

# Optimizing Deep Learning Efficiency: A Diagrammatic Approach to Understanding FlashAttention

## 1. Executive Summary: Optimizing Deep Learning with a Visual Language

The field of artificial intelligence, particularly the subfield of deep learning, has witnessed remarkable advancements in recent years, powering applications ranging from sophisticated language models to intricate image recognition systems. However, these advancements come at a cost: deep learning models are often computationally intensive and require significant energy and memory resources. Imagine trying to assemble a complex piece of furniture without a clear set of instructions; the process would be inefficient and prone to errors. Similarly, optimizing the complex inner workings of deep learning systems can be challenging. New research offers a visual "map" for these intricate systems, specifically for "deep learning," which underpins technologies like ChatGPT. This map uses simple diagrams to illustrate how information flows, making it easier to identify ways to accelerate AI performance and reduce energy consumption. This approach has even been applied to understand and potentially enhance an already highly efficient AI technique known as "FlashAttention".<sup>1</sup>

The paper titled "FlashAttention on a Napkin: A Diagrammatic Approach to Deep Learning IO-Awareness" by Vincent Abbott and Gioele Zardini from MIT introduces an innovative method grounded in Neural Circuit Diagrams to improve the analysis of deep learning algorithms with a specific focus on input/output (I/O) efficiency on Graphics Processing Units (GPUs).<sup>1</sup> This work directly addresses the critical challenge of optimizing deep learning algorithms, particularly concerning the performance gains achievable through a more profound understanding and representation of these algorithms. A core contribution of this research is the extension of existing Neural Circuit Diagrams to incorporate detailed information about resource utilization and the distribution of computational tasks across the various levels of a GPU's architecture. This visual representation provides a clear pathway to understanding how data moves through the system, thereby enabling a more effective identification of strategies to maximize both speed and efficiency.<sup>1</sup> Furthermore, the paper demonstrates the practical application of this diagrammatic approach by analyzing and deriving insights into existing algorithms such as FlashAttention. This analysis reveals how the new methodology offers novel perspectives and frameworks for evaluating performance, potentially uncovering previously unseen opportunities for optimization.<sup>1</sup> Ultimately, this research establishes a robust foundation for a more principled study of deep learning optimizations, suggesting a shift from traditional,

often ad-hoc methods towards a more scientific and systematic approach.<sup>1</sup>

This research highlights that intricate deep learning optimizations can be simplified and made more comprehensible through visual representations. This newfound clarity has the potential to expedite the development and discovery of more resource-efficient algorithms. The very act of visualizing the complex interplay of data and operations can reveal inherent inefficiencies or potential improvements that might remain hidden within lines of code or mathematical equations. Moreover, by making the optimization process more transparent and visual, this research could broaden access to advanced optimization techniques. Traditionally, this domain requires a deep understanding of both the algorithms and the underlying hardware, often acting as a barrier for many practitioners. A visual language that simplifies this process could lower this barrier, empowering a wider range of individuals and teams to identify and implement crucial optimizations. The "napkin" metaphor used in the title further emphasizes this potential for broader accessibility and ease of understanding. The success of this diagrammatic approach in analyzing FlashAttention, a state-of-the-art optimization, also suggests its potential applicability to other complex computational systems beyond the realm of deep learning. The fundamental challenges of resource usage optimization and efficient data flow are not unique to AI; they are prevalent in various complex systems, from logistical networks to intricate manufacturing processes. If the diagrammatic language developed in this research can effectively model and optimize these diverse systems, its impact could extend far beyond the field of artificial intelligence. The grounding of this new language in category theory<sup>2</sup>, a highly abstract branch of mathematics concerned with structures and their relationships<sup>8</sup>, further suggests a potential for broad applicability and a strong theoretical underpinning.

## **2. Introduction: The Ever-Growing Demands of Deep Learning**

Deep learning models have become indispensable components of numerous cutting-edge AI applications, serving as the core technology behind large language models (LLMs) like ChatGPT, sophisticated image generation systems such as Midjourney, and a multitude of other complex coordinated systems.<sup>2</sup> These models have demonstrated remarkable proficiency in tasks that were once considered the exclusive domain of human intelligence, including accurate object detection in images, nuanced speech recognition, and fluent language translation.<sup>18</sup> The power of these AI systems stems from the sheer scale and complexity of the underlying deep learning models, which often comprise billions, and in some cases, even trillions of interconnected parameters. These parameters are meticulously trained on massive datasets, allowing the models to learn intricate patterns and relationships within the

data.<sup>2</sup>

The sheer scale of these state-of-the-art AI systems translates directly into an immense demand for computational resources.<sup>20</sup> Training these behemoth models necessitates the use of specialized hardware accelerators, most notably Graphics Processing Units (GPUs) and, increasingly, Tensor Processing Units (TPUs).<sup>18</sup> These specialized processors are designed to handle the massive parallel computations inherent in deep learning algorithms. However, the high cost associated with acquiring and operating such advanced hardware, coupled with their often-limited availability, presents significant challenges, particularly for smaller research organizations, individual researchers, and startups with constrained resources.<sup>20</sup>

Given the substantial investment required in terms of both computational resources and time to develop and deploy these advanced AI systems, the optimization of their underlying deep learning models has become a matter of critical importance.<sup>6</sup> Optimizing these models is crucial not only for enhancing their efficiency, thereby reducing the computational burden and associated costs, but also for accelerating the pace of innovation within the field.<sup>20</sup> As deep learning models continue to grow in both size and complexity, the demand for effective strategies that can streamline their execution across a diverse range of computing platforms, from powerful data center servers to resource-constrained edge devices, is constantly escalating.<sup>50</sup>

Despite the critical need for optimization, the prevailing methods for discovering and implementing improvements in deep learning algorithms often suffer from significant limitations.<sup>2</sup> Many current approaches rely heavily on extensive trial-and-error experimentation, manual derivation of optimized algorithms, and a deep, often intuitive, understanding of both the algorithmic structure and the intricacies of the underlying hardware.<sup>2</sup> This manual process of deriving optimized algorithms is frequently slow and labor-intensive, potentially leaving substantial performance gains undiscovered.<sup>3</sup> A prime example of this is the development of the highly successful FlashAttention algorithm, which required three distinct iterations spanning over three years to reach its current level of performance.<sup>3</sup> Furthermore, automated compilation methods, which aim to automatically optimize code for specific hardware, have consistently lagged behind these manually crafted optimizations in terms of achieved performance.<sup>3</sup>

The significant time and effort invested in developing optimizations like FlashAttention underscore a critical need for more efficient and systematic methodologies within the field of deep learning. The fact that such a high-impact optimization required years of iterative refinement suggests that the current process is far from optimal. This

inefficiency likely stems from the inherent complexity of these systems and a lack of readily available tools for systematically exploring the vast space of potential optimizations. Moreover, the persistent performance gap between manually derived optimizations and those achieved by automated compilation techniques further highlights a fundamental challenge. It suggests that our current automated systems might not be effectively capturing the subtle interplay between software and hardware that human experts are often able to exploit. This gap strongly indicates the need for innovative theoretical frameworks and methodologies that can guide the development of more effective and efficient automated optimization techniques. The "FlashAttention on a Napkin" research represents a promising step in this direction, offering a novel approach to tackle these challenges.

### **3. The Bottleneck: Memory and Input/Output in Deep Learning**

As deep learning models have grown in size and complexity, a critical performance bottleneck has emerged, often referred to as the "memory wall".<sup>76</sup> This term signifies that the rate at which data can be moved to and from the processing units, primarily GPUs, has become a more significant limiting factor than the raw computational speed of the processors themselves.<sup>76</sup> In essence, the performance of deep learning algorithms is increasingly constrained by the efficiency of data transfer – the Input/Output (I/O) operations – rather than solely by the number of computations that can be performed per second.<sup>76</sup> This shift in the performance bottleneck is a crucial consideration for optimizing deep learning models.

Deep learning computations are predominantly executed on GPUs, which are characterized by a sophisticated hierarchical memory architecture specifically designed to facilitate parallel processing.<sup>1</sup> This hierarchy typically comprises several levels of memory with varying characteristics. At the top is a large but relatively slow High Bandwidth Memory (HBM), which provides the main storage for the vast datasets and model parameters.<sup>78</sup> Closer to the processing cores is a smaller but significantly faster on-chip Static Random-Access Memory (SRAM).<sup>78</sup> Registers, the fastest and smallest memory units, reside within the processing cores themselves and are used to hold the data being actively manipulated.<sup>82</sup>

The transfer of data between these different levels of memory, particularly the movement of data between the slower HBM and the faster SRAM, incurs substantial performance overhead and consumes a significant amount of energy.<sup>1</sup> Studies have indicated that the bandwidth required for these data transfers can account for a considerable portion, estimated at around 46%, of the total energy expenditure associated with GPU computations in deep learning.<sup>4</sup> This stark statistic underscores

the critical importance of optimizing data movement for the development of energy-efficient and cost-effective algorithms.<sup>4</sup>

To achieve optimal performance on modern hardware, it is essential to design algorithms that are "IO-aware".<sup>1</sup> This means that the algorithm's design must explicitly consider the cost of data transfers between the different levels of the GPU's memory hierarchy. An IO-aware approach necessitates a strategic plan for effectively utilizing the available memory resources and minimizing any unnecessary movement of data.<sup>78</sup> Algorithms that fail to account for these memory transfer costs will likely suffer from performance bottlenecks, regardless of their computational complexity.

The increasing significance of memory transfer costs compared to pure computation marks a fundamental shift in the landscape of deep learning optimization. For a considerable period, the primary focus of optimization efforts was on reducing the number of floating-point operations (FLOPs) required by an algorithm. However, with the rapid advancements in hardware, particularly the exponential growth in computational power outpacing the improvements in memory bandwidth, the primary performance constraint has transitioned to the speed at which data can be shuttled between memory and processing units. Consequently, optimization strategies must now adapt to this evolving reality and prioritize I/O efficiency to fully harness the capabilities of contemporary GPUs. Furthermore, the observation that standard deep learning frameworks like PyTorch and TensorFlow often lack the necessary fine-grained control over memory access highlights a notable limitation in the current software tools available for tackling this memory bottleneck.<sup>78</sup> While these frameworks provide high-level abstractions that simplify the development and training of deep learning models, they may not expose the lower-level mechanisms required for precise management of memory access patterns. This limitation suggests a pressing need for either enhancements to these existing frameworks or the development of specialized tools and methodologies, such as the diagrammatic approach presented in the "FlashAttention on a Napkin" research, that can offer more granular control over GPU memory management and data flow.

#### **4. "FlashAttention on a Napkin": A Diagrammatic Revolution**

In response to the growing challenges of optimizing deep learning algorithms, particularly concerning the critical bottleneck of memory and I/O efficiency, researchers at the Massachusetts Institute of Technology (MIT) have pioneered a novel approach. This groundbreaking work is spearheaded by Vincent Abbott, an incoming doctoral student, and Professor Gioele Zardini, a faculty member of MIT's Laboratory for Information and Decision Systems (LIDS).<sup>1</sup> Their research centers on

the development and application of simple, intuitive diagrams, specifically termed Neural Circuit Diagrams, as a powerful tool to uncover more effective strategies for software optimization within deep learning models.<sup>1</sup> In essence, they have conceived a new "language" for describing and analyzing these intricate systems, drawing heavily on the principles of category theory.<sup>2</sup>

The details of this innovative approach are presented in their paper, "FlashAttention on a Napkin: A Diagrammatic Approach to Deep Learning IO-Awareness," which has been published in the esteemed Transactions on Machine Learning Research (TMLR).<sup>2</sup> TMLR is a peer-reviewed journal that places emphasis on the novelty and rigor of research findings, even if those findings do not necessarily represent state-of-the-art performance.<sup>3</sup> The title of the paper itself, "FlashAttention on a Napkin," is a deliberate and evocative choice, intended to convey the simplicity and accessibility of their method in addressing complex optimization problems. The researchers suggest that their approach makes the often-daunting tasks of optimizing deep learning algorithms so straightforward that they can be conceptualized and even derived through a drawing that could fit on the back of a napkin – perhaps a generously sized one.<sup>2</sup>

A primary objective of this research is to establish a more formal and principled framework for studying deep learning optimizations.<sup>1</sup> The authors aim to move away from the prevailing reliance on often-opaque manual coding and ad-hoc optimization techniques towards a more scientific and systematic methodology.<sup>1</sup> Their work has already garnered significant attention and positive feedback from leading experts in the field of AI. For instance, Jeremy Howard, a highly respected figure as the founder and CEO of Answers.ai, expressed his strong approval of the research, stating that he was "very impressed by the quality" and suggesting that this new diagramming approach could be a "very significant step" forward.<sup>2</sup> Similarly, Petar Velickovic, a senior research scientist at Google DeepMind and a lecturer at Cambridge University, described the paper as a "beautifully executed piece of theoretical research" that also manages to be highly accessible to readers, even those who may not be deeply familiar with the technical details.<sup>2</sup>

The title "FlashAttention on a Napkin" is more than just a catchy phrase; it encapsulates the transformative potential of this new method to simplify and accelerate the often-arduous process of optimizing deep learning algorithms. It hints at a fundamental shift from lengthy, intricate development cycles, often requiring years of dedicated effort, to a more intuitive and visually-driven approach. The very notion that a complex algorithm like FlashAttention could be "derived on a napkin" underscores the potential for a significant reduction in the complexity and effort typically associated with understanding and improving deep learning models.



Furthermore, the fact that this diagrammatic language is rooted in category theory suggests that it is not merely a superficial visual aid but is instead built upon a robust and highly abstract mathematical framework. Category theory, which provides a way to describe the different components of a system and their interactions in an abstract and generalized manner <sup>2</sup>, could lend significant power and generality to this method, potentially enabling formal analysis and reasoning about deep learning algorithms and their optimizations in a way that was previously not feasible.

## **5. Decoding the Diagrams: A New Language for Deep Learning**

A key innovation of the "FlashAttention on a Napkin" research is the extension of existing Neural Circuit Diagrams to create a more powerful and informative visual language for deep learning.<sup>1</sup> The authors have expanded these diagrams to explicitly include information about the usage of computational resources, such as memory and processing units, as well as the distribution of computational tasks across the different levels of a GPU's hierarchical architecture.<sup>1</sup> This enhancement is crucial as it allows for a more holistic and hardware-aware analysis of deep learning algorithms.

One of the primary strengths of this new diagrammatic language lies in its ability to provide a clear and intuitive visual representation of how data flows through the deep learning system during the execution of an algorithm.<sup>1</sup> This visualization makes it significantly easier to identify potential bottlenecks in the movement of data between different memory levels and to understand how the algorithm's structure and operations impact overall efficiency.<sup>1</sup> By tracing the path of data through the diagram, researchers and engineers can gain a more immediate understanding of the algorithm's behavior and its interaction with the underlying hardware.

The diagrammatic approach offers a level of abstraction that is specifically tailored to represent algorithms in a manner that aligns well with the architecture of multi-level GPUs.<sup>1</sup> This means that algorithms can be visually depicted in terms of the operations they perform and the data they manipulate in a way that directly corresponds to how these algorithms are executed on parallel hardware. This alignment between the visual representation and the hardware execution can facilitate the development of algorithms that are inherently more efficient on GPU architectures.

The diagrams themselves employ a distinctive alternating column structure. One set of columns represents the various data types being processed by the algorithm, while the adjacent columns depict the functions or operations that are performed on that data.<sup>1</sup> This visual convention greatly simplifies the understanding of the input and output shapes of each function and provides a clear and easy-to-follow

representation of the overall execution sequences within the algorithm. By visually tracing the data as it moves from one column to the next, one can readily understand the transformations it undergoes at each step.

Furthermore, the methodology incorporates a two-level hierarchical model to represent the different types of memory available on a GPU, such as the high-bandwidth global memory (HBM) and the faster on-chip shared memory (SRAM).<sup>1</sup> Functions and data types can be systematically organized within this hierarchical structure in the diagrams, making it straightforward to track the movement of data and the operations performed on it as it traverses the different levels of the GPU's memory hierarchy. The framework allows for the visual depiction of various data structures, including arrays and tuples, as well as the operations performed on them.<sup>1</sup>

The alternating column structure of data types and functions in these diagrams provides a very natural and intuitive way to follow the computational steps of a deep learning algorithm. It is akin to visually narrating the algorithm's execution, where the data acts as the subject and the functions as the verbs, describing the transformations that occur in sequence. This visual clarity can be particularly beneficial for individuals who may not have a strong background in the complex mathematical formalisms that often underpin deep learning algorithms, making the underlying logic more accessible. Moreover, the explicit modeling of the GPU's memory hierarchy within the diagrammatic framework underscores a significant emphasis on hardware-aware algorithm design. By visually mapping an algorithm's operations and data to the different levels of GPU memory, the diagrams can effectively highlight potential memory bottlenecks and guide optimization efforts towards a more efficient utilization of the available memory resources. Understanding this interaction is crucial for achieving high performance on GPUs, and the framework's ability to make this interaction visually apparent is a key strength.

## **6. Tackling the Complexity: How the Diagrams Simplify Optimization**

The diagrams developed in this research offer a significant advantage in understanding and optimizing deep learning algorithms by effectively representing the intricate details of parallelized operations, which are fundamental to the efficiency of these models on modern GPUs.<sup>2</sup> By visually depicting these parallel computations, the diagrams clearly illustrate the relationships between the algorithms and the parallel processing capabilities of the GPU hardware, often provided by industry leaders like NVIDIA, on which these models are typically executed.<sup>2</sup> This visual clarity allows researchers and engineers to gain a more intuitive grasp of how the algorithm



is being processed in parallel.

Unlike traditional methods of representing algorithms through code or mathematical equations, these diagrams explicitly showcase important aspects of the algorithms, including the specific computational operators being used, the estimated energy consumption associated with different operations, and the allocation of memory resources at various stages of the computation.<sup>2</sup> This explicit representation of resource usage makes it considerably easier to identify areas within the algorithm where resource consumption might be inefficient or suboptimal. By visually pinpointing these potential problem areas, the diagrams enable a more targeted approach to optimization efforts.

To further enhance memory usage and computational efficiency, the framework incorporates established techniques such as stream partitioning and group partitioning.<sup>1</sup> Stream partitioning is a strategy used to reduce memory demands by processing data in smaller, more manageable chunks or batches, while ensuring that intermediate results are kept readily accessible on the GPU's fast on-chip memory.<sup>1</sup> Group partitioning, on the other hand, focuses on optimizing both memory and computational efficiency by allowing the algorithm to divide its operations across the multiple computational cores available within a GPU.<sup>1</sup> The diagrams visually illustrate how these partitioning strategies can be applied to a given algorithm, making it easier to understand their impact on the overall execution.

The diagrammatic approach also facilitates the quantification of key performance metrics, such as the total costs associated with data transfers between different memory levels and the efficiency of memory usage per computational core.<sup>1</sup> Moreover, the researchers introduce specific statistical models that can be used in conjunction with the diagrams to estimate the amount of data that needs to be transferred during the algorithm's execution.<sup>1</sup> This quantitative analysis, guided by the visual representation of the algorithm, provides valuable insights for designing algorithms that minimize unnecessary data movement, which is a critical factor in improving overall performance and energy efficiency on GPUs.

The ability to simultaneously visualize the algorithmic structure, its mapping to parallel hardware, and the associated resource consumption offers a powerful and integrated perspective that is often lacking in traditional optimization workflows. This holistic view allows for a more informed and targeted approach to developing optimization strategies. For example, a block in the diagram that represents a computationally intensive operation with high energy consumption and significant data transfer requirements might be identified as a prime candidate for optimization. Furthermore,

the explicit integration of partitioning techniques within the diagrammatic framework suggests that this research aims not only to provide a tool for analyzing existing algorithms but also a guide for designing new, more efficient ones from the outset. By visually representing how an algorithm can be decomposed and distributed across hardware resources, the framework empowers practitioners to develop algorithms that are inherently better suited for parallel execution on GPU architectures, potentially leading to a more systematic and less ad-hoc approach to high-performance deep learning model development.

## **7. The Case Study: Understanding FlashAttention Through Diagrams**

To effectively demonstrate the power and versatility of their novel diagrammatic approach, the researchers applied it as a comprehensive case study to the FlashAttention algorithm.<sup>1</sup> FlashAttention represents a significant advancement in the optimization of the "attention" mechanism, a core component of Transformer models that underpin many of the most powerful large language models currently in use.<sup>2</sup> This algorithm is renowned for its ability to significantly speed up attention computations and reduce memory usage, making it a crucial optimization for training and deploying large AI models.

Professor Zardini emphasized the remarkable simplification offered by their diagrammatic method, stating that they were able to derive the FlashAttention algorithm – which originally took over four years of dedicated research and development to achieve its impressive sixfold performance improvement over standard PyTorch implementations<sup>3</sup> – in a considerably more straightforward manner. He metaphorically described this process as being achievable "literally, on a napkin," albeit acknowledging that it might be a rather large napkin.<sup>2</sup> This dramatic simplification underscores the potential of their visual language to demystify and streamline the often-complex world of algorithmic optimization.

By applying their diagrammatic framework to FlashAttention, the researchers were able to gain novel insights into the algorithm's performance characteristics and identify potential avenues for further optimization that might not have been readily apparent through traditional analysis methods.<sup>1</sup> Their methodology provides a new and intuitive lens through which to evaluate the algorithm's efficiency and its intricate interactions with the underlying GPU hardware architecture. This visual approach allows for a deeper understanding of the factors contributing to FlashAttention's superior performance.

Furthermore, the diagrams generated using their method can effectively illustrate how

algorithms like FlashAttention leverage specific hardware features of GPUs to achieve their performance gains.<sup>4</sup> These features include techniques such as coalesced memory access, where multiple threads access memory in a coordinated way to maximize bandwidth; the utilization of specialized tensor core units that accelerate matrix multiplications, which are fundamental to deep learning computations; and the overlapping of computational tasks with data movement to minimize idle time. By visually representing how these hardware-level optimizations are employed within the algorithm, the diagrams make it easier to comprehend their contribution to the overall performance enhancement.

The successful "derivation" of FlashAttention using this diagrammatic method serves as a powerful validation of the framework's effectiveness and its potential to provide a more accessible and intuitive way to understand complex algorithmic optimizations. It suggests that the fundamental principles underpinning FlashAttention's efficiency might be more readily grasped and reasoned about through a visual representation compared to the more traditional, code-intensive development process. Moreover, the ability to analyze and potentially further optimize an already highly tuned algorithm like FlashAttention indicates the potential of this diagrammatic approach to push the boundaries of deep learning performance. It suggests that even algorithms considered to be state-of-the-art might have further room for improvement that can be revealed through this novel method of analysis.

## **8. Technical Insights: Performance Modeling and Hardware-Aware Algorithms**

A significant technical contribution of the "FlashAttention on a Napkin" research is the development of performance models directly derived from the Neural Circuit Diagrams.<sup>1</sup> These models are meticulously constructed to account for critical factors that influence the performance of deep learning algorithms on GPUs, including the effects of quantization (the process of reducing the precision of numerical representations to save memory and potentially increase speed) and the complexities inherent in the hierarchical memory architectures of modern GPUs.<sup>1</sup> By leveraging the diagrammatic framework, the authors can create mathematical representations that enable a rigorous evaluation of various algorithm and hardware configurations. This modeling capability provides valuable guidance to users in making informed decisions regarding the optimization of their algorithms under different constraints and hardware setups.<sup>1</sup>

Furthermore, the paper meticulously outlines a systematic methodology that allows practitioners to transition from the abstract theoretical representation of an algorithm in the diagrams to concrete, practical pseudocode implementations.<sup>1</sup> This

step-by-step derivation process is designed to facilitate the creation of algorithms that are inherently aware of the underlying hardware architecture on which they will be executed. By guiding the translation from the visual representation to the code, this methodology helps in exploiting specific hardware features and capabilities to achieve enhanced performance and efficiency.

The high-level performance models developed within this research are specifically designed to readily incorporate the impact of multi-level GPU hierarchies on the performance of deep learning algorithms.<sup>3</sup> This consideration is of paramount importance because the efficiency of these algorithms on GPUs is heavily dependent on how effectively they utilize the different levels of memory available, such as registers, SRAM, and HBM, and how efficiently data is moved between these levels. The models allow for the analysis of how data movement and memory access patterns affect the overall execution time and resource consumption of the algorithm.

A key strength of this work lies in its establishment of a theoretical framework that explicitly links assumptions about the behavior of GPU hardware to claims about the performance of the algorithms represented in the diagrams.<sup>3</sup> This provides a more rigorous and scientific foundation for understanding and predicting how deep learning algorithms will perform on specific GPU architectures. By formalizing the relationship between hardware characteristics and algorithmic performance, this framework allows for more evidence-based optimization strategies.

The ability to generate performance models directly from the visual representation of an algorithm offers a significant advantage for rapid prototyping and optimization. Researchers and engineers can explore different algorithmic structures and hardware configurations at a higher level of abstraction, potentially identifying optimal solutions before investing substantial time and resources in detailed implementation and benchmarking. This early-stage performance prediction can greatly accelerate the development cycle and lead to more efficient resource allocation. Moreover, the systematic methodology for deriving hardware-aware pseudocode from the diagrams suggests a promising pathway towards the automation of generating highly optimized code for various GPU architectures. This automation could significantly lower the barrier to entry for developing efficient deep learning algorithms and accelerate the deployment of AI solutions on platforms with diverse hardware constraints.

## **9. Expert Perspectives and Significance**

The "FlashAttention on a Napkin" research has garnered significant attention and high praise from leading figures in the field of Artificial Intelligence, underscoring its

potential to make a substantial impact. Jeremy Howard, the esteemed founder and CEO of Answers.ai, expressed his strong approval, stating that he was "very impressed by the quality of this research" and suggesting that the new diagramming approach could represent a "very significant step" forward for the field.<sup>2</sup> Similarly, Petar Velickovic, a senior research scientist at Google DeepMind and a lecturer at Cambridge University, lauded the work as a "beautifully executed piece of theoretical research" that also excels in its accessibility to a broad audience.<sup>2</sup>

Vincent Abbott, one of the researchers behind this work, articulated his belief that the domain of optimized deep learning models is "quite critically unaddressed." He further emphasized that their diagrammatic approach is particularly exciting because it "opens the doors to a systematic approach to this problem," suggesting a fundamental shift in how deep learning optimizations are conceived and implemented.<sup>2</sup> The core of this shift lies in the paper's emphasis on moving away from the often-opaque and intuition-driven practices of manual coding and optimization. Instead, the research strongly advocates for the adoption of a more scientific methodology rooted in formal principles and visual representations.<sup>1</sup>

A particularly significant aspect of this research is its direct engagement with the increasingly critical issue of I/O-awareness in the optimization of deep learning algorithms.<sup>1</sup> As the computational capabilities of hardware continue to advance at a rapid pace, the efficiency of data transfer between memory and processing units has emerged as a primary bottleneck in achieving higher performance in deep learning tasks.<sup>76</sup> By focusing on visualizing and optimizing data flow, this research directly tackles this crucial challenge.

The strong positive feedback from prominent experts in both the industry and academic spheres underscores the novelty and potential transformative impact of this diagrammatic approach on a long-standing and critical problem within deep learning. Their endorsements lend significant credibility to the research and suggest that the "FlashAttention on a Napkin" framework addresses a substantial unmet need for a more systematic and understandable methodology for optimization. Furthermore, the emphasis on transitioning from "tinkering" – an approach often characterized by trial and error – towards a more "scientific methodology" reflects a growing recognition within the AI community of the need for more rigorous and explainable approaches to tackling complex computational problems. This shift has the potential to lead to more reliable, reproducible, and generalizable advancements in the optimization of deep learning models.

## 10. Future Horizons: New Avenues in Deep Learning Optimization

Looking ahead, the researchers anticipate that their newly developed diagram-based language holds significant potential for future advancements in deep learning optimization. One promising avenue is the further development of this language to enable the automation of detecting potential improvements and optimizations within deep learning algorithms.<sup>2</sup> Such automation could dramatically accelerate the process of discovering more efficient ways to train and deploy AI models, reducing the current reliance on manual identification of optimization opportunities.

Professor Zardini also highlights the prospect of leveraging a robust framework for analyzing the intricate relationship between deep learning algorithms and the underlying hardware resources to facilitate a more systematic and integrated approach to the co-design of both hardware and software.<sup>2</sup> This tight coupling between algorithm design and hardware architecture could pave the way for the development of specialized hardware platforms that are optimally tailored to the specific computational demands of certain classes of deep learning algorithms, potentially leading to significant gains in performance and energy efficiency.

While the initial focus of the research has been on the FlashAttention algorithm, the authors propose that their diagrammatic framework possesses the versatility to extend beyond attention-based models. They suggest that it could be effectively applied to represent and optimize other fundamental deep learning operations, such as those found in convolutional neural networks (CNNs).<sup>1</sup> This indicates a potentially broad applicability of their approach across the wider field of deep learning, offering a unified visual language for optimization.

By providing a more formal and principled lens for studying deep learning optimizations, this research aims to pave the way for more rigorous empirical experiments designed to test the validity of their models.<sup>1</sup> The results obtained from these experiments can then be used to drive an iterative process of refinement for both the models and the diagrammatic framework itself. This feedback loop, based on real-world performance data, promises to lead to continuous improvement and a more robust understanding of deep learning performance.

The potential for automating the detection of algorithmic improvements through this visual language represents a significant step towards accelerating the pace of innovation in the field. By reducing the need for manual discovery of optimizations, researchers and engineers could focus on exploring a wider range of algorithmic designs and hardware configurations. Furthermore, the prospect of systematic



hardware-software co-design suggests a future where the development of AI algorithms and the underlying computing infrastructure are more deeply intertwined, potentially leading to breakthroughs in performance and efficiency that are currently unattainable.

## 11. Conclusion: A Scientific Approach to GPU Optimization

The research presented in "FlashAttention on a Napkin: A Diagrammatic Approach to Deep Learning IO-Awareness" makes a significant contribution to the field by establishing a solid foundation for a more formal and principled study of deep learning optimizations, particularly in the context of GPU hardware.<sup>1</sup> This work marks a crucial step towards transitioning from the more empirical and intuition-driven approaches that have historically characterized the field to a more rigorous and scientific methodology.

The diagrammatic representation introduced in this research effectively addresses the often-opaque nature of manual coding and optimization processes in deep learning.<sup>1</sup> By providing a clear and intuitive visual language for representing algorithms and their interaction with hardware, it makes these complex processes more transparent and understandable to a broader audience of researchers, engineers, and practitioners. The research strongly advocates for a fundamental shift from traditional ad-hoc methods of deep learning optimization to a more systematic and scientific methodology.<sup>1</sup> This transition is essential for achieving more reliable, reproducible, and generalizable advancements in the field, ultimately leading to more efficient and powerful AI systems.

Ultimately, this work aims to lay a robust groundwork for a more scientific approach to the optimization of deep learning algorithms on GPUs.<sup>3</sup> By providing a theoretical framework that explicitly links assumptions about hardware behavior to claims about algorithmic performance, it sets the stage for future research where empirical experiments can be designed to test clear hypotheses. This shift promises to move the field beyond post-hoc rationalizations towards a more predictive and evidence-based understanding of deep learning performance, paving the way for continued innovation and the development of increasingly sophisticated and efficient AI technologies.

Feature	Standard Attention	FlashAttention
Time Complexity	$O(n^2)$	Remains $O(n^2)$ but with significantly reduced constant

		factor
Memory Complexity	$O(n^2)$	$O(n)$
I/O Operations	High number of reads/writes to HBM for large N	Reduced number of reads/writes between HBM and SRAM through tiling
Key Techniques	Direct computation of the attention matrix	Tiling, recomputation, kernel fusion, IO-awareness

Memory Level	Size	Speed (Bandwidth)	Use Cases
Registers	Very Small (KB)	Very High	Holding actively used data for computations
L1 Cache	Small (Tens of KB)	Very High	Caching frequently accessed data for individual cores
Shared Memory (SRAM)	Medium (Tens of KB)	High	Fast on-chip memory shared by threads within a block
L2 Cache	Medium (MB)	Medium-High	Larger cache shared by multiple SMs
High Bandwidth Memory (HBM)	Large (Tens of GB)	Medium	Main memory for storing large datasets and model parameters

Expert Name	Affiliation	Key Quote/Opinion
Jeremy Howard	Answers.ai (Founder and CEO)	"I'm very impressed by the quality of this research. The new approach to

		diagramming deep-learning algorithms used by this paper could be a very significant step... This paper is the first time I've seen such a notation used to deeply analyze the performance of a deep-learning algorithm on real-world hardware. The next step will be to see whether real-world performance gains can be achieved." <sup>2</sup>
Petar Velickovic	Google DeepMind / Cambridge University (Senior Research Scientist / Lecturer)	"This is a beautifully executed piece of theoretical research, which also aims for high accessibility to uninitiated readers — a trait rarely seen in papers of this kind... These researchers, he says, "are clearly excellent communicators, and I cannot wait to see what they come up with next!" <sup>2</sup>

## Works cited

1. [Literature Review] FlashAttention on a Napkin: A Diagrammatic Approach to Deep Learning IO-Awareness - Moonlight, accessed May 10, 2025, <https://www.themoonlight.io/review/flashattention-on-a-napkin-a-diagrammatic-approach-to-deep-learning-io-awareness>
2. Designing a new way to optimize complex coordinated systems | MIT News, accessed May 10, 2025, <https://news.mit.edu/2025/designing-new-way-optimize-complex-coordinated-systems-0424>
3. FlashAttention on a Napkin: A Diagrammatic Approach to Deep Learning IO-Awareness, accessed May 10, 2025, <https://openreview.net/forum?id=pF2ukh7HxA>
4. [2412.03317] FlashAttention on a Napkin: A Diagrammatic Approach to Deep Learning IO-Awareness - arXiv, accessed May 10, 2025, <https://arxiv.org/abs/2412.03317>
5. Designing a new way to optimize complex coordinated systems - IDSS, accessed May 10, 2025, <https://idss.mit.edu/news/designing-a-new-way-to-optimize-complex-coordinate-d-systems/>

6. From The Massachusetts Institute of Technology: Deep Learning – “Designing a new way to optimize complex coordinated systems” - sciencesprings, accessed May 10, 2025, <https://sciencesprings.wordpress.com/2025/04/29/from-the-massachusetts-institute-of-technology-deep-learning-designing-a-new-way-to-optimize-complex-coordinated-systems/>
7. AI Trends - Artificial Intelligence - Library and Academic Success at Georgian College, accessed May 10, 2025, <https://library.georgiancollege.ca/c.php?g=720607&p=5151225>
8. A Gentle Introduction to Category Theory - GitHub Pages, accessed May 10, 2025, <https://maartenfokkinga.github.io/utwente/mmf92b.pdf>
9. Category Theory (Stanford Encyclopedia of Philosophy), accessed May 10, 2025, <https://plato.stanford.edu/entries/category-theory/>
10. Category theory - Wikipedia, accessed May 10, 2025, [https://en.wikipedia.org/wiki/Category\\_theory](https://en.wikipedia.org/wiki/Category_theory)
11. What is Category Theory Anyway? - Math3ma, accessed May 10, 2025, <https://www.math3ma.com/blog/what-is-category-theory-anyway>
12. What even is category theory anyway? : r/math - Reddit, accessed May 10, 2025, [https://www.reddit.com/r/math/comments/ft0tdw/what\\_even\\_is\\_category\\_theory\\_anyway/](https://www.reddit.com/r/math/comments/ft0tdw/what_even_is_category_theory_anyway/)
13. Basic category theory, I | Todd and Vishal's blog, accessed May 10, 2025, <https://topologicalmusings.wordpress.com/2008/06/22/basic-category-theory-i/>
14. What is category theory? - YouTube, accessed May 10, 2025, <https://www.youtube.com/watch?v=eXBwU9ieLL0&pp=0gcJCdgAo7VqN5tD>
15. Basic Category Theory - arXiv, accessed May 10, 2025, <https://arxiv.org/pdf/1612.09375>
16. Can someone explain the basics of Category Theory to me? : r/math - Reddit, accessed May 10, 2025, [https://www.reddit.com/r/math/comments/2n0b5q/can\\_someone\\_explain\\_the\\_basics\\_of\\_category\\_theory/](https://www.reddit.com/r/math/comments/2n0b5q/can_someone_explain_the_basics_of_category_theory/)
17. A gentle introduction to category theory - YouTube, accessed May 10, 2025, <https://www.youtube.com/watch?v=yP2RjVD-cZO>
18. Deep Learning - NVIDIA Developer, accessed May 10, 2025, <https://developer.nvidia.com/deep-learning>
19. How GPUs Supercharge AI and ML for Breakthroughs - Hyperstack, accessed May 10, 2025, <https://www.hyperstack.cloud/blog/thought-leadership/how-gpus-supercharge-ai-and-ml-for-breakthroughs>
20. Large-Scale AI Model Training: Key Challenges and Innovations - AiThORITY, accessed May 10, 2025, <https://aithority.com/natural-language/large-scale-ai-model-training-key-challenges-and-innovations/>
21. Deep Learning GPU: Making the Most of GPUs for Your Project - Run:ai, accessed May 10, 2025, <https://www.run.ai/guides/gpu-deep-learning>
22. GPU Acceleration in AI: How Graphics Processing Units Drive Deep Learning -

- Gcore, accessed May 10, 2025, <https://gcore.com/blog/deep-learning-gpu>
23. The role of GPU architecture in AI and machine learning - Telnyx, accessed May 10, 2025, <https://telnyx.com/resources/gpu-architecture-ai>
  24. An Introduction to GPU Performance Optimization for Deep Learning | DigitalOcean, accessed May 10, 2025, <https://www.digitalocean.com/community/tutorials/an-introduction-to-gpu-optimization>
  25. Demystifying GPU Architectures For Deep Learning – Part 1 | - LearnOpenCV, accessed May 10, 2025, <https://learnopencv.com/demystifying-gpu-architectures-for-deep-learning/>
  26. GPU Performance Background User's Guide - NVIDIA Docs, accessed May 10, 2025, <https://docs.nvidia.com/deeplearning/performance/dl-performance-gpu-background/index.html>
  27. Which graphics card should I get for deep learning? : r/learnmachinelearning - Reddit, accessed May 10, 2025, [https://www.reddit.com/r/learnmachinelearning/comments/1364ktp/which\\_graphics\\_card\\_should\\_i\\_get\\_for\\_deep\\_learning/](https://www.reddit.com/r/learnmachinelearning/comments/1364ktp/which_graphics_card_should_i_get_for_deep_learning/)
  28. ELI5: What about GPU Architecture makes them superior for training neural networks over CPUs? : r/explainlikeimfive - Reddit, accessed May 10, 2025, [https://www.reddit.com/r/explainlikeimfive/comments/zpso6w/eli5\\_what\\_about\\_gpu\\_architecture\\_makes\\_them/](https://www.reddit.com/r/explainlikeimfive/comments/zpso6w/eli5_what_about_gpu_architecture_makes_them/)
  29. Distributed Data Parallel: Speeding Up Deep Learning - Acceldata, accessed May 10, 2025, <https://www.acceldata.io/blog/how-distributed-data-parallel-transforms-deep-learning>
  30. How does parallel processing improve the performance of deep learning models?, accessed May 10, 2025, <https://massedcompute.com/faq-answers/?question=How+does+parallel+processing+improve+the+performance+of+deep+learning+models%3F>
  31. How parallel training works in PyTorch and Deep Learning? The comprehensive guide., accessed May 10, 2025, <https://www.corpnce.com/how-parallel-training-works-in-pytorch-and-deep-learning-the-comprehensive-guide/>
  32. What is the role of parallel processing in deep learning model performance?, accessed May 10, 2025, <https://massedcompute.com/faq-answers/?question=What%20is%20the%20role%20of%20parallel%20processing%20in%20deep%20learning%20model%20performance?>
  33. Deep Learning At Scale: Parallel Model Training | Towards Data Science, accessed May 10, 2025, <https://towardsdatascience.com/deep-learning-at-scale-parallel-model-training-d7c22904b5a4/>
  34. Parallel and Distributed Deep Learning - Stanford University, accessed May 10, 2025,

- [https://web.stanford.edu/~rezab/classes/cme323/S16/projects\\_reports/hedge\\_usmani.pdf](https://web.stanford.edu/~rezab/classes/cme323/S16/projects_reports/hedge_usmani.pdf)
35. Scale Up Deep Learning in Parallel, on GPUs, and in the Cloud - MathWorks, accessed May 10, 2025, <https://la.mathworks.com/help/deeplearning/ug/scale-up-deep-learning-in-parallel-on-gpus-and-in-the-cloud.html>
  36. How useful is knowledge of parallel programming in ML? [D] : r/MachineLearning - Reddit, accessed May 10, 2025, [https://www.reddit.com/r/MachineLearning/comments/s64ozi/how\\_useful\\_is\\_knowledge\\_of\\_parallel\\_programming/](https://www.reddit.com/r/MachineLearning/comments/s64ozi/how_useful_is_knowledge_of_parallel_programming/)
  37. Is Parallel Computing useful for ML? : r/MLQuestions - Reddit, accessed May 10, 2025, [https://www.reddit.com/r/MLQuestions/comments/nfff20/is\\_parallel\\_computing\\_useful\\_for\\_ml/](https://www.reddit.com/r/MLQuestions/comments/nfff20/is_parallel_computing_useful_for_ml/)
  38. GPU accelerated deep learning: Real-time inference - KX, accessed May 10, 2025, <https://kx.com/blog/gpu-accelerated-deep-learning-real-time-inference/>
  39. How GPUs Enhance Machine Learning and AI Performance - Aethir, accessed May 10, 2025, <https://blog.aethir.com/blog-posts/how-gpus-enhance-machine-learning-and-ai-performance>
  40. The Role of GPUs in AI: Accelerating Innovation | TRG Datacenters, accessed May 10, 2025, <https://www.trgdatacenters.com/resource/gpu-for-ai/>
  41. Inference: The Next Step in GPU-Accelerated Deep Learning | NVIDIA Technical Blog, accessed May 10, 2025, <https://developer.nvidia.com/blog/inference-next-step-gpu-accelerated-deep-learning/>
  42. GPU Cloud - VMs for Deep Learning - Lambda, accessed May 10, 2025, <https://lambda.ai/service/gpu-cloud>
  43. Top 10 Best GPUs for Deep Learning in 2025 - Cherry Servers, accessed May 10, 2025, <https://www.cherryservers.com/blog/best-gpus-for-deep-learning>
  44. Best GPU for Deep Learning: Considerations for Large-Scale AI, accessed May 10, 2025, <https://www.run.ai/guides/gpu-deep-learning/best-gpu-for-deep-learning>
  45. Accelerate AI & Machine Learning Workflows | NVIDIA Run:ai, accessed May 10, 2025, <https://www.nvidia.com/en-us/software/run-ai/>
  46. GPU Accelerated Solutions for Data Science - NVIDIA, accessed May 10, 2025, <https://www.nvidia.com/en-us/deep-learning-ai/solutions/data-science/>
  47. Deep Learning – What Is It and Why Does It Matter? - NVIDIA, accessed May 10, 2025, <https://www.nvidia.com/en-us/glossary/deep-learning/>
  48. The Best GPUs for Deep Learning in 2023 — An In-depth Analysis - Tim Dettmers, accessed May 10, 2025, <https://timdettmers.com/2023/01/30/which-gpu-for-deep-learning/>
  49. 15 Best GPUs for Machine Learning for Your Next Project - ProjectPro, accessed May 10, 2025, <https://www.projectpro.io/article/gpus-for-machine-learning/677>
  50. Deep Learning Model Optimization Methods - neptune.ai, accessed May 10,



- 2025, <https://neptune.ai/blog/deep-learning-model-optimization-methods>
51. Addressing Challenges in Large-scale Distributed AI Systems | Open Research Commons, accessed May 10, 2025, <https://bcommons.berkeley.edu/addressing-challenges-large-scale-distributed-ai-systems>
  52. Optimizing Resource Allocation in Cloud for Large-Scale Deep Learning Models in Natural Language Processing | Journal of Electrical Systems, accessed May 10, 2025, <https://journal.esrgroups.org/jes/article/view/652>
  53. Large-Scale Deep Learning Optimizations | Restackio, accessed May 10, 2025, <https://www.restack.io/p/deep-learning-answer-large-scale-optimizations-cat-ai>
  54. A Systematic Survey of Resource-Efficient Large Language Models - arXiv, accessed May 10, 2025, <https://arxiv.org/pdf/2401.00625>
  55. Rethinking Large Language Models for Efficiency and Performance - Allganize, accessed May 10, 2025, <https://www.allganize.ai/en/blog/rethinking-large-language-models-for-efficiency-and-performance>
  56. Efficient Large Language Models: A Survey | OpenReview, accessed May 10, 2025, <https://openreview.net/forum?id=bsCCJHbO8A>
  57. [2401.00625] Beyond Efficiency: A Systematic Survey of Resource-Efficient Large Language Models - arXiv, accessed May 10, 2025, <https://arxiv.org/abs/2401.00625>
  58. [Literature Review] Optimization Strategies for Enhancing Resource Efficiency in Transformers & Large Language Models - Moonlight, accessed May 10, 2025, <https://www.themoonlight.io/en/review/optimization-strategies-for-enhancing-resource-efficiency-in-transformers-large-language-models>
  59. tiingweii-shii/Awesome-Resource-Efficient-LLM-Papers - GitHub, accessed May 10, 2025, <https://github.com/tiingweii-shii/Awesome-Resource-Efficient-LLM-Papers>
  60. Paper page - A Survey of Resource-efficient LLM and Multimodal Foundation Models, accessed May 10, 2025, <https://huggingface.co/papers/2401.08092>
  61. Beyond Efficiency: A Systematic Survey of Resource-Efficient Large Language Models, accessed May 10, 2025, <https://www.semanticscholar.org/paper/Beyond-Efficiency%3A-A-Systematic-Survey-of-Large-Bai-Chai/6348aeb405a496ca1729f3cc6e5eb6f12d0bc151>
  62. Towards Resource Efficient and Interpretable Bias Mitigation in Natural Language Generation | OpenReview, accessed May 10, 2025, <https://openreview.net/forum?id=PjUoztugza>
  63. UbiquitousLearning/Paper-list-resource-efficient-large-language-model - GitHub, accessed May 10, 2025, <https://github.com/UbiquitousLearning/Paper-list-resource-efficient-large-language-model>
  64. Intuitive System Enhances Developers' Ability to Create More Efficient Simulations and AI Models - BIOENGINEER.ORG, accessed May 10, 2025, <https://bioengineer.org/intuitive-system-enhances-developers-ability-to-create-more-efficient-simulations-and-ai-models/>

65. FlashAttention on a Napkin: A Diagrammatic Approach to Deep Learning IO-Awareness, accessed May 10, 2025, <https://arxiv.org/html/2412.03317v1>
66. FlashAttention on a Napkin: A Diagrammatic Approach to Deep Learning IO-Awareness, accessed May 10, 2025, [https://www.researchgate.net/publication/386454403\\_FlashAttention\\_on\\_a\\_Napkin\\_in\\_A\\_Diagrammatic\\_Approach\\_to\\_Deep\\_Learning\\_IO-Awareness](https://www.researchgate.net/publication/386454403_FlashAttention_on_a_Napkin_in_A_Diagrammatic_Approach_to_Deep_Learning_IO-Awareness)
67. Challenges in Deep Learning | GeeksforGeeks, accessed May 10, 2025, <https://www.geeksforgeeks.org/challenges-in-deep-learning/>
68. 30 Major Machine Learning Limitations, Challenges & Risks - Onix-Systems, accessed May 10, 2025, <https://onix-systems.com/blog/limitations-of-machine-learning-algorithms>
69. Deep Learning Optimization Algorithms - neptune.ai, accessed May 10, 2025, <https://neptune.ai/blog/deep-learning-optimization-algorithms>
70. Optimization in Deep Learning- Learn with examples - E2E Networks, accessed May 10, 2025, <https://www.e2enetworks.com/blog/optimization-in-deep-learning-learn-with-examples>
71. On Optimization Methods for Deep Learning, accessed May 10, 2025, [https://icml.cc/2011/papers/210\\_icmlpaper.pdf](https://icml.cc/2011/papers/210_icmlpaper.pdf)
72. What are the current limitations of deep learning algorithms, and what advancements are needed to achieve more generalized, human-level artificial intelligence? - Quora, accessed May 10, 2025, <https://www.quora.com/What-are-the-current-limitations-of-deep-learning-algorithms-and-what-advancements-are-needed-to-achieve-more-generalized-human-level-artificial-intelligence>
73. The Limitations of Deep Learning - Hacker News, accessed May 10, 2025, <https://news.ycombinator.com/item?id=14790251>
74. [D] Deep Learning has a size problem. We need to focus on state-of-the-art efficiency, not state-of-the-art accuracy. : r/MachineLearning - Reddit, accessed May 10, 2025, [https://www.reddit.com/r/MachineLearning/comments/ds1xvc/d\\_deep\\_learning\\_has\\_a\\_size\\_problem\\_we\\_need\\_to/](https://www.reddit.com/r/MachineLearning/comments/ds1xvc/d_deep_learning_has_a_size_problem_we_need_to/)
75. (PDF) Examining the Optimization Challenges in Deep Models Learning - ResearchGate, accessed May 10, 2025, [https://www.researchgate.net/publication/379652427\\_Examining\\_the\\_Optimization\\_Challenges\\_in\\_Deep\\_Models\\_Learning](https://www.researchgate.net/publication/379652427_Examining_the_Optimization_Challenges_in_Deep_Models_Learning)
76. FlashAttention on a Napkin: A Diagrammatic Approach to Deep Learning IO-Awareness - OpenReview, accessed May 10, 2025, <https://openreview.net/pdf?id=pF2ukh7HxA>
77. Conference Talk 5: Napkin Math For Fine Tuning with Johnno Whitaker - Christian Mills, accessed May 10, 2025, <https://christianjmill.com/posts/mastering-llms-course-notes/conference-talk-005/>
78. Accelerating AI: A Dive into Flash Attention and Its Impact - DZone, accessed May 10, 2025, <https://dzone.com/articles/accelerating-ai-flash-attention-impact>

79. Understanding GPU Architecture: Structure, Layers, and Performance Explained, accessed May 10, 2025, <https://www.scalecomputing.com/resources/understanding-gpu-architecture>
80. Designing Hardware-Aware Algorithms: FlashAttention - DigitalOcean, accessed May 10, 2025, <https://www.digitalocean.com/community/tutorials/flashattention>
81. Flash Attention - Hugging Face, accessed May 10, 2025, [https://huggingface.co/docs/text-generation-inference/conceptual/flash\\_attention](https://huggingface.co/docs/text-generation-inference/conceptual/flash_attention)
82. FlashAttention-3: Fast and Accurate Attention with Asynchrony and Low-precision | Tri Dao, accessed May 10, 2025, <https://tridao.me/blog/2024/flash3/>
83. Papers by Gíoele Zardini - AIModels.fyi, accessed May 10, 2025, <https://www.aimodels.fyi/author-profile/gioele-zardini-f727f4d2-cf3e-4349-81e6-fa646d6a9241>
84. FLASHATTENTION: Fast and Memory-Efficient Exact Attention with IO-Awareness - OpenReview, accessed May 10, 2025, <https://openreview.net/pdf?id=H4DqfPSibmx>
85. [2205.14135] FlashAttention: Fast and Memory-Efficient Exact Attention with IO-Awareness, accessed May 10, 2025, <https://arxiv.org/abs/2205.14135>
86. FlashAttention: Revolutionizing Transformers by Overcoming Hardware Performance Bottlenecks - AIFT, accessed May 10, 2025, <https://hkaift.com/flashattention-revolutionizing-transformers-by-overcoming-hardware-performance-bottlenecks/>
87. FlashAttention: Fast and Memory-Efficient Exact Attention with IO-Awareness - deepsense.ai, accessed May 10, 2025, <https://deepsense.ai/wp-content/uploads/2023/04/2205.14135.pdf>
88. AI optimizes complex coordinated systems in groundbreaking approach - CO/AI, accessed May 10, 2025, <https://getcoai.com/news/ai-optimizes-complex-coordinated-systems-in-groundbreaking-approach/>
89. Publications - Zardini Lab, accessed May 10, 2025, <https://zardini.mit.edu/publications/>