

# Advancing AI: Enhancing Large Language Model Performance through GPU Optimization Techniques

Sriram Sagi

NetApp

**Abstract:** *This study delves into optimizing GPU utilization for supporting Large Language Models LLMs within Generative AI frameworks. Focusing on dynamic resource allocation, kernel optimization, and memory management, our investigation reveals significant improvements in LLM efficiency and performance. By integrating NVIDIA's advanced AI technologies, we propose a scalable, cost-effective approach for deploying AI applications at the enterprise level. The findings underscore the pivotal role of GPU optimization in enhancing AI accessibility and fostering innovation across diverse sectors.*

**Keywords:** Large Language Models (LLMs), GPU Optimization, Generative AI, Artificial Intelligence Deployment

## 1. Introduction

Large language models have become really popular and are being utilized in a bunch of applications. Thanks to their amazing ability to understand language, create text, hold onto context, and figure out problems. The ability of these large language models, or LLMs, to look into complicated language patterns and give the right solutions within a given situation. Models such as OpenAI's GPT series have truly walked a whole mile in terms of technological advancements and the applications they can be used for. These advancements really do mark a big thing in the advancement of artificial intelligence. These models have been trained on a large amount of text data. They've shown impressive abilities in knowing human language. Not just that, but also in creating content and having the ability to chat with it. Their versatility allows them to be used in a bunch of areas like creative writing and content development. There's more you know like sentiment analysis, legal document evaluation, and coding activities. Every day we see people more dependent on these LLMs in all sorts of sectors. It's because of their power to analyze data and produce language which kind of sounds exactly like human speech. Because of these attributes, they are considered super-effective tools for making things more efficient, inspiring innovation and helping us make better decisions. This growth really tells us how we are moving towards automation and smart systems in this digital age, like you can imagine a future where AI has a big role, it can take on complex problems and make innovative solutions. These language models are everyone's go-to tool for language processing and that is surely going to boost our usage of digitization.

The purpose of this article is to explore and present effective strategies for optimizing GPU performance to support the deployment and efficiency of Large Language Models LLMs in generative AI applications.

## 2. Importance of optimizing GPU computations and resource allocation for LLMs

Efficiently maximizing the use of GPU processing and managing resources is essential to harnessing the capabilities of Language Models (LLMs) used in tasks like natural language processing, machine learning and artificial intelligence. GPUs are key in accelerating the computations required for training and deploying these models as they can handle multiple tasks simultaneously. By optimizing GPU utilization, we ensure that these resource intensive operations should not only run faster but, in a cost-effective manner facilitating the widespread and affordable adoption of LLMs. Enhancing GPU resources has an impact on critical areas. It reduces the time needed for model training and inference cuts costs by enhancing power efficiency and enables the use of sophisticated models capable of processing larger datasets with increased accuracy. Moreover, during a period where AI applications are highly sought after efficient resource management plays a role in mitigating the impact associated with these activities. This involves minimizing the carbon footprint resulting from power consumption and cooling requirements.

Furthermore, incorporating optimization methods such as precision training, kernel optimization and memory management strategies can significantly boost the performance of Large Language Models (LLM) models. These techniques expedite the training process. Also improve the efficiency of inference tasks enabling real-time applications and interactions that were previously unattainable. As LLM models evolve and uncover applications, the importance of optimizing GPU computations and resource allocation will grow to support the development of advanced, efficient and ecofriendly artificial intelligence (AI) systems. The continuous advancements in LLMs necessitate an evolution in the strategies and technologies used to enhance GPU computations and resource allocation. Successfully implementing AI is not a challenge but also a critical strategic objective for businesses looking to leverage AI for competitive advantage. With

LLMs becoming more complex their data requirements and computational demands are escalating. This underscores the significance of employing approaches to optimize GPU performance, such as dynamic resource allocation. This approach involves dynamically assigning GPU resources based on workload demands in time to ensure efficiency during various stages of model training and deployment.

Additionally, the rise of hardware and software solutions tailored for AI tasks opens possibilities for improving effectiveness. These solutions are specifically crafted to boost data processing speed and reduce latency leading to iterations in developing and implementing models. Dedicated hardware accelerators for AI computations offer advantages over general purpose GPUs in terms of energy efficiency and computational performance.

Efficiently managing GPU resources is crucial to making robust LLM capabilities more accessible to an audience. By improving the availability of training and deploying models we can foster an inclusive AI innovation environment that caters to smaller teams and businesses. The democratization of AI not only accelerates research and development progress but also ensures a wider array of applications and perspectives.

With the increasing research and development of AI, also, there is needs to maximize the GPU power and handling resource appropriately. There are many LLMs in play today. These advancements enable for the esteemed models being used with effectiveness in scenarios of the real world, and you know, pushing the invisible boundaries of AI abilities. It is essential to, focus trying, to refining strategies for optimization in moving forward to the potential for realizing the potential, of the LLMs. Opportunities for integration among, several sectors Such as healthcare, education, finance, entertainment will be arising shaping technology, and societal landscapes. Exploration for methods of efficiency, GPU computations and allocation of resources, when Large Language, Models. For GenAI are hosted, is straight the focus of this paper. Overcoming challenges, such as optimizing GPU resources for hosting LLMs, in GenAI cases, super critical outlined in this study.

### 3. Challenges in optimizing GPU computations and resource allocation for LLMs

Efficiently managing GPU computations and resource allocation for Large Language Models (LLMs) is crucial to ensure the scalability and effectiveness of these AI systems. However, tackling this task presents challenges due to the practical complexities involved in deploying LLMs on a large scale.

One major hurdle is the high computational complexity and resource requirements of LLMs. Training these models involves processing datasets using neural networks demanding substantial computing power and memory bandwidth. With models there arises a need for more powerful GPU memory and processing capabilities that often outstrip current hardware resources. This mismatch can lead to data transmission bottlenecks and slower computational

speeds impacting the efficiency and cost effectiveness of LLM implementations.

One of the challenges lies in distributing tasks across multiple GPUs to achieve optimal parallelization. While GPUs are designed for processing achieving peak performance, Large Language Models (LLMs) require careful tuning and coordination of these operations. Maintaining performance involves ensuring the distribution of data and computational tasks while minimizing communication overhead between GPUs, which can be complex and require intricate programming and configuration.

Managing energy consumption and thermal regulation presents challenges as well. The high computational demands of training and deploying LLMs result in increased energy usage leading to costs and environmental concerns. Additionally, the intensive workload of GPUs generates heat that necessitates advanced cooling solutions to maintain performance and prevent hardware damage.

The rapid advancements in LLM techniques and GPU technology pose optimization challenges as developers must continually adapt their approaches to leverage the hardware capabilities and algorithmic enhancements. Proficiency in LLMs requires an understanding of both the foundations and practical aspects of GPU design and programming to stay abreast of evolving technologies.

### 4. Strategies for optimizing GPU computations and resource allocation

Optimizing the performance of Large Language Models (LLMs) by managing GPU computations and allocating resources is essential for improving results, cutting down on costs and making the most of hardware resources. It requires a strategy that considers both how hardware is used and the efficiency of algorithms. This involves tactics like parallelization and distribution adaptable resource allocation training with mixed precision optimizing kernels, managing memory effectively fine - tuning software and frameworks ensuring algorithmic efficiency promoting energy efficient computing practices and keeping an eye on hardware usage through monitoring and profiling.

Data parallelism entails spreading data across GPUs to enable each GPU to process a portion of the data. On the other hand, model parallelism divides the model among GPUs when it's too large to fit into a single GPUs memory. Dynamic resource allocation adjusts GPU resource distribution in time to meet demands promptly while maximizing utilization and minimizing periods. Mixed precision training involves using both half precision floats in calculations to reduce memory usage and accelerate operations without compromising accuracy significantly. Kernel optimization focuses on enhancing GPU kernels for functions to improve execution speed and memory access times. Memory management techniques such as checkpointing and in place operations help decrease memory usage.

Software optimization takes advantage of libraries and frameworks optimized for computations, on GPUs like CUDA Deep Neural Network (cuDNN) tailored for deep

learning tasks. Efficiency in algorithms means tweaking network design to cut down on demands while maintaining model accuracy. Strategies for energy computing include adapting GPU voltage and frequency according to workload thereby lowering power usage and heat output during computational workloads. Keeping track of hardware utilization and performance profiles assists in spotting areas of congestion and inefficiencies, in resource allocation.

## 5. Literature Review

Multiple studies have delved into the way to distribute resources, in GPUs. Jooya (2012). Punyala (2018) both propose methods to optimize resource allocation with Jooya focusing on exploring design options and Punyala concentrating on minimizing interference and maximizing throughput. The research introduces an approach to reduce interference in shared resources for execution on GPUs leading to a significant boost in system throughput compared to default and profile-based optimization techniques. The ILP SMRA method showcases improvements in throughput for specific workloads.

Yunjoo (2019) and Gelado (2019) delve into resource utilization strategies and memory allocation techniques to further enhance efficiency. The article categorizes GPU workloads. Assesses the bottleneck resources to aid in efficient workload distribution. When the grained bottleneck resources vary, among identical computing bound workloads it becomes feasible to allocate these workloads together. The research introduces programming methods and synchronization primitives to achieve high levels of concurrency on GPUs. This leads to the development of a memory allocator with improved allocation speeds that outperform those of CUDA 9. A throughput focused memory allocator achieves high memory allocation rates while reducing memory fragmentation levels, on GPUs effectively. Cebrian (2012) and Li (2016) both emphasize the importance of energy efficiency with Cebrian noting the underutilization of resources and Li proposing a strategy, for managing kernels. Improving resource utilization can boost the energy efficiency of computations carried out on GPUs. There is room for resource utilization in GPU based processing, which can be optimized to enhance energy efficiency. Their research delves into approaches and suggestions aimed at enhancing energy efficiency in GPU designs. The study introduces a framework designed to improve performance and energy conservation when running kernels on GPUs. This framework enhances performance by a factor of 1.42X. Boosts energy efficiency by an average factor of 1.33X compared to the default concurrent kernel execution setup.

Moreover Ibrahim (2016) expands on this discussion by exploring the context and introducing a resource allocation technique tailored specifically for GPUs in a cloud setting. The study introduces a resource allocation algorithm for GPUs in cloud environments focusing on providing GPUs within a Software as a Service (SaaS) environment through load balancing mechanisms. The paper showcases the systems scalability and its enhancements, in Service Level Agreement (SLA). These studies show how important it is to allocate resources to improve the performance and energy efficiency of GPUs. The research also emphasizes the need to

account for the changing nature of cloud environments when assigning resources, to GPUs. By using real time monitoring and adaptive algorithms the suggested method can optimize resource allocation effectively according to varying workload requirements.

## 6. Results

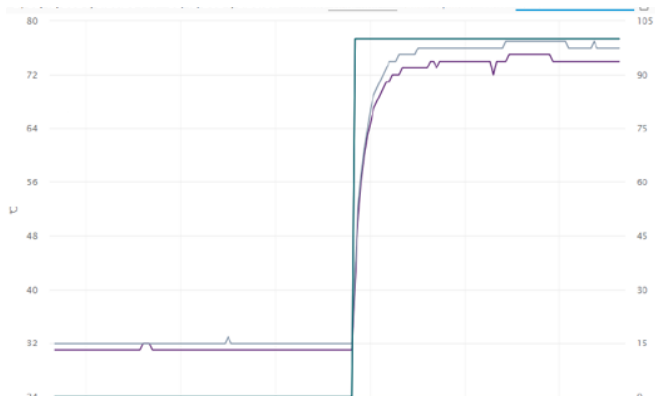
During this validation process we assessed Generative AI Inference Models. One method for executing Generative AI Inference involves setting up an inference server with an AI model. The server has a listener that remains inactive until it receives inference requests, to which it promptly responds. Another approach is to transfer the model to the GPU (s), perform the inference request (s), provide the response (s) and then remove the model. In this validation we used NetApp Astra Trident provided storage to store the AI models thereby eliminating the need to download them every time a pod was created or recreated. The validation process included obtaining AI software and models from sources, like NVAIE, Hugging Face, Github and other platforms. NVAIE provides a range of inference servers, AI frameworks and AI models tailored for NVIDIA GPUs.

### GPU Monitoring with nvidia - smi:

The NVIDIA System Management Interface, commonly referred to as Nvidia smi is a command line tool that comes bundled with the NVIDIA GPU driver package. It serves as a tool for monitoring and managing NVIDIA GPU devices on both Windows and Linux systems. By using nvidia smi, users and system administrators can access real time information about the status of the GPU and its driver as they manage various important GPU settings. Within the NVAIE environment the nvidia smi utility allows for monitoring directly on the interface of the tool. This can be achieved by employing either the loop command line parameter or utilizing the Linux watch command. Additionally, this application provides a set of command line parameters to gather information about the GPU performance. It also has the ability to generate CSV files containing data at specified intervals, which can then be utilized for creating representations of GPU data through graphs.

### GPU Monitoring with VMware vCenter GPU Statistics

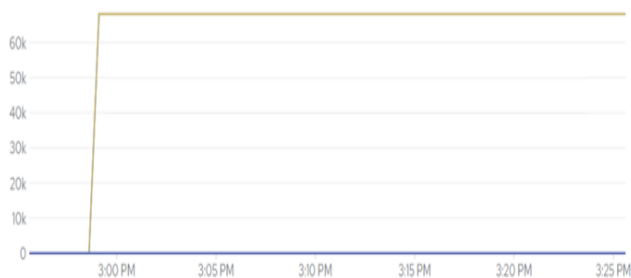
The NVAIE vibs installed in VMware ESXi enable VMware vCenter to gather and display real - time data on the physical GPU. VMware vCenter collects and displays data on GPU memory consumption, memory consumed, temperature, and utilization. Users can choose up to two of these metrics to be charted and shown. The diagram below illustrates a specific instance of executing GPU Burn.



Picture 1: GPU Burn with VMware vCenter GPU Statistics

### GPU Monitoring with NVIDIA DCGM Exporter Dashboard

The NVIDIA GPU Operator provides access to vGPU telemetry data for Prometheus through the utilization of the NVIDIA DCGM Exporter. The NVIDIA DCGM Exporter Grafana Dashboard can be installed on the OpenShift Container Platform (OCP) and accessed using the OCP console.



Picture 2: NVIDIA DCGM dashboard while running GPU Burn

A Python benchmark script included in the NeMo Inference container was used for each of the models tested, utilizing one or two vGPUs. The results are displayed in the table below. As the models get bigger, the advantage of the second GPU is displayed, along with the average latency and throughput for each model.

Table 1: GPU Benchmark Results

Model	Avg Latency (ms)		Avg Throughput (Sentence/s)	
	1 GPU	2 GPUs	1 GPU	2 GPUs
Llama-2-7B	7026.116	6341.28	1.139	1.262
Llama-2-13B	11851.359	9662.157	0.675	0.828
NeMo GPT 2B	2735.752	2902.535	2.924	2.756
Nemotron 3 8B	7071.629	6269.211	1.131	1.276

## 7. Conclusion

The deployment and hosting of Large Language Models (LLMs), for Generative AI pose challenges, due to the significant computational resources and energy consumption they require. Our research demonstrates the importance of optimizing GPU resources to address these challenges. By utilizing techniques like dynamic resource allocation, training, and memory management we can greatly enhance

the performance and efficiency of LLMs. Collaborating with NVIDIA's technology, such as AI Enterprise software further strengthens our ability to create an efficient environment for AI applications. This study emphasizes how optimized GPU computations can revolutionize LLM hosting promoting accessible and effective AI solutions. With the increasing demand for AI technology the methodologies and insights detailed in this paper provide guidance for advancements in AI infrastructure laying the groundwork for an era of intelligent computing.

### Acknowledgements

This article was supported by the technical team from Cisco and NetApp, we sincerely thank them for their contribution and valuable input.

### References

- [1] Jooya, A. Z. et al. "Efficient Design Space Exploration of GPGPU Architectures." Euro - Par Workshops (2012).
- [2] Punyala, Srinivasa Reddy et al. "Throughput optimization and resource allocation on GPUs under multi - application execution." 2018 Design, Automation & Test in Europe Conference & Exhibition (DATE) (2018): 73 - 78.
- [3] yunjoo, park et al. "Analyzing Fine - Grained Resource Utilization for Efficient GPU Workload Allocation." The Journal of the Institute of Webcasting, Internet and Telecommunication 19 (2019): 111 - 116.
- [4] Gelado, Isaac and Michael Garland. "Throughput - oriented GPU memory allocation." Proceedings of the 24th Symposium on Principles and Practice of Parallel Programming (2019): n. pag.
- [5] Cebrian, Juan M. et al. "Energy Efficiency Analysis of GPUs." 2012 IEEE 26th International Parallel and Distributed Processing Symposium Workshops & PhD Forum (2012): 1014 - 1022.
- [6] Li, Xiuhong and Yun Liang. "Efficient kernel management on GPUs." 2016 Design, Automation & Test in Europe Conference & Exhibition (DATE) (2016): 85 - 90.
- [7] Ibrahim, Ahmed H. et al. "Resource allocation algorithm for GPUs in a private cloud." Int. J. Cloud Comput.5 (2016): 45 - 56.
- [8] Li, Teng et al. "Efficient Resource Sharing Through GPU Virtualization on Accelerated High Performance Computing Systems." ArXiv abs/1511.07658 (2015): n. pag.