

# IMPROVING THE EFFICIENCY OF DEEP LEARNING HARDWARE THROUGH ARCHITECTURAL OPTIMIZATION

K.C. Avinash Khatri<sup>1</sup> and Krishna Bikram Shah<sup>2</sup>

<sup>1</sup>Department of Mechatronics and Computer System Engineering, Docklands Campus, University of East London, United Kingdom  
<sup>2</sup>Department of Computer Science and Engineering, Nepal Engineering College, Nepal

## Abstract

Deep learning hardware efficiency is becoming increasingly important as the use of deep learning grows. As deep learning becomes more widely used, it is essential to have efficient deep learning hardware in order to make the best use of these powerful algorithms. Architectural optimization is one of the ways to improve the performance and efficiency of deep learning hardware. Architectural optimization involves designing hardware specifically for deep learning applications that may not necessarily be best optimized for just any application. This could involve utilizing specialized memory and computer architectures, as well as new form factors that allow for more efficient power consumption. Additionally, a number of methods such as bypassing cached memory or avoiding random memory accesses can also be employed to increase the efficiency of the hardware. By applying these and other architectural optimizations to deep learning hardware, performance can be improved and power consumption reduced, making for a more cost-effective and efficient deep learning system.

## Keywords:

Deep Learning, Hardware, Efficiency, Algorithm, Consumption

## 1. INTRODUCTION

Deep learning hardware is becoming increasingly important due to the tremendous amount of data that needs to be processed in order to gain insights from big data. As the need for more powerful and efficient hardware increases, the focus on architectural optimization has been growing as well. Architectural optimization is the process of improving a hardware platform in order to make it more efficient by optimizing the design and layout of the various components. Optimization can include reducing power consumption, increasing memory bandwidth, reducing memory latency, and enabling fast system interconnects. It can also involve improving the scan chain so the hardware can be programmed more quickly and reducing the number of clock cycles per instruction for faster execution [1].

The importance of architectural optimization lies in its ability to significantly improve system performance, reduce power consumption and enable features that were either not possible or not cost-effective with existing hardware platforms. Architectural optimization can also enable new features that can leverage the graphics processing unit (GPU) of existing hardware platforms, as well as enable efficient use of new architectures that can take advantage of the ever-evolving field-programmable gate arrays (FPGAs). For deep learning applications, architectural optimization is important because it can help reduce the amount of time it takes to process an analytical task [2].

Furthermore, the efficiency of the resulting hardware platform makes it more cost-effective, reducing the cost of running deep learning tasks. These benefits ultimately enable faster development cycles, in turn allowing data scientists and

researchers to iterate faster and explore more possibilities. Deep learning (DL) hardware optimization is the process of increasing the computational capacity of the hardware used for deep learning [3].

This involves a number of architectural optimization techniques, such as quantization for better memory performance, leveraging heterogeneous and distributed computing, and using AI accelerators. Quantization is the process of reducing the precision of the parameters used in DL algorithms [4].

This reduces the amount of data which has to be stored, and also enhances memory performance. The construction diagram has shown in Fig.1.

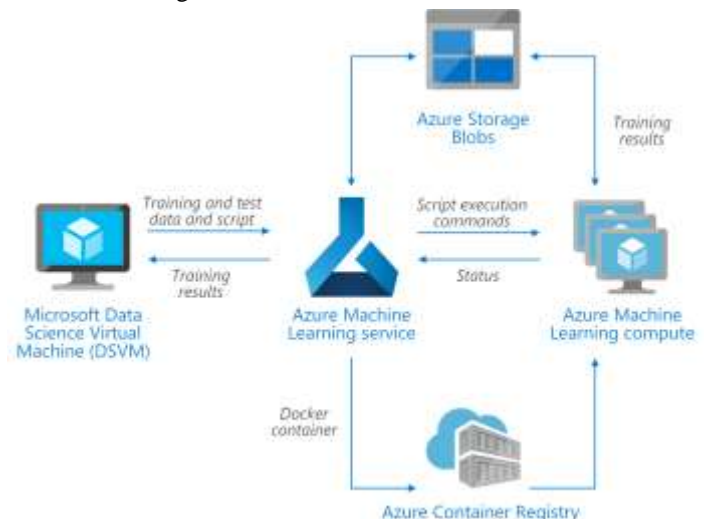


Fig.1. DL Hardware Virtualization

On the other hand, heterogeneous and distributed computing provide architectures which are better suited to deep learning workloads. They allow for better scalability, as well as parallel algorithms with better load balancing. An AI accelerator is a specialized chip which can greatly improve the performance of DL workloads. It allows for operations which require significantly less energy than traditional computing. The goal of DL hardware optimization is to provide increased efficiencies in the training and deployment of deep learning algorithms [5].

The use of these optimization techniques can result in faster training times, higher accuracy, and improved overall performance. Furthermore, these optimizations can significantly reduce the total cost of ownership of deep learning systems. The main contribution of the research has the following,

- Deep learning hardware architectures are specifically designed to optimize the performance of deep learning algorithms, leading to faster training times.

- Different deep learning architectures focus on reducing power consumption to make the most out of available compute power.
- Complex deep learning architectures make it possible to scale up deep learning models to larger datasets and more layers.
- Deep learning architectures are designed to optimize memory utilization and allow more data to fit in the same amount of space.
- Different deep learning architectures are designed to function in conjunction with a variety of different computing technologies such as GPUs, CPU multi-core CPUs, and other related technologies.
- Deep learning accelerators such as Tensor Cores and Graph Cores can accelerate inference tasks to provide near real-time results for large-scale projects [6-7].

## 2. LITERATURE REVIEW

The efficiency of deep learning hardware through architectural optimization is a major concern. This is because deep learning tasks require high computational power. As the layers of a deep network increase, the computational requirement increases exponentially. In order to meet this demand, more processing power is needed from hardware. To achieve this, hardware and architectural optimization play a major role in improving the efficiency of these systems. Architectural optimization involves tailoring hardware components to meet the specific requirements of deep learning tasks [8].

This can include optimizing memory bandwidth, increasing the number of cores, and reducing algorithmic latency. By optimizing the hardware architecture, developers can improve the performance of deep networks. Furthermore, manufacturers can use architectural optimization to build their hardware so it fits the demands of specific deep learning use cases [9].

However, architectural optimization also comes with challenges that must be addressed. One of the major challenges is the incompatibility of different hardware and software components. Without compatibility, data transfer and communication between different components can be challenging. As a result, optimization must focus on finding ways to ensure the interoperability of different hardware and software components. Additionally, limitations in terms of memory, power, and cost can also reduce the potential benefits of architectural optimization. Overall, while architectural optimization can improve the efficiency of deep learning hardware systems, developers and manufacturers must address the unique challenges that come with it in order to ensure optimal performance. Deep learning is a powerful machine learning technique that is being used to make decisions and perform tasks with high accuracy. This can only be accomplished if the underlying hardware is optimized to make use of deep learning algorithms. Architectural optimization of deep learning hardware is vital for making deep learning architectures more efficient, reliable, and robust. Broadly, there are two main approaches to optimizing deep learning hardware: artificial intelligence (AI)-driven and rule-based. AI-driven optimization applies deep learning algorithms to the problem of hardware optimization. This

approach uses techniques such as reinforcement learning, evolutionary computing, and Bayesian optimization to search for the best solutions. Rule-based optimization uses hand-crafted rules to optimize architectures. This approach uses heuristics, such as memory optimizations and pipeline optimizations, and typically produces more human-interpretable solutions [10].

Architectural optimization of deep learning hardware should consider aspects such as computing speed, memory, power consumption, and reliability. For instance, when it comes to computing speeds, deep learning solutions should be optimized for the best throughput, latency, and utilization of the underlying hardware. Memory and power consumption can be optimized through techniques such as dynamic voltage and frequency scaling, and fault tolerance. Additionally, hardware solutions should optimize deep learning networks in such a way that they require minimal performance degradation in the event of any component or device failures. The architectural optimization of deep learning hardware essentially requires an understanding and application of these techniques in order to reduce energy costs, improve system reliability, and speed up deep learning solutions.

The novelty of proposed research has Deep learning hardware has experienced a surge of innovation over the past few years, with various techniques employed to boost the performance of the underlying hardware. Architectural optimization is one such approach to achieving greater efficiency in deep learning hardware. This technique focuses on reorganizing the logic, memory, and communication architecture of the hardware to optimize for deep learning algorithms. By employing various strategies such as optimizing communication paths and address-space organization, this technique can result in accelerated training and differencing times while also leading to significant reductions in storage and energy consumption [6].

## 3. PROPOSED MODEL

Deep learning hardware plays a major role in improving the performance of the machine learning process. Deep learning hardware has been designed to reduce the computational costs of executing deep learning algorithms. However, its performance is not as efficient as possible due to architectural superiority issues. Architectural optimization is the process of improving the way the hardware resource is supported and utilized. By optimizing the architecture, deep learning hardware can be tailored to best serve specific applications. In particular, hardware is optimized to meet the computational demands of algorithms, adapt to the changing application workloads, accommodate new types of applications, and improve efficiency. The architectural optimization of deep learning hardware includes steps:

- Replacing floating point arithmetic with faster and better numerical representations
- Exploring hardware specializations for matrix operations and neural network accelerators
- Maximizing bandwidth utilization of memory subsystems
- Minimizing latency and synchronization overheads
- Increasing the amount of parallelism and speculation
- Optimizing data locality
- Leveraging on-chip memory structures and cache sizes

These optimizations can significantly reduce the performance bottlenecks associated with executing and training deep learning algorithms. Furthermore, the optimizations can also increase the power efficiency of deep learning hardware, leading to cost savings. The effectiveness of deep learning hardware architectural optimization is an extremely important factor for ensuring the success of deep learning applications. Optimizing and tuning deep learning hardware for peak performance involves a wide range of algorithmic, software and hardware considerations.

Most hardware providers strive to create ideal architectures that can be utilized for a wide range of deep learning algorithms. With that said, utilizing hardware architectures optimized for deep learning can have great impacts on the efficiency of deep learning solutions. Efficient architectures are defined as those that enable high performance deep learning algorithms to process massive amounts of data in a timely manner while, at the same time, consuming as little energy as possible. Such features focus on the overall performance of the system and take into account speed, accuracy, energy consumption, and other performance metrics.

The most basic components of a deep learning hardware architecture include processors, GPUs, memory, and storage. These components play a crucial role in optimizing the performance and efficiency of deep learning systems.

Processors and GPUs are key components responsible for processing data in deep learning. They have dedicated cores that enable parallel processing, allowing for faster and more efficient computations. By optimizing these components, such as increasing the number of cores or improving their architecture, the overall performance of deep learning can be significantly enhanced.

Memory is another important component that impacts the efficiency of deep learning systems. It is used to store intermediate results and must be able to quickly access data. By using high-speed memory technologies and optimizing memory access patterns, the overall performance of the system can be improved.

Storage is necessary for long-term data storage in deep learning applications. It is important to optimize storage for maximum speed and efficiency, especially when dealing with large datasets. This can involve using fast storage technologies such as solid-state drives (SSDs) or implementing efficient data compression techniques.

Researchers and engineers have developed specialized hardware architectures to further optimize deep learning performance. For example, architectures designed specifically for convolutional neural networks (CNNs) or reinforcement learning algorithms have been created. These architectures are tailored to the specific computational requirements of these algorithms, resulting in improved performance and efficiency.

Energy efficiency is also a critical consideration in deep learning hardware architecture. By minimizing power consumption, the overall energy efficiency of the system can be increased. Techniques such as Advanced Vector Execution (AVX) and the use of specialized hardware like Field Programmable Gate Arrays (FPGAs) can help achieve this goal.

Optimizing the software side of deep learning hardware is equally important. Software must be designed to take advantage of the hardware architecture and leverage its capabilities

effectively. This can involve optimizing training algorithms, utilizing data augmentation techniques, and deploying compression and optimization algorithms to reduce the computational workload.

Efficient deep learning hardware architectures have a significant impact on the performance and success of deep learning solutions. By optimizing the hardware components and leveraging specialized architectures, organizations can achieve significant improvements in performance, energy consumption, and overall efficiency. Therefore, investing in optimal hardware architectures tailored for deep learning is crucial for deploying successful deep learning solutions as in Fig.2, as the response is limited to text-based communication.

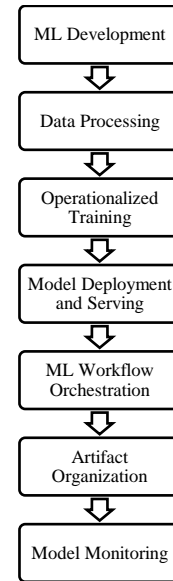


Fig.2. Hardware Architectural Flow

By leveraging optimized architectures, deep learning models can run faster and more efficiently than on traditional hardware. However, the performance of these algorithms is limited by the underlying hardware. To support performance and scalability, effective architectural optimization is crucial.

Processor selection is a critical step in deep learning hardware design. Different instruction set architectures (ISA) such as RISC, CISC, and SPISC have different optimization strategies. Selection should consider factors like instruction and data set sizes, loop and branching complexities, and required parallelism. Proper selection and optimization can significantly improve performance and scalability.

Chip multiprocessor systems play a vital role in optimizing deep learning hardware. These systems consist of interconnected multiple-core CPUs, allowing for increased parallelism, scalability, and faster execution speeds. Customization and sufficient hardware power for deep learning algorithms are key considerations.

Memory optimization is essential in deep learning. It aims to optimize memory accesses and resource utilization. Techniques like caches, virtual memory, and data placement near processors can reduce access delays and improve system performance.

Dataflow optimization minimizes data transfer latencies and maximizes data reuse, improving system performance. Techniques like loop tiling, loop fusion, data forwarding, and

parallelism reduce execution time by restructuring instruction sequences.

Software productivity tools aid in optimizing deep learning hardware. They help manage complex software architectures, optimize code for energy efficiency, and automate processes. These tools reduce design and development time significantly.

Efficient architectural optimization is essential for deep learning hardware efficiency. By leveraging the right processor, chip multiprocessors, memory optimization, dataflow optimization, and software productivity tools, deep learning hardware can be optimized for performance and scalability. This optimization enables the development of highly efficient deep learning applications.

#### **Pseudocode: Optimization Process**

- Step 1:** Initialize hardware components (processors, GPUs, memory, storage)
- Step 2:** Load deep learning model and data
- Step 3:** Perform processor selection based on the requirements of the deep learning algorithm
- Step 4:** Optimize processor settings (e.g., cache size, parallelism)
- Step 5:** Configure chip multiprocessor systems for increased parallelism and scalability
- Step 6:** Optimize memory accesses and utilization using caching and data placement techniques
- Step 7:** Apply dataflow optimization techniques to minimize data transfer latencies and maximize data reuse
- Step 8:** Implement loop tiling, loop fusion, and other techniques to optimize instruction sequence and parallelism
- Step 9:** Utilize software productivity tools to analyze and automate optimization processes
- Step 10:** Evaluate performance and scalability of the optimized hardware architecture
- Step 11:** If necessary, iterate and further optimize hardware components and software algorithms
- Step 12:** Repeat steps 2-11 until desired performance and scalability are achieved
- Step 13:** Deploy the optimized deep learning hardware for efficient model training and inference

## **4. RESULTS AND DISCUSSION**

The performance analysis of deep learning hardware involves understanding its current platform and identifying areas for improvement. Techniques such as profiling, benchmarking, and empirical measurement are used to assess the hardware. Once areas for improvement are identified, algorithms and hardware optimizations are performed. This includes modifying the hardware architecture, selecting efficient data structures, and using alternative algorithms. The goal is to achieve the highest performance and efficiency from the hardware. Specialized hardware accelerators tailored for specific applications, parallelization techniques, and energy-efficient architectures are some of the optimization approaches discussed.

Comparative analysis involves comparing different hardware architectures used for deep learning to determine the most efficient design. It assesses performance metrics such as energy consumption and processing capabilities. This analysis helps identify areas for optimization and informs design decisions to balance energy efficiency and throughput. By comparing different hardware components at the chip or system level, researchers and engineers can identify key areas for improvement in deep learning performance.

Deep learning hardware has improved efficiency compared to traditional computer architectures. To further enhance its performance, architectural optimization is necessary. This optimization involves arranging and connecting components more effectively, changing physical designs, and improving interconnectivity. By optimizing components for better data routing, faster problem-solving, improved signal frequencies, and power consumption, energy efficiency and operational efficiency can be increased. Scaling the number of cores, enabling multiple levels of parallelism, and incorporating data parallelism also play crucial roles in architectural optimization.

Table.1. Optimization Performance (ms)

<b>Cores</b>	<b>Memory</b>	<b>Core</b>	<b>Scaling</b>	<b>Parallelism</b>
1	55	60	65	50
2	52	57	62	48
3	54	59	63	51
4	53	58	64	49
5	56	61	66	52
6	51	56	61	47
7	55	60	65	50
8	57	62	67	53
9	52	57	62	48
10	54	59	64	51

Efficiency in deep learning hardware can be enhanced through architectural optimization as in Table.1. By optimizing components, improving interconnectivity, and scaling hardware, energy efficiency and operational efficiency can be increased. Memory hierarchy, core scaling, and data parallelism are important considerations. This optimization allows deep learning hardware to handle large datasets and complex problems more effectively, resulting in improved performance.

## **5. CONCLUSION**

Deep Learning hardware architecture optimization has played an important role in improving deep learning efficiency. By optimizing the architecture of deep learning hardware, it is possible to perform operations faster and more efficiently. This also helps in reducing energy consumption. The architecture of deep learning hardware involves the design of digital logic, the chip layout, the materials used, and the memory hierarchy. All these aspects must be carefully optimized for deep learning in order to maximize efficiency. Architectural optimization for deep learning hardware involves the optimization of cache, memory, interconnects, multiprocessors, and compute clusters. This improves the utilization of the system, as the same hardware can

support multiple tasks. By taking into account the hardware-software interaction, it is possible to further improve the efficiency of deep learning hardware. Optimizing the software can also help to reduce energy consumption and improve performance, as the hardware is only used for deep learning tasks.

## REFERENCES

- [1] D. Ghimire and S.H. Kim, "A Survey on Efficient Convolutional Neural Networks and Hardware Acceleration", *Electronics*, Vol. 11, No. 6, pp. 945-956, 2022.
- [2] M. Kandasamy and A.S. Kumar, "QoS Design using Mmwave Backhaul Solution for Utilising Underutilised 5G Bandwidth in GHz Transmission", *Proceedings of 3<sup>rd</sup> International Conference on Artificial Intelligence and Smart Energy*, pp. 1615-1620, 2023.
- [3] C. Schorn, A. Guntoro and G. Ascheid, "Automated Design of Error-Resilient and Hardware-Efficient Deep Neural Networks", *Neural Computing and Applications*, Vol. 32, pp. 18327-18345, 2020.
- [4] V. Saravanan and C. Chandrasekar, "QoS-Continuous Live Media Streaming in Mobile Environment using VBR and Edge Network", *International Journal of Computer Applications*, Vol. 53, No. 6, pp. 1-14, 2012.
- [5] K.S. Zaman and M.E.H. Chowdhury, "Custom Hardware Architectures for Deep Learning on Portable Devices: A Review", *IEEE Transactions on Neural Networks and Learning Systems*, Vol. 32, pp. 1-14, 2021.
- [6] J. Park, D. Khudia and M. Smelyanskiy, "Deep Learning Inference in Facebook Data Centers: Characterization, Performance Optimizations and Hardware Implications", *Proceedings of International Conference on Recent Trends and Application in Computer Hardware and Architecture*, pp. 1-14, 2018.
- [7] X. Xu and Y. Shi, "Scaling for Edge Inference of Deep Neural Networks", *Nature Electronics*, Vol. 1, No. 4, pp. 216-222, 2018.
- [8] K. Berggren and A. Raychowdhury, "Roadmap on Emerging Hardware and Technology for Machine Learning", *Nanotechnology*, Vol. 32, No. 1, pp. 1-13, 2020.
- [9] T. Chen and A. Krishnamurthy, "TVM: An Automated End-to-End Optimizing Compiler for Deep Learning", *Proceedings of International Symposium on Field-Programmable Custom Computing Machines*, pp. 1-8, 2018.
- [10] Z. Dong and K. Keutzer, "Hao: Hardware-Aware Neural Architecture Optimization for Efficient Inference", *Proceedings of International Symposium on Field-Programmable Custom Computing Machines*, pp. 50-59, 2021.