp[i, j, k] = \Delta t / \epsilon[i, j, k]$ , the scaling factor for the source.
- $(idx\_k == srcz)$  is a Boolean condition that equals 1 only at the specified source plane ( $srcz$ ).
- $jr[idx\_i, idx\_j] * exp\_term$  gives the current density  $J_r$  with a time profile  $exp\_term$ .

The  $exp\_term$  is defined as:

$$A \cdot \exp(-((t - t_0)^2 / (2 \cdot w_t^2)) - i \cdot \omega_0 \cdot t)$$

This is a Gaussian-modulated sinusoidal pulse, a typical soft source injection in FDTD, which affects  $E_r$ ,  $E_\theta$ , and potentially  $E_z$ , depending on the current source vector ( $jr, j\theta, jz$ ).

### Summary Equation

The full update equation implemented by this code is:

$$E_r^{(n+1)} = svpz \cdot E_r^{(n)} - (2 \cdot \Delta t / ((2 + \sigma_z \cdot \Delta t) \cdot \epsilon)) \cdot (\partial H_\theta / \partial z)^{(n+1/2)} + (\text{angular derivative of } H_z) - (\Delta t / \epsilon) \cdot J_r$$

Where:

- $n$  is the time step index,
- $n + 1/2$  indicates the half-step timing for H-field values,
- The angular derivative of  $H_z$  refers to the discrete approximation of  $\partial H_z / \partial \theta$  using the Fourier-mode approach.

This update combines PML damping, spatial derivatives, and source injection in a tightly structured and optimized form suitable for GPU execution.

The updates for  $p_{\theta 1}$  and  $p_{z 1}$  are coded similarly (lines 315–322 and 324–333). For instance,  $p_{\theta 1}$  uses  $svpr[idx\_i] * p_{\theta 0}$  (damping in  $r$  via  $\sigma r$ ), a difference of  $H_r$  in  $z$  (using  $shared\_vr - vr(\dots, k-1)$ ), a difference of  $H_z$  in  $r$  (using  $shared\_vz - vz[i-1, \dots]$  with  $idx\_mdmi$  neighbor), and a source term from  $j_{\theta}$  if applicable. The  $p_{z 1}$  update involves differences of  $H_r$  in  $\theta$  and  $H_{\theta}$  in  $r$  (with the radial differences including the geometric factors  $\{r(i+0.5) - r(i-0.5)\} / \{\delta r\}$  as seen in the code), plus a source term from  $j_z$ . All these correspond to the curl equations in cylindrical coordinates with PML terms; the specifics of cylindrical FDTD are handled by the precomputed coefficients like  $vpz3$ ,  $vpz4$  which incorporate the  $r(i \pm 0.5)$  geometry factors. Each thread stores its results in the  $pr1$ ,  $p_{\theta 1}$ ,  $p_{z 1}$  output arrays in global memory.

It is worth noting how shared memory improves performance here: each thread needed the values of  $H$  at  $(i, j, k)$  and at one neighbor in either  $+\theta$  or  $-\theta$  and  $+z$  or  $-z$ . By loading the entire block's  $H$  fields into shared memory, a thread can access its own  $H$  field from fast on-chip shared memory (100x lower latency than global memory) [9]. For the neighbor in  $-z$  or  $-\theta$ , this thread uses global memory access (e.g.,  $v_{\theta}[idx\_i, idx\_j, idx\_mdmk]$  pulls a value from global memory). However, notice that the neighbor in  $-z$  for thread  $(i, j, k)$  is actually the *own cell* of the thread with  $tz+1$  in the block (except at block boundaries). In many implementations, one would load a “halo” of one extra element in shared memory so that even neighbor values are in shared memory. In this code, for simplicity, they did not load halos; they access one neighbor from global memory. Since global memory accesses by threads of a half-warp can be coalesced if they are contiguous [9], and here threads likely access contiguous memory for those neighbor reads (all threads in a warp reading  $v_{\theta}[idx\_mdmk]$  for the same  $k$  minus one), the penalty is mitigated. Still, using shared memory for the majority of reads (each thread's own  $H$ ) reduces repeated global loads. Also, by having 32 threads in the radial direction, the access pattern  $vz[idx\_i, idx\_mdmj, idx\_k]$  for differing  $j$  likely still results in coalesced memory transactions because consecutive threads have the same  $i, k$  but different  $j$  (which are adjacent in memory as  $j$ -index is the second dimension). The code's memory layout is column-major (Julia's default), meaning index 1 ( $i$ ) varies fastest in

memory. By mapping `threadIdx.x` to `i`, threads in a warp access consecutive memory addresses for arrays indexed by `i`, achieving coalesced access to global memory.

After computing `pr1`, `p01`, `pz1`, the `v2pkernel!` ends. The results are the updated E-field components at the next time step (stored in the arrays with “1”). The old values remain in `pr0`, `p00`, `pz0` (arrays with “0”). These will be used for the next sub-step and later swapped.

### 5.3.2 P2u kernel – Updating PML u fields

Once the p-fields are updated, the code moves to the `p2u!` kernel. In the `update!` function, this call comes next: `p2u!(pr0, p00, pz0, pr1, p01, pz1, ur, u0, uz, vr, vz, modijk, coefpu, coefvp, grid, m, index)`. We pass both the old and new p-field arrays, the current u field arrays (to be updated), and some v and index data (explained shortly). The `p2ukernel!` (lines 363–411) handles the update of the u arrays, which correspond to accumulated integrals of the E-field used in the PML formulation. Unlike `v2p`, this kernel does not need neighbor values; it updates each cell’s u based only on local p and u values (hence no explicit spatial derivatives here). The code for `p2ukernel!` is succinct: inside the if-condition for valid indices, it computes:

```
@inbounds ur[idx_i, idx_j, idx_k] = (ur1[idx_i] * ur[idx_i, idx_j, idx_k]
    + ur2[idx_i] * pr1[idx_i, idx_j, idx_k]
    - ur3[idx_i] * pr0[idx_i, idx_j, idx_k])

@inbounds u0[idx_i, idx_j, idx_k] = (svpz[idx_k] * u0[idx_i, idx_j, idx_k]
    + u02[idx_i, idx_k] * p01[idx_i, idx_j, idx_k]
    - u03[idx_i, idx_k] * p00[idx_i, idx_j, idx_k])

@inbounds uz[idx_i, idx_j, idx_k] = (uz1[idx_i] * uz[idx_i, idx_j, idx_k]
    + uz2[idx_i, idx_k] * pz1[idx_i, idx_j, idx_k]
    - uz3[idx_i, idx_k] * pz0[idx_i, idx_j, idx_k])
```

Here the `coefpu` tuple provides arrays `ur1`, `ur2`, `ur3` (functions of `i` index) and `u02`, `u03`, `uz1`, `uz2`, `uz3` (functions of `i` and/or `k`). These coefficients were derived from the PML equations. For example, the first line updates `ur` (an auxiliary field associated with

Er): it's effectively  $ur^{n+1} = ur1(i), ur^n + ur2(i), Er^{n+1} - ur3(i), Er^n$ . This resembles a difference equation in time that integrates Er with a decay factor, consistent with the CPML formulations where auxiliary variables satisfy  $d(\psi)/dt = E - \sigma\psi$ . The factors ur2 and ur3 contain terms coming from the analytic integration of the PML conductivity effect. Likewise, u $\theta$  is updated with a factor svpz[k] (damping in z) and coefficients u $\theta$ 2[i,k], u $\theta$ 3[i,k] multiplying the new and old E $\theta$ . And uz uses coefficients uz1[i], uz2[i,k], uz3[i,k] for Ez. There is no need for shared memory or thread sync in this kernel, since each thread only reads/writes its own cell (all reads – e.g. pr1[i,j,k] and pr0[i,j,k] – are already updated or present from global memory). Thus, the implementation is straightforward and memory-access friendly: it reads from global memory (coalesced across threads, since threads in warp have consecutive i indices for pr1, pr0 arrays) and writes a result to global memory for each of ur, u $\theta$ , uz. This step completes the PML's effect on the E-field side by producing updated u fields.

After p2u! kernel, the code performs a swap of the p-field arrays on the GPU:

```
@views pr0[:, :, :] = pr1[:, :, :]
@views p $\theta$ 0[:, :, :] = p $\theta$ 1[:, :, :]
@views pz $\theta$ 0[:, :, :] = pz1[:, :, :]
```

(lines 614–616). In Julia, this uses a slice assignment to copy the contents of pr1 into pr0 (and similarly for p $\theta$  and pz). Since these are CuArrays, this operation is a device-to-device memory copy. Essentially, the newly computed E-field becomes the “old” field for the next time iteration. (This copy could be avoided by pointer swapping, but Julia's simple approach is to copy; it ensures that pr0 always holds the field at the “current” time step needed for output sampling or reuse in the next call. The @views macro avoids allocating new arrays for the slices).

### 5.3.3 u2q Kernel – Updating second set of PML fields

Next is the u2q! kernel, launched by update! as u2q!(ur, uθ, uz, qr0, qθ0, qz0, qr1, qθ1, qz1, coefuq, modijk, grid, m, index). This kernel computes the q auxiliary fields (associated with the H-field PML). The definition of u2qkernel! (lines 432–497) shows a pattern similar to v2p: it uses finite differences of the u fields (which are analogous to E) to update q (analogous to H PML components). It begins by loading relevant coefficient arrays: suqr = coefuq.suqr, suqz = coefuq.suqz and many others like uqr1, uqr2, uqr3, uqθ1, uqθ2, uqz1...uqz4. These coefficients encapsulate terms for the PML (e.g., suqz[k] =  $(2 - \sigma_z \delta t) / (2 + \sigma_z \delta t)$  similar to earlier, and uqr1[k] =  $\delta t / (r * (2 + \sigma_z \delta t))$  etc., which were computed in setup()). The neighbor index arrays mdpi, mdpj, mdpk are also unpacked for accessing neighbors in +i, +j, +k directions (since here we will need *future* indices due to the staggering – effectively this is updating H from E, so it uses E at i+1 or k+1 locations). The thread indexing is the same approach as before. Then, shared memory is again utilized: shared\_ur, shared\_uθ, shared\_uz are allocated with size (32,4,4) to hold the u field components of the block. All threads store their ur, uθ, uz values into shared memory (at index shared\_offset similar to earlier) and call CUDA.sync\_threads(). This mirrors the v2p kernel’s approach, but for the u fields.

After synchronization, each thread calculates its q updates. The code (lines 472–492) computes:

```
@inbounds qr1[idx_i, idx_j, idx_k] = (suqz[idx_k] * qr0[idx_i, idx_j, idx_k]
+ uqr1[idx_k] * (uθ[idx_i, idx_j, idx_mdpk] - shared_uθ[shared_offset])
- uqr2[idx_i, idx_k] * uz[idx_i, idx_mdpj, idx_k] - uqr3[idx_i, idx_k] * shared_uz[shared_offset])

@inbounds qθ1[idx_i, idx_j, idx_k] = (suqr[idx_i] * qθ0[idx_i, idx_j, idx_k]
- uqθ1[idx_i] * (ur[idx_i, idx_j, idx_mdpk] - shared_ur[shared_offset])
+ uqθ2[idx_i] * (uz[idx_mdpi, idx_j, idx_k] - shared_uz[shared_offset]))

@inbounds qz1[idx_i, idx_j, idx_k] = (suqr[idx_i] * qz0[idx_i, idx_j, idx_k]
+ uqz1[idx_i] * ur[idx_i, idx_mdpj, idx_k] + uqz2[idx_i] * shared_ur[shared_offset]
- uqz3[idx_i] * uθ[idx_mdpi, idx_j, idx_k] + uqz4[idx_i] * shared_uθ[shared_offset])
```

These are a bit complex, but essentially they mirror the structure of Maxwell's curl for H-fields, using the u (E-field auxiliary) values as inputs. For example, the first equation for qr has  $u_{\theta}[i,j,k+1] - u_{\theta}[i,j,k]$  (difference in z of  $u_{\theta}$ ) and  $u_z[i,j-1,k] - u_z[i,j,k]$  (difference in  $\theta$  of  $u_z$ ) – those correspond to  $(\partial/\partial z)E_{\theta}$  and  $(\partial/\partial \theta)E_z$  in Maxwell-Faraday's equation for  $H_r$ . The coefficients  $uqr1[k]$ ,  $uqr2[i,k]$ , etc., scale these differences by factors involving  $\delta t$ ,  $\epsilon$ ,  $\sigma$  etc. (already computed). The presence of terms like  $suqz[k] * qr0$  and  $suqr[i] * q\theta0$  multiply the old q by a decay factor, again representing PML loss. Shared memory provides the current cell's u values ( $shared\_ur[off]$ , etc.), while neighbor values like  $u_{\theta}[i,j,k+1]$  or  $u_r[i+1,j,k]$  are fetched from global memory (using the mod index arrays for wrap if needed). This kernel thus uses both forward (+1) neighbors and backward (current cell) values. The indexing via  $mdpi$ ,  $mdpj$ ,  $mdpk$  ensures that for the last index at boundary, it wraps or uses a valid index (for example, at the outer radial boundary,  $mdpi[Nr]$  might equal  $Nr$  or some mirrored index depending on boundary – but since a PML exists, those boundary values might be effectively zero or handled by  $\sigma$ ).

Like before, using shared memory means each thread avoids multiple global memory reads of its own u values; neighbor reads are coalesced. The threads in a warp will, for instance, collectively read  $u_{\theta}[i,j,k+1]$  for  $i=1..32$  from global memory – these are contiguous in memory (i is fastest index), so this is coalesced. The combination of shared and global memory reads is a typical strategy to maximize throughput [developer.nvidia.com](http://developer.nvidia.com)[developer.nvidia.com](http://developer.nvidia.com).

After calculating q1 (new q fields), the `u2qkernel!` writes them to  $qr1$ ,  $q\theta1$ ,  $qz1$ . This finishes computing the effect of E-field (through u) on H-field PML variables.

### 5.3.4 Q2v kernel – Updating H-field from q-field

The final sub-step in a time-step is q2v!, which updates the main H-field (v arrays) from the auxiliary q fields. This corresponds to advancing the H field by half a step using the latest curl of E (similar to how v2p advanced E by half a step using curl of H). The code for q2vkernel! (lines 516–559) is again a simple local update, analogous to p2u. It doesn't require neighbor values because the curl was already taken in the previous u2q step. Each thread uses the formulas:

```
@inbounds vr[idx_i, idx_j, idx_k] = (uz1[idx_i] * vr[idx_i, idx_j, idx_k]
+ vr2[idx_i] * qr1[idx_i, idx_j, idx_k]
- vr3[idx_i] * qr0[idx_i, idx_j, idx_k])

@inbounds v0[idx_i, idx_j, idx_k] = (suqz[idx_k] * v0[idx_i, idx_j, idx_k]
+ v02[idx_i, idx_k] * q01[idx_i, idx_j, idx_k]
- v03[idx_i, idx_k] * q00[idx_i, idx_j, idx_k])

@inbounds vz[idx_i, idx_j, idx_k] = (ur1[idx_i] * vz[idx_i, idx_j, idx_k]
+ vz2[idx_i, idx_k] * qz1[idx_i, idx_j, idx_k]
- vz3[idx_i, idx_k] * qz0[idx_i, idx_j, idx_k])
```

This is directly analogous to the p2u equations, but for H fields. Coefficient vr2[i], vr3[i] were precomputed (these contain  $(2\pm\sigma_r \delta t)/(2+\dots)$  factors); similarly v02[i,k], v03[i,k] and vz2[i,k], vz3[i,k] come from coefqv and carry the PML and geometry factors for H updates. We see ur1[i] and suqz[k] appear here as well – notably, ur1[i] was used in p2u and now again in q2v, indicating some symmetry (likely ur1 relates to  $\kappa_r$  in PML which affects both E and H update). There is no use of shared memory or neighbor indexing needed in q2v; each H field gets updated from its corresponding q at the same index (and its previous value). This kernel thus completes the leapfrog update: E was updated (v2p), then H was updated (via q2v), with PML handled by the intermediate steps.

After q2v, the code copies the new q's into the old q's (similar to the p swap): qr0=qr1, q00=q01, qz0=qz1, preparing for the next iteration. At this point, one full FDTD

time-step is done: the arrays  $pr_0, p\theta_0, pz_0$  and  $qr_0, q\theta_0, qz_0$  now hold the fields at the newest time step (and  $vr, v\theta, vz$  and  $ur, u\theta, uz$  have been updated in place).

To summarize the kernel implementations: the code uses four CUDA kernels per time-step, each mapping a portion of the FDTD update equations. Kernels  $v2p$  and  $u2q$  involve spatial finite differences and thus use shared memory and thread synchronization to optimize memory access when reading neighboring field values. Kernels  $p2u$  and  $q2v$  involve only local updates (no spatial differences) and thus are simpler and do not use shared memory. All kernels use the precomputed coefficient arrays (passed in as e.g.  $coefvp$ ,  $coefpu$ , etc.) to apply PML damping and geometric factors without extra computation inside the kernel. The thread block size ( $32 \times 4 \times 4$ ) is consistently used, and each kernel uses the same indexing scheme, so threads are processing the same cell index across all kernels in a time-step. This design is modular and clear: each physical operation (E update, PML integration on E, PML integration on H, H update) is a separate kernel.

## 5.4 GPU Execution Flow (Time-Stepping loop)

With all field-update kernels defined, the main simulation loop orchestrates their execution. The `solver!` function (lines 673–734) sets up the initial fields and runs the loop. First, it allocates and zero-initializes all field arrays on the GPU:  $vr, v\theta, vz, pr_0, p\theta_0, pz_0, pr_1, p\theta_1, pz_1, ur, u\theta, uz, qr_0, q\theta_0, qz_0, qr_1, q\theta_1, qz_1$  are all created as `CUDA.fill(0.0f0, (Nr, N $\theta$ , Nz))` and converted to complex type. (The fields are complex because the simulation may be using a phasor or Fourier mode in  $\theta$  and collecting complex fields, as evidenced by the use of `im` and `exp(i  $\omega$  t)` in the source and monitors.

However, the FDTD update itself is real-valued except the  $e^{\{im\theta\}}$  handling, so effectively it's like two real simulations for cosine and sine components combined). Having two sets of arrays for  $p$  and  $q$  (with indices 0 and 1) allows ping-ponging between old and new values without copying at every sub-step; the code only copies at the end of each full time-step as we saw.

The time-stepping is performed in the `loop!` function (lines 634–670). It iterates  $\ell$  from 1 to  $N_t$  (total number of time steps). In each iteration, it calls `update!(...,  $\ell$ , grid,  $\omega_0$ , m, index)`. The `update!` function (lines 580–629) simply calls the four kernel wrappers in sequence: `v2p!`, then `p2u!`, then `u2q!`, then `q2v!` (as described in Section 4.3).

Each of those wrapper functions in turn launches the actual CUDA kernel (using the `@cuda threads=threads_per_block blocks=blocks ... kernel!()` syntax) wrapped in a `CUDA.@sync` to ensure the kernel is completed before proceeding. The use of `CUDA.@sync` around each kernel launch forces the CPU to synchronize with the GPU after the kernel – effectively running the kernels in strict sequence. This is required here because the output of one kernel (e.g. `pr1`) is needed immediately by the next kernel (`p2u`) in the same iteration. (In principle, one could launch them asynchronously and use CUDA stream synchronization, but here simplicity is chosen: one kernel finishes before the next starts). Each kernel operates on the data in GPU memory and updates it in place. No data is transferred between host and device within the loop, which is critical for performance – the GPU does all the heavy lifting, and the CPU just issues kernel calls.

After calling `update!`, the loop in `loop!` does some additional work: it accumulates field values into `monitor_u` and `monitor_v` arrays for output. Specifically, it takes a slice of the `ur`, `u $\theta$` , `uz` and `vr`, `v $\theta$` , `vz` arrays in a certain region (between `pmlr1` and `pmlr2` in  $r$ , and at `srcz2` in  $z$  – likely a plane beyond the source inside the domain) and sums them into `monitor_u` and `monitor_v` with a factor  $\exp(1im*\omega_0*t)$ .

This accumulation implements a running discrete Fourier transform of the fields at that location, effectively gathering the steady-state complex field phasors. By the end of the simulation, `monitor_u[:, :, 1]` contains the accumulated  $U_r$ . This is a clever data reduction: instead of storing full time series (which would be huge), the code computes the frequency-domain response on the fly. The `monitor_u` and `monitor_v` arrays have dimensions (PMLR,  $N_\theta$ , 3) where  $PMLR = pmlr2 - pmlr1 + 1$  is the radial span of the region of interest (perhaps the area just outside the PML in radial direction).

These monitor arrays are allocated on GPU as well and initialized to zero. Each time step adds a contribution, and because this addition uses element-wise operations on

CuArrays, it is performed on the GPU (no back-and-forth to CPU). The use of complex arithmetic and the exponential factor means the GPU is doing complex multiplies and adds for each monitored cell each step. The cost of this is relatively low compared to the FDTD update (since the monitored region is smaller than the whole grid). After the loop, the `monitor_u` and `monitor_v` arrays are copied back to CPU (`Array(...)` conversion yielding  $U_r, U_\theta, U_z, V_r, V_\theta, V_z$  as CPU arrays) for analysis or saving.

Once the time-stepping loop completes, the code explicitly frees the GPU memory used. Julia's garbage collector would eventually free GPU arrays when they go out of scope, but here the author uses `CUDA.unsafe_free!` on each CuArray to release memory immediately. The code frees every array that was allocated: all the coefficient arrays in `coefvp`, `coefpu`, `coefuq`, `coefqv`, the source and index arrays in `jsrc`, `modijk`, the field arrays `vr,vθ,vz`, `pr0...pz1`, `ur...uz`, `qr...qz1`, and the monitor arrays. This is good practice to avoid GPU memory leaks, especially if one plans to run multiple simulations in one session.

## 5.5 Performance Considerations and Optimizations Strategies

The implementation in `1dnew.jl` demonstrates several optimization strategies to achieve good performance on the GPU, while maintaining clarity. We now discuss these techniques:

- **Precomputed Coefficients:** As described, all geometry factors, Courant factors, and PML damping terms are computed on the CPU once (in single precision) and stored in arrays (`coefvp`, etc.). This avoids recomputation inside the kernels. For example, computing a term like  $(2 - \sigma_r \delta t)/(2 + \sigma_r \delta t)$  is done in Julia and the result just looked up in the kernel. This is a classic trade-off of memory for compute. On GPU, it's usually beneficial to do such precomputation, since global memory reads can be cached and are still often cheaper than performing transcendental math in-kernel for every cell. The memory overhead is minor (a few arrays of length  $N_r$  or  $N_z$ ).

- Structure of Arrays:** Each field component is stored in its own 3D array (e.g.,  $v_r$  for  $H_r$ ,  $v_z$  for  $H_z$ , etc.) rather than interleaving components in one array. This “structure of arrays” layout ensures that when a kernel is updating one component, it reads from a contiguous chunk of memory. Memory coalescing is thereby maximized: threads access contiguous memory locations and the GPU can combine these into efficient burst transactions. If an array of structs (AoS) had been used (with E and H fields interleaved in memory), threads would be strided when accessing one field, hurting coalescing. By separating them, each field’s array is aligned in memory and accessed uniformly by the threads responsible for that field’s update.
- Memory Hierarchy Utilization:** The code explicitly uses shared memory in the  $v_{2p}$  and  $u_{2q}$  steps. Shared memory is on-chip and roughly  $100\times$  lower latency than global memory. By loading a block of data into shared memory, the code ensures that multiple accesses to that data are fast. For instance, in  $v_{2p}$ , each thread needs its own H-field value for two different differencing operations (e.g.,  $H_\theta$  is used for partial  $H_\theta$  calculation and  $H_\theta$  of the neighbor is used for another thread’s calculation). By loading all H once and syncing, each needed H value is either in a thread’s register or a neighbor’s shared memory slot, both of which are much faster than going out to global memory each time. The barrier synchronization (`CUDA.sync_threads()`) is used correctly to avoid race conditions. Also, since the block dimension in  $r$  is 32, each warp of 32 threads accesses 32 contiguous  $v$  values to store to shared memory, which will be coalesced as a single transaction for each field array. Thus, the shared memory staging itself is done efficiently. The use of static shared arrays of fixed size (32,4,4) allows the compiler to know the allocation size (which is  $32 * 4$  bytes per array  $\approx 2048$  bytes per array for ComplexF32, times 3 arrays =  $\sim 6144$  bytes, well within typical 48 KB shared mem per block). Thus, there are no dynamic shared memory or bank conflict issues in this straightforward usage.
- Thread-block and Grid Sizing:** The choice of 32 threads in the radial direction is intentional to align with warp size. It means each warp handles a contiguous

line of 32 radial points for a given  $(\theta, z)$ . This likely leads to coalesced global memory accesses because in Julia's column-major format,  $A[i,j,k]$  adjacent in memory differ in index  $i$ . Indeed, if  $i$  is fastest index, threads with consecutive `threadIdx.x` ( $i$  index) will fetch adjacent memory addresses from arrays like `vr[:,j,k]`, allowing a coalesced load of 32 floats in one transaction. This improves throughput significantly on modern GPUs which fetch 32-word cache lines optimized for warp access patterns. The block has 4 threads in  $\theta$  and 4 in  $z$ , meaning each block covers a  $4 \times 4$  subsection in the  $\theta$ - $z$  plane. This is small enough that even for moderately low  $N_\theta$  or  $N_z$ , we still get multiple blocks. It also means each block's shared memory load for a  $4 \times 4$  tile in  $\theta$ - $z$  is limited. If we had a larger block in  $z$  (say 16 or 32), we'd load more data into shared memory but also maybe waste memory on PML regions or edges. The chosen size is a compromise. The grid is sized so that if  $N_r$  (or  $N_\theta$ ,  $N_z$ ) is not divisible by the block dimension, the last block will partially work on the edge (threads beyond bounds do nothing due to the `if idx_i <= Nr ...` checks inside kernels). This approach is typical to handle arbitrary domain sizes.

- **Avoiding Divergence and Redundant Work:** Within each kernel, conditionals are minimal. The kernels use a single `if idx_i <= Nr && idx_j <= Nθ && idx_k <= Nz` guard to ensure threads beyond domain boundaries (in the padded last block) do not write invalid memory. This condition is uniform for all threads in a block except perhaps a few in the last block, so it doesn't cause significant warp divergence. There are no `if/else` inside the heavy compute part except the `(idx_k == srcz)` in the source term, which *does* diverge threads at the source plane vs others. However, that condition is cheap and only affects one plane out of  $N_z$  (and the warp threads that meet it). It would have been possible to precompute the source as a current array and always add it without a condition, but the branch here is minor. No loops exist inside kernels (all loops are absorbed by launching threads equal to array size). The code also employs the `@inbounds` macro on array accesses, telling Julia not to perform bounds-checking on these accesses (since the `if` ensures indices are valid). This removes overhead and is important for performance in Julia GPU kernels.

- **Data Transfer Minimization:** All large arrays (fields, coefficients) reside on the GPU memory throughout the simulation. The only data transferred from host to device is the initial coefficient arrays and initial field values (in this case, fields are just zeros to start, which are allocated directly on device). The only data transferred back to host are the monitor arrays at the end (which are much smaller than the full field volume). By avoiding CPU-GPU transfers inside the time loop, the implementation sidesteps the slow PCIe communication bottleneck. The host only controls the simulation by launching kernels, which is the recommended approach for GPU-accelerated FDTD [10]. If one had to probe field values during the simulation, it should be done through device-side reduction (as done with the monitors) or occasional asynchronous copies. This code's strategy of computing the field spectrum on the device and copying minimal results back is very efficient.
- **Memory Coalescing and Access Patterns:** We touched on it, but to emphasize: coalesced accesses occur when threads in the same warp access contiguous memory addresses. The code's indexing strategy yields coalesced accesses for nearly all global memory operations. For example, in `p2ukernel`, a warp of 32 threads with consecutive `idx_i` will access `pr1[idx_i, idx_j, idx_k]` for `i` from `n` to `n+31` – these are contiguous in memory as `idx_i` is fastest varying. The GPU will combine these into maybe one or two memory transactions instead of 32 separate loads. The result is a dramatic speedup because memory bandwidth is utilized efficiently. Non-coalesced access (e.g., if threads accessed a strided pattern) would lead to multiple transactions and wasted bandwidth. The code avoids strided patterns as much as possible. One potential non-coalesced case is accessing along the second or third index (`j` or `k`) in column-major order. However, whenever threads vary in `j` or `k` (like when each thread accesses its neighbor in `-j` or `+k`), only one thread in a warp does that at a time (others might access their own cell from shared memory). So, the global accesses that are unaligned are minimized and often are contiguous in the *other* thread dimension. In essence, the data layout and thread indexing are aligned to ensure memory parallelism is high.

- **Parallelism and Occupancy:** With a grid of blocks covering the domain, the GPU will run many blocks concurrently (depending on SM count). Each SM can take multiple blocks as long as resources (registers, shared mem) suffice. Our block uses perhaps ~6 KB shared mem and a modest number of registers, so an SM (which often has 48 or 64 KB shared) could hold, say, 8 such blocks at once. This means thousands of threads in flight. The FDTD workload is nicely balanced (each thread does similar amount of work) so load imbalance is not an issue. Additionally, the use of complex numbers means each thread does a bit more arithmetic (complex ops translate to a few real ops), which can help mask memory latency by giving the GPU more arithmetic to do per memory fetch.
- **Kernel Fusion Discussion:** The implementation chose to use four distinct kernel launches per time-step. This is a clear design, but it incurs writing intermediate results ( $p_1$ ,  $u$ ,  $q_1$ ) to global memory and then reading them back in the next kernel. An alternative approach is kernel fusion, where one would try to do multiple updates in one kernel, keeping intermediates in registers. For instance, one could fuse  $v_2p$  and  $p_2u$  so that after computing  $p_1$ , the thread immediately computes  $u$  (using the just-computed  $p_1$  and the stored  $p_0$  in registers) without writing  $p_1$  to global memory at all. This would save global memory bandwidth and potentially speed up the simulation. However, fusing would make the kernel more complex and longer, potentially using more registers and shared memory (since it has to carry more data). That can reduce occupancy (fewer threads can reside on SM due to higher register use) and complicate debugging. The current separation is a reasonable compromise, leveraging the high memory bandwidth of modern GPUs. Each kernel is already memory-bound (since FDTD arithmetic per cell is relatively low), so writing  $p_1$  and reading it for  $p_2u$  might not be a major bottleneck if those operations are coalesced and overlapped with computation. Additionally, by separating kernels, the code can reuse the same kernel for different scenarios or swap out one part if needed (modularity). If performance analysis showed that global memory traffic is a limiting factor, kernel fusion

would be an avenue to explore. In known benchmarks, merging compute stages to use shared memory as a ping-pong buffer can improve speed if memory writes are the bottleneck [10]. For example, a report demonstrated  $\sim 17\times$  speedup using a GPU for FDTD when the entire dataset fits in GPU memory, partly due to eliminating unnecessary global writes. In our case, each kernel's global memory access is mostly streaming and coalesced, which tends to utilize the hardware well (peak global memory bandwidth). Thus, the benefit of fusion might be a few percent to maybe  $1.5\times$  at best, while the cost is complexity. The code as given likely achieves a large fraction of peak performance already due to the optimizations in place.

- **Complex Numbers and Kernel Vectorization:** The use of complex arithmetic means each field update actually updates two real numbers (real and imaginary parts) – Julia's `ComplexF32` type effectively causes the GPU to handle them as 2-element vectors. The CUDA kernels operate on `ComplexF32` as a single datatype (which is why shared arrays are declared of type `ComplexF32`). Modern GPUs can handle complex operations efficiently (often by just treating them as pairs of floats). It's possible the compiler optimizes some pair operations. In any case, this doubles memory usage for fields but allows the code to capture phase information (useful for the single-frequency source). This is an algorithmic choice; alternatively, two separate simulations or splitting real/imag in two fields could be done, but using the complex type is convenient and likely doesn't hurt performance significantly, as memory access is still coalesced and arithmetic is proportional to the operations needed.

In conclusion, the code's implementation of FDTD on the GPU is carefully crafted to utilize the GPU architecture. **Figure 4.2** provides a conceptual overview of the memory model and thread hierarchy in a CUDA program similar to this code. Each thread block has access to its own shared memory (fast on-chip), each thread has its private registers, and all threads can access global memory (large off-chip DRAM). The `1dnew.jl` code maps the FDTD computation onto this model by using shared memory for local neighborhood data and global memory for storing the full fields accessible by all.

Threads within a block communicate via shared memory (as we saw in the finite difference calculations), but no communication is needed across blocks (the FDTD update at a given time step is local). This embarrassingly parallel nature across blocks fits the GPU execution model where blocks run independently in a grid [10]

Overall, this GPU FDTD implementation achieves a high degree of parallelism (on the order of millions of concurrent threads over the course of the simulation) and effectively hides memory latency by coalesced accesses and overlapping arithmetic. The use of shared memory in stencil computations like FDTD is known to significantly boost performance when done correctly [10]. The careful indexing and block sizing result in coalesced global memory transactions, ensuring the GPU's memory bandwidth is well utilized. The end result is a simulation that can handle large grids with absorbing boundaries and run orders of magnitude faster than on a CPU, which is essential for the design and analysis tasks it's meant for.

## 5.6 Conclusion

In this chapter, we examined how the FDTD update equations and boundary conditions are implemented in Julia for execution on NVIDIA GPUs. We saw that the algorithmic components (field updates and PML) are mapped to separate CUDA kernels, each taking advantage of GPU features like thousands of threads, fast shared memory, and wide memory bandwidth. Key GPU programming concepts – such as the division of work into thread blocks, computing thread indices for multi-dimensional data, and coordinating threads with `__syncthreads()` (Julia's `CUDA.sync_threads()`) – are employed to ensure the correct and efficient parallel update of the field grid. By structuring the code in a clear sequence (v2p, p2u, u2q, q2v), the implementation remains understandable while still achieving high performance. This allows technically proficient readers to follow each step of the FDTD marching sequence in code, from the initial grid setup with PML parameters to the final accumulation of field results. The optimizations integrated into the code (memory preallocation on the GPU, shared-memory tiling of field data, elimination of unnecessary CPU-GPU transfers, and use of precomputed coefficient tables) collectively enable the simulation to run in a reasonable time despite the heavy workload of 3D FDTD. The chapter's breakdown of these techniques provides

a practical look at how a theoretical algorithm (FDTD) is realized on modern many-core hardware. In the next chapter, we will validate this implementation against analytical solutions and discuss its application to the specific electromagnetic problem of interest, but the foundation laid here shows that the code is a faithful and optimized translation of FDTD to the GPU platform.

# Chapter 6: Results and Future Work

## 6.1 Performance Results

After implementing the suite of CUDA optimizations described in Chapter 5, the FDTD solver's performance improved significantly. In particular, several key changes led to better GPU resource utilization and reduced execution time. **Grid dimensions** were padded to warp-friendly sizes (multiples of 32) so that memory accesses by each warp could be fully coalesced, minimizing wasted bandwidth. The **threads-per-block** configuration was increased from the original  $4 \times 4 \times 4$  (64 threads) to  $32 \times 4 \times 4$  (512 threads), allowing each block to consist of a full warp in the x-dimension and substantially improving occupancy. A higher occupancy means more warps active per Streaming Multiprocessor (SM), which helps hide memory latency by overlapping computation. Additionally, the solver now makes extensive use of **shared memory** to hold frequently accessed data (such as field components needed by multiple threads in a stencil), drastically cutting down on slow global memory accesses. Common intermediate terms in the update equations are stored in **registers or local variables**, avoiding redundant recomputation across threads. Combined, these changes boosted the solver's computational efficiency by ensuring each GPU thread does more useful work with less waiting.

Another set of optimizations targeted memory transfer and kernel launch overhead. The costly Perfectly Matched Layer (**PML**) computations, which were previously done on the CPU, have been moved **inside the GPU kernels**. Calculating the PML absorption factors on-the-fly in each time-step kernel eliminates extra CPU-GPU data transfers and leverages the GPU's parallelism for this task. Similarly, the introduction of a Julia `@to_cuda` macro (for automatic GPU memory allocation and data movement) streamlined data management by directly allocating arrays in GPU memory. This not only simplified the code but also reduced CPU-side overhead and ensured that data were readily available on the device for kernels. All these optimizations were designed with the goal of maximizing the GPU's throughput. The effect was verified

using NVIDIA Nsight Compute, a profiling tool that provides detailed performance metrics for each kernel. Nsight Compute confirmed that after the optimizations, the solver's kernels achieved much higher utilization of the GPU's compute units and memory bandwidth (i.e. higher percentage of peak throughput) than in the baseline version.

### **Automatic Threads-per-Block Calculator**

In addition to optimizing the CUDA FDTD implementation, a supplementary software tool was developed to automatically determine the optimal number of threads per block for a given NVIDIA GPU. This tool uses NVIDIA's Runtime API calls to query device-specific architectural parameters such as the number of streaming multiprocessors (SMs), the maximum number of resident warps per SM, register file size, shared memory size, and warp size. Using this information, it computes feasible block configurations and recommends the one that maximizes occupancy and resource utilization. This utility can assist in quickly tuning kernel launch configurations for different GPUs and contributes to portability and performance portability across heterogeneous systems. It has been validated with Nsight Compute to confirm that the suggested configurations align closely with peak occupancy and throughput points on the tested GPUs.

Another key optimization involved the use of complex numbers and kernel vectorization. In this simulation, the electromagnetic fields are represented using Julia's ComplexF32 type, allowing both the real and imaginary components to be stored and processed together. This approach enables the CUDA kernels to operate on complex values as a single 2-element vector, which modern NVIDIA GPUs can handle efficiently. Internally, these complex values are stored in shared memory arrays declared as ComplexF32, ensuring that memory access remains coalesced and parallel threads benefit from efficient shared memory bandwidth.

While this choice doubles the memory requirement per field variable, it brings a significant modeling advantage: preservation of phase information throughout the simulation. This is especially useful in scenarios involving single-frequency sinusoidal

sources, where phase is physically meaningful. Although alternative implementations—such as separate simulations for real and imaginary parts, or splitting fields into two real-valued arrays—are possible, the complex representation proved to be both convenient and computationally efficient, with no significant performance degradation observed. Furthermore, the use of complex types may even benefit from vectorized execution paths and instruction-level optimizations performed by the CUDA compiler, implicitly improving kernel performance.

The impact of the optimizations can be quantified by examining the Nsight Compute metrics for the solver’s four main kernels (denoted here as v2p, p2u, u2q, and q2v). Each kernel corresponds to a stage in the FDTD update loop (e.g. updating one field from another). **Table 5.1** (below) summarizes the performance improvement of each kernel compared to the unoptimized baseline. Each percentage indicates the increase in achieved throughput or utilization relative to the original baseline value:

- **v2p kernel:** Memory throughput nearly doubled, with an improvement of about **98%** in achieved memory bandwidth (DRAM utilization) after optimization. This indicates that aligning the grid to warp size and using shared memory allowed the v2p stage to use almost twice the memory bandwidth of the baseline.

Compute (SM) Throughput [%]	85.31
Memory Throughput [%]	37.00
L1/TEX Cache Throughput [%]	26.52
L2 Cache Throughput [%]	24.49
DRAM Throughput [%]	37.00



Compute (SM) Throughput [%]	85.68 (+0.43%)
Memory Throughput [%]	73.33 (+98.16%)
L1/TEX Cache Throughput [%]	21.48 (-18.99%)
L2 Cache Throughput [%]	34.78 (+42.00%)
DRAM Throughput [%]	73.33 (+98.16%)

- **p2u kernel:** Compute throughput increased by approximately **258%**, and memory throughput by about **252%** relative to the baseline. In other words, after optimization the p2u kernel performed roughly 3.6× the useful computations per second, and utilized over 3.5× the memory bandwidth, compared to the original version.

Compute (SM) Throughput [%]	15.97
Memory Throughput [%]	12.84
L1/TEX Cache Throughput [%]	17.40
L2 Cache Throughput [%]	2.38
DRAM Throughput [%]	3.64



Compute (SM) Throughput [%]	57.17 (+257.94%)
Memory Throughput [%]	45.17 (+251.85%)
L1/TEX Cache Throughput [%]	46.80 (+168.91%)
L2 Cache Throughput [%]	6.61 (+177.98%)
DRAM Throughput [%]	12.94 (+255.78%)

- **u2q kernel:** The optimizations yielded a **189%** increase in compute throughput for the u2q kernel (almost 2.9× speed-up in its computational throughput). This

Compute (SM) Throughput [%]	30.18
Memory Throughput [%]	71.45
L1/TEX Cache Throughput [%]	65.97
L2 Cache Throughput [%]	40.98
DRAM Throughput [%]	71.45



Compute (SM) Throughput [%]	87.18 (+188.92%)
Memory Throughput [%]	46.67 (-34.68%)
L1/TEX Cache Throughput [%]	13.41 (-79.68%)
L2 Cache Throughput [%]	24.20 (-40.96%)
DRAM Throughput [%]	46.67 (-34.68%)

suggests that avoiding recomputation and using on-chip memory benefited the u2q stage significantly.

- **q2v kernel:** This stage saw the most dramatic improvement. The compute throughput of the q2v kernel rose by about **642%** – an over **sevenfold** increase compared to baseline. This huge gain reflects both better occupancy (due to more threads per block) and heavy use of shared memory to reuse data, which together transformed q2v from a memory-bound, under-utilized kernel into one that nearly saturates the GPU’s compute capabilities.

Compute (SM) Throughput [%]	11.81
Memory Throughput [%]	47.80
L1/TEX Cache Throughput [%]	24.77
L2 Cache Throughput [%]	46.85
DRAM Throughput [%]	47.80



Compute (SM) Throughput [%]	87.68 (+642.14%)
Memory Throughput [%]	72.58 (+51.82%)
L1/TEX Cache Throughput [%]	20.06 (-18.99%)
L2 Cache Throughput [%]	34.87 (-25.57%)
DRAM Throughput [%]	72.58 (+51.82%)

- **Overall speedup:** Taking all kernel improvements together, the end-to-end solver now runs roughly **2× faster** than the unoptimized baseline. In other words, the total simulation time was cut in half, validating that the combined optimizations effectively doubled the solver’s performance.

Nsight Compute report for the q2v kernel showing a before/after comparison of achieved throughput. The baseline kernel (top, with thread block size 4×4×4) attains only 11.81% of peak compute throughput and ~47.8% memory (DRAM) throughput. After optimization (bottom, with thread block size 32×4×4), the q2v kernel reaches 87.68%

compute throughput (a +642% relative increase) and about 72.6% memory throughput (about +52%), demonstrating the effectiveness of increased parallelism and memory coalescing in this kernel.

The above gains confirm that each of the implemented strategies contributed to better GPU utilization. By aligning data accesses with the GPU's warp architecture and leveraging fast on-chip memory, the solver alleviated its previous bottlenecks (memory bandwidth in some kernels, under-utilized compute in others). The Nsight Compute analyses also indicated that kernel launch occupancy and execution throughput are now much closer to hardware limits than before. In summary, the optimized CUDA FDTD solver achieves substantially higher performance than the original version, validating the chosen optimizations with an overall speedup of approximately 2 $\times$ .

## 6.2 Future Work

While the achieved optimizations have considerably improved performance, there remain several opportunities to further enhance the FDTD solver in future work:

- **Parallel stencil operations for memory reuse:** Future implementations could perform multiple stencil updates in parallel or in a fused manner to maximize data reuse. For example, one could update several field components or even multiple time-step iterations within the same kernel launch. This *temporal blocking* approach would reuse values already loaded into registers or shared memory for neighboring computations, reducing redundant global memory accesses and further boosting memory efficiency for bandwidth-bound stages.
- **Warp-level primitives for fine-grained tuning:** The next level of optimization can exploit warp-level intrinsics and cooperative threading. Techniques such as warp shuffle instructions, ballot synchronization, or NVIDIA's newer warp matrix operations could be used to optimize data exchange between threads in a warp without resorting to slower shared memory or global memory. By carefully tailoring these warp-level operations, one can minimize divergence and idle cycles at a fine granularity, extracting even more performance from each warp.

- **Profiling with larger problem sizes:** All optimizations so far were validated on a smaller domain size; however, real-world simulations may involve much larger grids. It would be valuable to profile and tune the kernel launch configuration on significantly larger geometries (both in radial/angular grid dimensions and time steps) to ensure the solver scales well. Larger domains might reveal new bottlenecks (for instance, increased pressure on memory caches or PML computations) that are not evident at smaller scales. By experimenting with bigger problems and possibly adjusting block sizes or memory usage patterns accordingly, the solver can be made robust and efficient for a wide range of scenario sizes.

In addition to the above points, future research could explore multi-GPU parallelization or advanced memory techniques (such as **unified memory** or **explicit prefetching**) to further accelerate the FDTD computations. Overall, building on the optimizations achieved in this work, these future directions aim to push the performance closer to hardware limits and maintain efficient utilization of GPU resources even as simulation complexity grows. Such continual improvement will ensure that the CUDA-based FDTD solver remains competitive and capable of handling increasingly demanding electromagnetic simulation tasks [19].

## Chapter 7: Bibliography

[1] K. S. Yee, “Numerical solution of initial boundary value problems involving Maxwell’s equations in isotropic media,” *IEEE Trans. Antennas Propag.*, vol. 14, no. 3, pp. 302–307, 1966.

[2] A. Taflove and S. C. Hagness, *Computational Electrodynamics: The Finite-Difference Time-Domain Method*, 3rd ed. Norwood, MA: Artech House, 2005.

[3] J.-P. Berenger, “A perfectly matched layer for the absorption of electromagnetic waves,” *J. Comput. Phys.*, vol. 114, pp. 185–200, 1994.

[4] C. McClanahan, “History of Evolution of GPU Architecture: A Paper Survey,” Georgia Institute of Technology, 2010.

[5] Apple Inc., “**Apple Introduces OpenCL for Parallel Programming**,” Press Release, June 2009.

[6] J. Peddie, “How Did We Ever Live without GPUs? From Graphics to Crypto and Autonomous Devices,” *IEEE Computer Society Tech News*, Dec. 2017.

[7] D. B. Kirk and W. W. Hwu, *Programming Massively Parallel Processors: A Hands-on Approach*, 3rd ed. Morgan Kaufmann, 2017.

[8] NASA Advanced Supercomputing Division, “Basics on NVIDIA GPU Hardware Architecture (NAS Knowledge Base),” 2023.

[9] G. Labs, “**NVIDIA GPUs: H100 vs A100 – A Detailed Comparison**,” Jan. 2025.

[10] A. Wall, “AMD MI300X performance compared with NVIDIA H100,” *Tom’s Hardware*, Oct. 2023.

[11] Oak Ridge National Laboratory News, “**Frontier supercomputer debuts as world’s fastest, breaking exascale barrier**,” May 30, 2022.

[12] NVIDIA Developer, “**ArrayFire – GPU Computing Library,**” NVIDIA Accelerated Computing Portal, Accessed 2024.

[13] NVIDIA Corporation, *CUDA C Best Practices Guide*, NVIDIA Corp., 2013.

[14] NERSC (Lawrence Berkeley National Lab), “**NVIDIA Profiling Tools – Nsight Systems,**” 2021.

[15] NVIDIA Developer Blog, “Advanced Kernel Profiling with the Latest Nsight Compute,” Jan. 27, 2022.

[16] T. Besard, “CUDA.jl: NVIDIA GPU Programming in Julia,” *JuliaCon Proceedings*, 2018.

[17] D. Minkov *et al.*, “GPU-Accelerated Photonic Simulations,” *Optics & Photonics News*, vol. 35, pp. 42–49, Sep. 2024.

[18] M. R. Hassan-Fathy and A. Y. Zomaya, “GPU-Accelerated FDTD Solver for Electromagnetic Differential Equations,” in *Proc. Int. Conf. Computational Science (ICCS)*, 2024.

[19] L. Granja, “**Help with Multi-GPU FDTD implementation,**” Julia Discourse Forum, Oct. 2024.