

SUPERPIPELINE: A UNIVERSAL APPROACH FOR REDUCING GPU MEMORY USAGE IN LARGE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

The rapid growth in size and complexity of machine learning models, particularly in natural language processing and computer vision, has led to significant challenges in model execution on hardware with limited resources. This paper introduces Superpipeline, a novel framework designed to optimize the execution of large-scale AI models on constrained hardware for both training and inference phases. Our approach focuses on dynamically managing model execution by partitioning models into individual layers and efficiently transferring these partitions between GPU and CPU memory. Superpipeline achieves substantial reductions in GPU memory consumption—up to 60% in our experiments—while maintaining model accuracy and acceptable processing speeds. This enables the execution of models that would otherwise exceed available GPU memory capacity. Unlike existing solutions that primarily target inference or specific model types, Superpipeline demonstrates broad applicability across large language models (LLMs), vision-language models (VLMs), and vision-based models. We evaluate Superpipeline’s effectiveness through comprehensive experiments on diverse models and hardware configurations. Our method is characterized by two key parameters that allow fine-tuning of the trade-off between GPU memory usage and processing speed. Importantly, Superpipeline does not require model retraining or parameter modification, ensuring full preservation of the original model’s output fidelity. The simplicity and flexibility of Superpipeline make it a valuable tool for researchers and practitioners working with state-of-the-art AI models under hardware constraints. It enables the use of larger models or increased batch sizes on existing hardware, potentially accelerating innovation across various machine learning applications. This work represents a significant step towards democratizing access to advanced AI models and optimizing their deployment in resource-constrained environments.

1 INTRODUCTION

The field of machine learning has undergone unprecedented growth in recent years, with neural network models at the forefront of this revolution. These models, spanning domains from natural language processing to computer vision, have demonstrated remarkable capabilities in tackling complex tasks. However, their increasing size and complexity present significant challenges for execution, particularly in resource-constrained environments. State-of-the-art models such as LLaMA-3 Dubey et al. (2024) and PaLM 2 Anil et al. (2023) now comprise hundreds of billions of parameters, pushing the boundaries of what’s possible in language understanding and generation. While these models achieve unprecedented performance across a wide range of tasks, they also demand substantial computational resources, straining the limits of current hardware capabilities. As model parameters reach into the hundreds of billions, the constraints of GPU memory become a critical bottleneck, especially during both training and inference tasks on consumer-grade hardware. This growing disparity between model size and available computational resources presents a pressing challenge for the machine learning community, necessitating innovative solutions for efficient model execution, training, and deployment.

The machine learning community has made significant strides in optimizing model training on high-performance computing clusters. Techniques such as model parallelism Shoeybi et al. (2019), which distributes model layers across multiple devices, and data parallelism, which processes different

054 batches of data on separate devices, have been crucial in scaling up model sizes. Recent advance-
055 ments like Fully Sharded Data Parallel (FSDP) Zhao et al. (2023) and Distributed Data Parallel
056 (DDP) Li et al. (2020) have further improved training efficiency by optimizing memory usage and
057 communication patterns. FSDP, in particular, allows for training larger models by sharding param-
058 eters, gradients, and optimizer states across data parallel workers. However, while these techniques
059 have revolutionized training capabilities, they primarily address the needs of institutions with access
060 to substantial computational resources. For the broader user base, both training and inference —
061 the process of deploying trained models to make predictions on new data — have become increasingly
062 challenging, especially when such models must run on consumer-grade hardware or edge devices
063 with constrained computational resources. Even when high-end hardware is available, AI practition-
064 ers often run into out of memory (OOM) issues when dealing with large batch sizes that are critical
065 for producing high-performance models.

066 Recent advances in model optimization have addressed the challenges of working with large models,
067 focusing on both efficient training and inference. Model segmentation and partitioning techniques,
068 such as GPipe Huang et al. (2019) and Megatron-LM Shoeybi et al. (2019), enable the distribu-
069 tion of large models across multiple accelerators. Dynamic memory management strategies, like
070 the Zero Redundancy Optimizer (ZeRO) Rajbhandari et al. (2020) and SuperNeurons Wang et al.
071 (2018), optimize memory usage during training by minimizing data redundancy and efficiently man-
072 aging intermediate activations. Pipelined execution methods such as PipeDream Narayanan et al.
073 (2019) and TeraPipe Li et al. (2021) have shown considerable promise in improving throughput for
074 distributed training. In the realm of inference, recent work has made significant strides in address-
075 ing efficiency challenges. Alizadeh et al. Alizadeh et al. (2023) propose an innovative method to
076 run LLMs on devices with limited DRAM capacity by utilizing flash memory for model storage.
077 The FlexGen system by Sheng et al. Sheng et al. (2023) addresses the challenge of running LLMs
078 on a single commodity GPU with limited memory by utilizing a combination of GPU, CPU, and
079 disk storage. While these advancements represent significant progress, many existing techniques
080 are specifically tailored for LLMs and may not generalize well to other types of neural network ar-
081 chitectures. Additionally, some approaches may produce outputs that differ from the original model,
potentially affecting performance and reliability.

082 In this paper, we present Superpipeline, a novel approach designed to overcome the limitations
083 associated with executing and training large neural network models on limited hardware resources.
084 Our method synthesizes and extends existing concepts to formulate a comprehensive framework
085 that addresses both memory constraints and execution efficiency, while maintaining three crucial
086 advantages. First, our approach ensures perfect fidelity to the original model’s output, guaranteeing
087 that the results of both training and inference are identical to those produced by the unmodified
088 model. Second, our method is designed for versatility, easily adaptable to a wide range of neural
089 network architectures beyond just LLMs. This broad applicability makes our solution relevant across
090 various domains and model types. Third, we prioritize ease of use, allowing for straightforward
091 implementation without the need for complex model modifications or specialized hardware setups.
092 Referring to Figure 1, Superpipeline is reminiscent of the super pipelining technique in computer
093 architecture Gaudiot et al. (2005). Superpipeline breaks a model into units and load k units into
094 GPU memory initially. Once a preset number, k' , where $k' < k$, of units have been executed in
095 GPU, they are offloaded back to CPU to make space for the next k' units while $k - k'$ units are still
executing in GPU. The key contributions of Superpipeline can be summarized as follows:

- 096 1. **Efficient Training and Inference:** Our method enhances both training and inference
097 phases, ensuring optimized execution on single GPU environments. It addresses the critical
098 need for efficiently training and deploying large models in resource-constrained scenarios.
- 099 2. **No Model Retraining or Parameter Modification:** Our method works without intro-
100 ducing any new parameters to the model, ensuring that no retraining is required. This
101 guarantees that both the model structure and its output remain identical to the original.
- 102 3. **Universal Applicability:** We present a versatile approach that is easily adaptable to various
103 neural network architectures, from LLMs to image generation models like Stable Diffusion,
104 without requiring model-specific modifications.
- 105 4. **Simplified Implementation and Broad GPU Compatibility:** Our method is designed for
106 straightforward implementation, requiring no complex modifications or specialized hard-
107 ware setups. Additionally, unlike methods such as FlashAttention Dao et al. (2022), which

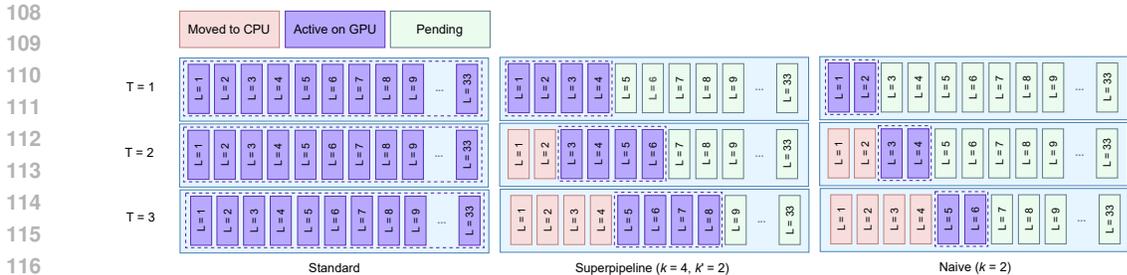


Figure 1: Superpipeline Diagram. Comparison of model execution strategies: Standard (all layers on GPU), Naive ($k = 2$), and Superpipeline ($k = 4, k' = 2$). k represents layers simultaneously on GPU. k' denotes layers transferred back to CPU after computation, and simultaneously, the number of next layers moved to GPU. Superpipeline optimizes GPU memory usage through this dynamic layer management.

are limited to Ampere GPUs, our approach is compatible with any GPU architecture, providing greater flexibility and accessibility across various hardware setups.

By focusing on these key aspects, Superpipeline offers a practical and efficient solution for training and deploying large models on memory-constrained devices, effectively balancing computational load and memory availability to maximize performance without sacrificing accuracy or generalizability. This has profound implications for various applications, including edge computing, mobile applications, large batch-sized training recipes, and other scenarios where access to high-end computing resources is limited. By enabling the training and deployment of advanced neural networks on such devices, our method can help bridge the gap between cutting-edge AI research and practical, everyday applications as well as ensuring equity among common AI practitioners and well-endowed institutions alike.

The remainder of this paper is organized as follows: Section 2 provides a comprehensive review of related work in model optimization and efficient training and inference techniques. Section 3 details our proposed method, emphasizing its universality and output fidelity preservation for both training and inference phases. Section 4 presents our experimental results across various model types and tasks, demonstrating the effectiveness of Superpipeline in both training and inference scenarios. We conclude in Section 6 with a summary of our findings and their potential impact on democratizing access to state-of-the-art AI models.

2 RELATED WORK

Recent advancements in neural network research have focused on enhancing the efficiency and scalability of large models, particularly in environments with limited hardware resources. This section reviews key developments in model compression, memory management, parallelism strategies, and data transfer optimization techniques relevant to our proposed method.

2.1 MODEL SEGMENTATION AND PARTITIONING

The concept of dividing large models into smaller, manageable units has gained prominence in recent years. GPipe Huang et al. (2019) introduced a scalable model-parallelism library that efficiently trains large neural networks using pipeline parallelism. By partitioning deep networks into smaller segments and distributing them across different accelerators, GPipe optimizes hardware utilization and reduces training time. To maintain the simplicity of the proposed method and ensure its generalizability across different models, we use the repetitive layers present in every deep model as the model’s partition for memory management.

Megatron-LM Shoeybi et al. (2019) proposed an intra-layer model parallelism technique that efficiently trains large-scale Transformer-based language models by distributing computations across multiple GPUs. While this method enhances scalability for training, our approach adapts these prin-

162 ciples for single-GPU environments, focusing on dynamic partitioning and memory management to
163 optimize inference and training processes.
164

165 2.2 MODEL COMPRESSION AND SELECTIVE EXECUTION 166

167 As large language models (LLMs) increase in size, reducing their computational and memory re-
168 quirements has become a critical area of research. Model compression techniques such as pruning
169 and quantization have been extensively explored to shrink models without significantly compromis-
170 ing performance Han et al. (2015); Jaiswal et al. (2023); Ahmadian et al. (2023); Li et al. (2024).
171 Additionally, selective execution methods, including sparse activations and conditional computa-
172 tion Zhang et al. (2024); Baykal et al. (2024), aim to reduce the computational overhead by limiting
173 operations to necessary components, which aligns with the broader goal of minimizing resource
174 usage during inference.

175 Selective weight loading is another related concept, where techniques have been developed to dy-
176 namically load a subset of weights based on activation patterns Liu et al. (2023); Sheng et al. (2023).
177 This strategy reduces the memory footprint required for model execution, complementing efforts to
178 manage memory transfers between different hardware components effectively.
179

180 2.3 DYNAMIC MEMORY MANAGEMENT AND HARDWARE OPTIMIZATION 181

182 Dynamic memory management strategies have been proposed to address GPU memory limitations
183 in training and deploying deep neural networks. The Zero Redundancy Optimizer (ZeRO) Rajb-
184 handari et al. (2020) optimizes memory usage by eliminating redundant copies of model states and
185 distributing them across devices. This method has parallels to dynamic memory management strate-
186 gies used to optimize memory allocation for inference, particularly in settings with limited hardware
187 resources.

188 Hardware optimization techniques, including efficient memory architectures Gao et al. (2019) and
189 dataflow optimizations Han et al. (2016), also contribute to more efficient LLM inference. These
190 methods can further enhance algorithmic improvements for memory management and model execu-
191 tion by leveraging hardware-specific optimizations.

192 2.4 PIPELINED EXECUTION AND SPECULATIVE TECHNIQUES 193

194 Pipelined execution has been a focus of several studies aimed at improving deep neural network
195 (DNN) training throughput. PipeDream Narayanan et al. (2019) and TeraPipe Li et al. (2021) ex-
196 plore combining intra-batch and inter-batch parallelism to optimize training processes across mul-
197 tiple GPUs. In contrast, adaptations of pipelining principles for single-GPU environments have
198 also been proposed to enhance inference efficiency, where models are partitioned and dynamically
199 transferred between memory hierarchies to optimize execution speed.

200 Speculative execution, a technique used to manage latency in model inference, has been explored
201 in various contexts, including speculative decoding for LLMs Zhang et al. (2023); He et al. (2023).
202 This approach utilizes draft models to predict outputs and verifies them with larger models, serving
203 as an orthogonal strategy to improve inference efficiency. Speculative techniques and adaptive exe-
204 cution methods contribute to the growing toolbox for managing the complexity of large models on
205 constrained hardware.
206

207 2.5 TRANSFER STRATEGIES AND PIPELINE OPTIMIZATION 208

209 Optimizing data transfer between different memory hierarchies is a critical yet underexplored area
210 for efficient large model inference. Research on minimizing memory usage through optimal check-
211 pointing and data movement Feng & Huang (2021) provides a foundation for strategies that aim to
212 reduce data transfer overhead during model execution. Techniques that streamline these transfers are
213 essential for executing large models effectively, particularly on devices with limited GPU or DRAM
214 capacity.

215 In contrast to previous works, which primarily target specific model types like LLMs or focus on
optimizing either the training or inference phase, Superpipeline is versatile and applicable across

216 a wide range of models. It seamlessly integrates into both the training and inference processes
217 without altering the original model’s output, making it a simple yet effective solution for enhancing
218 efficiency on resource-constrained hardware. Additionally, it provides AI researchers with an effi-
219 cient solution for developing their models on high-end GPUs, enabling the use of larger batch sizes
220 while optimizing resource utilization.

222 3 PROPOSED METHOD: SUPERPIPELINE

224 This section introduces Superpipeline, our novel approach for efficient execution of large neural
225 network models on constrained hardware resources. Superpipeline addresses the challenge of run-
226 ning memory-intensive models on limited GPU hardware through dynamic memory management
227 and optimized data transfer strategies.

229 3.1 CONCEPTUAL FRAMEWORK

231 Our method segments large models into manageable units based on their repetitive structure. This
232 approach, applicable to various neural network architectures, enables efficient processing and dy-
233 namic memory management. By exploiting the inherent repetition in modern models, we achieve
234 simplicity in implementation and universality across model types.

235 3.2 KEY COMPONENTS OF SUPERPIPELINE

237 3.2.1 MODEL SEGMENTATION STRATEGY

239 We partition neural networks along their natural repetitive boundaries, such as transformer layers in
240 language models or convolutional blocks in vision models. Each repetitive unit becomes a distinct
241 partition. This strategy requires minimal modification to the original architecture, adapts to differ-
242 ent model sizes, and preserves model behavior. For example, a model like LLaMA-2 7B with 32
243 repeating layers would yield 32 partitions. This approach forms the foundation for our subsequent
244 optimization techniques, allowing efficient resource management across diverse model types.

245 3.2.2 DYNAMIC GPU-CPU PARTITION TRANSFER

247 Superpipeline employs a dynamic approach to memory management. Only specific partitions are
248 loaded onto the GPU as needed, and once computation is complete, their outputs are transferred
249 back to CPU memory. This process frees up GPU memory for subsequent partitions, allowing for
250 the processing of models that would otherwise exceed available hardware capacity. This dynamic
251 transfer mechanism is crucial for optimizing GPU resource utilization. It allows larger models to be
252 run on more constrained hardware by effectively managing the limited GPU memory available.

253 3.3 THE SUPERPIPELINE ALGORITHM

255 Superpipeline introduces two critical hyperparameters: k , representing the number of partitions
256 simultaneously on the GPU, and k' , which denotes the number of partitions transferred back to the
257 CPU after computation, making room for the next k' partitions. Figure 1 illustrates this.

258 By adjusting these parameters, the Superpipeline framework achieves an optimal balance between
259 GPU memory usage and processing speed. Increasing k maximizes GPU utilization and accelerates
260 computation, but also raises memory requirements. On the other hand, decreasing k lowers memory
261 usage while slowing down execution. This flexibility allows the method to be tailored to specific
262 hardware constraints, optimizing the trade-off between speed and memory efficiency.

264 In the training phase, Superpipeline extends its benefits to both the forward and backward passes.
265 During the forward pass, it dynamically transfers partitions between GPU and CPU as needed. The
266 same process is repeated for the backward pass, where gradients are computed. This dual application
267 in both forward and backward passes results in even greater reductions in GPU memory usage while
268 maintaining acceptable performance.

269 By efficiently managing memory across both phases of training, Superpipeline significantly reduces
the overall GPU memory footprint, particularly in large-scale models. This method ensures that even

270 resource-constrained hardware can support models that would otherwise be unmanageable, without
271 sacrificing speed or accuracy.

272 273 274 4 EXPERIMENTS AND RESULTS

275
276 In this section, we present our experimental methodology and key findings. Our experimental setup
277 encompasses a diverse array of models, ranging from vision architectures to language models, all
278 implemented using Superpipeline. This broad selection demonstrates the versatility and wide appli-
279 cability of our proposed method. We begin by outlining the implementation details and experimental
280 parameters, followed by a comprehensive description of the models tested. Through these experi-
281 ments, we aim to demonstrate two critical aspects of our approach: first, its ability to reduce GPU
282 memory usage significantly, and second, its capacity to maintain acceptable inference times across
283 various model types. Our experiments are designed to illustrate not only the ease with which our
284 approach can be adapted to various model architectures and domains but also its effectiveness in
285 optimizing resource utilization without substantially compromising performance.

286 287 4.1 EXPERIMENTAL SETUP

288
289 **Models.** To demonstrate that the method presented in this work is applicable to any model, we have
290 conducted our evaluation across different categories of models. We have selected three different
291 models from three distinct domains. One is the llama2 model from the world of LLM (Large Lan-
292 guage Models), the SD model from the world of VLM (Vision Language Models), and ViT-bigG
293 from the world of vision models. We perform our evaluations of the proposed method in two sec-
294 tions: during inference time and during training time. The aim of these experiments is to show the
295 extent to which the proposed method helps in optimizing GPU consumption and how much faster it
296 is compared to the naive approach.

297 **Hardware Configuration.** We evaluated models on three distinct hardware configurations to ensure
298 the generalizability of our method across various devices. The first setup featured a Quadro 8000
299 graphics card with 50 GB of GPU memory. The second configuration utilized an NVIDIA GTX
300 3090 graphics card, offering 24 GB of GPU RAM. Our third setup employed an H20 graphics card
301 with a substantial 98 GB of GPU RAM. By conducting evaluations across these diverse hardware
302 environments, we aimed to validate the robustness and adaptability of our approach

303 304 4.2 RESULTS

305 Our experiments evaluated Superpipeline across four distinct modes of operation:

- 306
307 1. **Standard mode:** The entire model is loaded onto the GPU and processed, representing the
308 conventional approach for model execution.
- 309
310 2. **Naive method:** The model is loaded onto the GPU k layers at a time, offering a simple but
311 potentially inefficient way to reduce memory usage.
- 312
313 3. **CPU-only mode:** The entire model runs on the CPU without GPU acceleration, providing
314 a baseline for comparison in resource-constrained environments.
- 315
316 4. **Superpipeline method:** Our proposed approach for dynamic memory management, bal-
317 ancing GPU utilization and processing efficiency.

318 The key metrics we focused on were GPU Memory Usage and Processing Time. GPU Memory
319 Usage, measured in gigabytes (GB), shows how efficiently each method utilizes available GPU
320 memory. Processing Time, measured in milliseconds (ms) for inference tasks and iterations per
321 second for training tasks, reflects the speed of each method. It's important to note that Superpipeline,
322 by design, does not alter the model's computations or outputs in any way. The results produced by
323 Superpipeline are *identical to those of the standard mode*, ensuring perfect fidelity to the original
model's performance and accuracy.

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

Table 1: Superpipeline Performance During Inference

Model	Method	GPU Usage (GB)	Time (ms)	K	K'
ViT-bigG	Standard	13.7	37.5 ms /embed	-	-
	CPU-only	0	9250 ms /embed	-	-
	Naive	4.2	181.5 ms /embed	1	-
	Naive	5.2	175.5 ms /embed	8	-
	Naive	6.0	176.6 ms /embed	15	-
	Superpipeline	4.8	111.5 ms /embed	4	2
	Superpipeline	5.7	109.7 ms /embed	6	3
	Superpipeline	6.0	103.5 ms /embed	10	8
	Superpipeline	7.3	96.8 ms /embed	14	12
	Superpipeline	8.8	90.5 ms /embed	19	16
LlaMA2	Standard	15.0	26 ms /token	-	-
	CPU-only	0	29200 ms /token	-	-
	Naive	2.9	4520 ms /token	1	-
	Naive	3.8	4000 ms /token	8	-
	Naive	3.8	3980 ms /token	8	-
	Naive	6.5	3970 ms /token	15	-
	Superpipeline	4.7	2053 ms /token	4	2
	Superpipeline	5.7	1964 ms /token	5	3
	Superpipeline	8.0	1748 ms /token	8	3
	Superpipeline	9.2	1607 ms /token	10	2
Stable Diffusion	Standard	6.3	10 s /image	-	-
	CPU-only	2.5	529 s /image	-	-
	Naive	2.5	60 s /image	1	-
	Naive	3.8	59 s /image	8	-
	Superpipeline	2.9	33 s /image	5	3
	Superpipeline	3.3	27 s /image	5	4
	Superpipeline	4.0	27 s /image	8	6
	Superpipeline	4.1	22 s /image	7	5
	Superpipeline	4.4	19 s /image	8	2
	Superpipeline	4.8	16 s /image	9	3
Superpipeline	5.0	14 s /image	10	2	

4.2.1 INFERENCE TIME

One of the significant advantages of our proposed method is its ease of implementation across various existing models by making necessary changes in the forward pass. Since no parameters are added to or removed from the model, and no changes are made to the overall model structure, Superpipeline can be applied to many current models without the need for retraining.

The primary parameters in this approach are K and k' . These values can be easily optimized through a grid search, tailored to the hardware on which the model is running. This flexibility allows for adjusting GPU consumption during inference by simply modifying k and k' . Consequently, any remaining GPU space can be utilized for processing larger batch sizes if required.

The effectiveness of Superpipeline during inference is demonstrated in Table 1. These results highlight the method’s capability to optimize GPU usage without compromising model performance, making it a versatile solution for both training and inference stages.

As shown in Table 1, Superpipeline achieves significant reductions in GPU usage and inference time while maintaining the same accuracy as the standard and naive methods. This demonstrates the

method’s efficiency in resource utilization during the inference phase. Table 4 shows results from the Quadro GPU, with other GPU results in the appendix.

The adaptability of Superpipeline to different hardware configurations and model architectures, combined with its performance benefits in both training and inference, positions it as a valuable tool for optimizing deep learning workflows across various applications and deployment scenarios.

4.2.2 TRAINING TIME

Unlike some previous methods that are only applicable during the inference stage, Superpipeline can be used in both training and inference phases with minimal modifications.

Implementing Superpipeline involves applying this method to the forward section of each model. By utilizing `pre_backward_hook` and `post_backward_hook` functions, Superpipeline can be easily integrated into the model training phase. This capability is particularly significant during training, as gradients are calculated in addition to the usual computations. In these conditions, the efficiency of our proposed method in optimizing GPU usage becomes even more pronounced.

A key feature of Superpipeline is the preservation of model accuracy even when used in training. Since no changes are made to the computational values, the model’s output using our proposed method is identical to that of the standard approach. This distinguishes Superpipeline from methods that rely on predicting which neurons or layers will be used, which may lead to prediction errors and changes in output.

To rigorously evaluate the proposed method, we compared the ViT-BigG model on the imagenet-tiny dataset using identical hyperparameters (batch size, learning rate, number of epochs). The

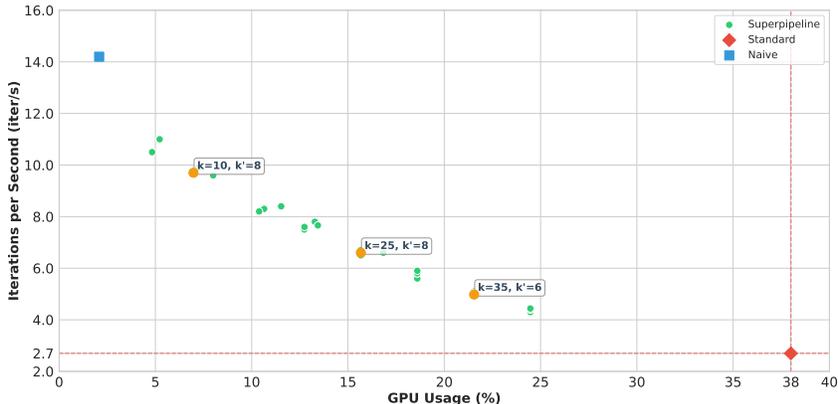


Figure 2: Comparison of memory usage and speed during ViT-BigG training on ImageNet-tiny.

results of this comparison are shown in Figure 2. As observed, the Superpipeline method not only significantly reduces GPU consumption but also provides a highly acceptable speed compared to the standard mode and the naive method.

Superpipeline offers several notable advantages. It’s implementation process for training remains straightforward, and can be applied to various types of models. Additionally, by adjusting the parameters k and k' , GPU consumption can be easily controlled. By optimizing GPU usage, it becomes possible to train models with larger batch sizes, which can lead to improved performance and faster model convergence.

4.2.3 BENEFIT OF DIFFERENT GPU USAGE

While Superpipeline offers significant benefits for general users, its impact on AI research and model development is particularly noteworthy. As deep learning models continue to grow in size and complexity, GPU memory constraints have become a critical bottleneck in the training process. Even with high-capacity GPUs boasting 50 to 100 gigabytes of memory, researchers face limitations in increasing batch sizes, a crucial factor for many advanced training techniques.

Table 2: GPU Usage for ViT-BigG and LLaMA2 Models w/ and w/o Gradient Checkpointing. OOM: Out-Of-Memory.

Model	Method	With Grad. Checkpointing				Without Grad. Checkpointing			
		BS	GPU	BS	GPU	BS	GPU	BS	GPU
ViT-BigG (Fully Trainable) (on Quadro)	Superpipe ($k=6, k'=3$)	16	10.1	64	25.8	4	23.0	10	42.3
		32	15.4	128	42	8	36.9	12	48.0
LLaMA2 (Fully Trainable) (on H20)	Superpipe ($k=6, k'=3$)	32	33.5	128	53.6	16	55.3	64	OOM
		64	39.5	256	81.8	32	85.6	128	OOM
LLaMA2 (Half of Layers Frozen) (on H20)	Superpipe ($k=6, k'=3$)	4k	21	16k	48	2k	29.8	8k	73
		8k	30	32k	88.8	4k	44	10k	88.3
	Standard	32	OOM	128	OOM	16	OOM	64	OOM
		64	OOM	256	OOM	32	OOM	128	OOM
	Standard	4k	36	16k	64	2k	64	8k	OOM
		8k	45	32k	OOM	4k	79	10k	OOM

As shown in Table 2, Superpipeline significantly expands the potential for larger batch sizes during model training. For instance, when training the LLaMA2 model without gradient checkpointing, the standard approach fails due to out-of-memory errors even at smaller batch sizes. In contrast, Superpipeline successfully trains the model with larger batch sizes, demonstrating its ability to handle scenarios infeasible with standard training methods.

To provide a more equitable comparison and further demonstrate Superpipeline’s advantages in enabling larger batch sizes, we conducted an additional experiment where half of the LLaMA2 model’s layers were frozen. This approach allowed the standard method to handle larger batch sizes, creating a more balanced comparison scenario. In this setting, Superpipeline continued to outperform, accommodating significantly larger batch sizes and achieving more efficient GPU utilization.

By alleviating memory constraints, Superpipeline enables the exploration of training regimes that were previously infeasible, potentially accelerating advancements in areas such as self-supervised learning, large-scale visual representation learning, and the training of Large Language Models. This adaptability is crucial in an era where model innovation often outpaces hardware advancement, allowing researchers with limited resources to work on cutting-edge models and training techniques previously exclusive to well-resourced institutions.

5 LIMITATIONS AND FUTURE WORK

An examination of Table 1 reveals that while the superpipeline method consistently outperforms the naive approach across all models, the performance gap between the proposed superpipeline method and the standard approach is notably smaller for the ViT-bigG model compared to models like Llama2. To investigate this discrepancy, we measured two distinct timings for both the ViT-bigG and Llama2 models: the time required to transfer a layer to the GPU and the time needed to transfer a layer to the CPU. The results are illustrated in Figure 3. Two key observations can be drawn from Figure 3. First, the transfer time of a layer to the CPU is slower than to the GPU. Second, and more importantly, we observe that the transfer speed of a single layer from the Llama model is significantly slower than that of the ViT-bigG model. This difference explains the larger performance gap between the superpipeline and standard approaches in the Llama model.

In essence, when considering a single forward pass, the superpipeline and standard methods do not differ significantly. However, since we calculate model speed based on an average of multiple consecutive forward passes, a limitation becomes apparent in the superpipeline approach. Although the model’s output is quickly generated in the first forward pass, it cannot immediately produce the second output as it must wait for the layers from the previous forward pass to complete their transfer to the CPU. This issue represents a key limitation of our work. In scenarios where the layer transfer speed to the CPU is slow for a particular model or hardware configuration, the superpipeline method, while still outperforming the naive approach, may not achieve performance parity with the standard method.

Several potential solutions to address this limitation could be explored in future work. One approach involves rewriting the model transfer function to the CPU using CUDA custom kernels. Another possibility is developing a faster method for creating and transferring a copy of each layer to the GPU. This approach would eliminate the need to transfer layers back to the CPU after GPU processing, instead overwriting the previous layer directly on the GPU. Currently, implementing this with existing PyTorch features is significantly more time-consuming than transferring a layer to the CPU, necessitating a more optimized implementation. In future research, we plan to explore these optimization strategies to further enhance the performance of the superpipeline method across a wider range of models and hardware configurations. Additionally, we aim to investigate the applicability of our approach to emerging model architectures and to develop adaptive strategies that can automatically adjust the superpipeline parameters based on the specific characteristics of the model and hardware in use.

6 CONCLUSION

In this paper, we introduced Superpipeline, a novel method for efficient execution of large neural network models on constrained hardware resources. Our approach addresses the critical challenge of deploying and training increasingly complex models in environments with limited GPU memory, without compromising model performance or accuracy. The key strengths of Superpipeline lie in its versatility and ease of implementation. Unlike previous methods that primarily focused on LLM models or were limited to inference stages, Superpipeline demonstrates broad applicability across various model architectures, including LLMs, VLMs, and vision-based models. Moreover, it can be seamlessly integrated into both inference and training pipelines, offering a comprehensive solution for resource optimization throughout the model lifecycle. A significant advantage of our method is its ability to substantially reduce GPU memory consumption while maintaining acceptable execution speeds. This is achieved without adding new parameters to the model or requiring retraining, ensuring that the model’s output in Superpipeline mode remains identical to that in standard mode. This preservation of accuracy sets Superpipeline apart from other optimization techniques that may introduce performance trade-offs.

Our experimental results across diverse model types and hardware configurations validate the effectiveness of Superpipeline. We demonstrated significant reductions in GPU usage during both inference and training, with minimal impact on processing speed. The method’s adaptability to different hardware setups further enhances its practical value, making it a viable solution for a wide range of deployment scenarios. The simplicity of Superpipeline’s implementation, coupled with its flexibility in fine-tuning through the k and k' parameters, positions it as a powerful tool for researchers and practitioners alike. By optimizing resource utilization, our method opens up new possibilities for working with larger models or increased batch sizes on existing hardware, potentially accelerating research and development in the field of deep learning.

In conclusion, Superpipeline represents a significant step forward in making advanced AI models more accessible and efficient to deploy. As the complexity of neural networks continues to grow, methods like Superpipeline will play a crucial role in bridging the gap between state-of-the-art model architectures and the practical constraints of real-world computing environments. Future work could explore further optimizations and extensions of this approach, potentially leading to even more efficient and scalable solutions for large-scale model deployment and training. You may include other additional sections here.

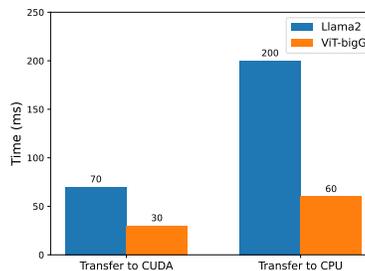


Figure 3: Comparison of layer transfer times between GPU and CPU for ViT-bigG and Llama2 models.

REFERENCES

- 540
541
542 Arash Ahmadian, Saurabh Dash, Hongyu Chen, Bharat Venkitesh, Zhen Stephen Gou, Phil Blun-
543 som, Ahmet Üstün, and Sara Hooker. Intriguing properties of quantization at scale. *Advances in*
544 *Neural Information Processing Systems*, 36:34278–34294, 2023.
- 545 Keivan Alizadeh, Iman Mirzadeh, Dmitry Belenko, Karen Khatamifard, Minsik Cho, Carlo C
546 Del Mundo, Mohammad Rastegari, and Mehrdad Farajtabar. Llm in a flash: Efficient large lan-
547 guage model inference with limited memory. *arXiv preprint arXiv:2312.11514*, 2023.
- 548 Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos,
549 Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. Palm 2 technical report.
550 *arXiv preprint arXiv:2305.10403*, 2023.
- 551 Cenk Baykal, Dylan Cutler, Nishanth Dikkala, Nikhil Ghosh, Rina Panigrahy, and Xin Wang. Al-
552 ternating updates for efficient transformers. *Advances in Neural Information Processing Systems*,
553 36, 2024.
- 554 Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-
555 efficient exact attention with io-awareness. *Advances in Neural Information Processing Systems*,
556 35:16344–16359, 2022.
- 557 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
558 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.
559 *arXiv preprint arXiv:2407.21783*, 2024.
- 560 Jianwei Feng and Dong Huang. Optimal gradient checkpoint search for arbitrary computation
561 graphs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recog-
562 nition*, pp. 11433–11442, 2021.
- 563 Fei Gao, Georgios Tziantzioulis, and David Wentzlaff. Computedram: In-memory compute using
564 off-the-shelf drams. In *Proceedings of the 52nd annual IEEE/ACM international symposium on
565 microarchitecture*, pp. 100–113, 2019.
- 566 Jean-Luc Gaudiot, Jung-Yup Kang, and Won Woo Ro. Techniques to improve performance beyond
567 pipelining: Superpipelining, superscalar, and vliw. volume 63 of *Advances in Computers*, pp.
568 1–34. Elsevier, 2005. doi: [https://doi.org/10.1016/S0065-2458\(04\)63001-4](https://doi.org/10.1016/S0065-2458(04)63001-4). URL [https://
569 www.sciencedirect.com/science/article/pii/S0065245804630014](https://www.sciencedirect.com/science/article/pii/S0065245804630014).
- 570 Song Han, Huizi Mao, and William J Dally. Deep compression: Compressing deep neural networks
571 with pruning, trained quantization and huffman coding. *arXiv preprint arXiv:1510.00149*, 2015.
- 572 Song Han, Xingyu Liu, Huizi Mao, Jing Pu, Ardavan Pedram, Mark A Horowitz, and William J
573 Dally. Eie: Efficient inference engine on compressed deep neural network. *ACM SIGARCH
574 Computer Architecture News*, 44(3):243–254, 2016.
- 575 Zhenyu He, Zexuan Zhong, Tianle Cai, Jason D Lee, and Di He. Rest: Retrieval-based speculative
576 decoding. *arXiv preprint arXiv:2311.08252*, 2023.
- 577 Yanping Huang, Youlong Cheng, Ankur Bapna, Orhan Firat, Dehao Chen, Mia Chen, Hyoungho
578 Lee, Jiquan Ngiam, Quoc V Le, Yonghui Wu, et al. Gpipe: Efficient training of giant neural
579 networks using pipeline parallelism. *Advances in neural information processing systems*, 32,
580 2019.
- 581 Ajay Jaiswal, Zhe Gan, Xianzhi Du, Bowen Zhang, Zhangyang Wang, and Yinfei Yang. Compress-
582 ing llms: The truth is rarely pure and never simple. *arXiv preprint arXiv:2310.01382*, 2023.
- 583 Liang Li, Qingyuan Li, Bo Zhang, and Xiangxiang Chu. Norm tweaking: High-performance low-
584 bit quantization of large language models. In *Proceedings of the AAAI Conference on Artificial
585 Intelligence*, volume 38, pp. 18536–18544, 2024.
- 586 Shen Li, Yanli Zhao, Rohan Varma, Omkar Salpekar, Pieter Noordhuis, Teng Li, Adam Paszke, Jeff
587 Smith, Brian Vaughan, Pritam Damania, et al. Pytorch distributed: Experiences on accelerating
588 data parallel training. *arXiv preprint arXiv:2006.15704*, 2020.
- 589
590
591
592
593

- 594 Zhuohan Li, Siyuan Zhuang, Shiyuan Guo, Danyang Zhuo, Hao Zhang, Dawn Song, and Ion Stoica.
595 Terapipe: Token-level pipeline parallelism for training large-scale language models. In *International
596 Conference on Machine Learning*, pp. 6543–6552. PMLR, 2021.
597
- 598 Zichang Liu, Jue Wang, Tri Dao, Tianyi Zhou, Binhang Yuan, Zhao Song, Anshumali Shrivastava,
599 Ce Zhang, Yuandong Tian, Christopher Re, et al. Deja vu: Contextual sparsity for efficient llms
600 at inference time. In *International Conference on Machine Learning*, pp. 22137–22176. PMLR,
601 2023.
- 602 Deepak Narayanan, Aaron Harlap, Amar Phanishayee, Vivek Seshadri, Nikhil R Devanur, Gre-
603 gory R Ganger, Phillip B Gibbons, and Matei Zaharia. Pipedream: Generalized pipeline par-
604 allelism for dnn training. In *Proceedings of the 27th ACM symposium on operating systems
605 principles*, pp. 1–15, 2019.
- 606 Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yuxiong He. Zero: Memory optimizations
607 toward training trillion parameter models. In *SC20: International Conference for High Perfor-
608 mance Computing, Networking, Storage and Analysis*, pp. 1–16. IEEE, 2020.
609
- 610 Ying Sheng, Lianmin Zheng, Binhang Yuan, Zhuohan Li, Max Ryabinin, Beidi Chen, Percy Liang,
611 Christopher Ré, Ion Stoica, and Ce Zhang. Flexgen: High-throughput generative inference of
612 large language models with a single gpu. In *International Conference on Machine Learning*, pp.
613 31094–31116. PMLR, 2023.
- 614 Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan
615 Catanzaro. Megatron-lm: Training multi-billion parameter language models using model par-
616 allelism. *arXiv preprint arXiv:1909.08053*, 2019.
- 617 Linnan Wang, Jinmian Ye, Yiyang Zhao, Wei Wu, Ang Li, Shuaiwen Leon Song, Zenglin Xu,
618 and Tim Kraska. Superneurons: Dynamic gpu memory management for training deep neural
619 networks. In *Proceedings of the 23rd ACM SIGPLAN symposium on principles and practice of
620 parallel programming*, pp. 41–53, 2018.
621
- 622 Jun Zhang, Jue Wang, Huan Li, Lidan Shou, Ke Chen, Gang Chen, and Sharad Mehrotra. Draft &
623 verify: Lossless large language model acceleration via self-speculative decoding. *arXiv preprint
624 arXiv:2309.08168*, 2023.
- 625 Zhengyan Zhang, Yixin Song, Guanghui Yu, Xu Han, Yankai Lin, Chaojun Xiao, Chenyang Song,
626 Zhiyuan Liu, Zeyu Mi, and Maosong Sun. Relu 2 wins: Discovering efficient activation functions
627 for sparse llms. *arXiv preprint arXiv:2402.03804*, 2024.
628
- 629 Yanli Zhao, Andrew Gu, Rohan Varma, Liang Luo, Chien-Chin Huang, Min Xu, Less Wright,
630 Hamid Shojanazeri, Myle Ott, Sam Shleifer, et al. Pytorch fsdp: experiences on scaling fully
631 sharded data parallel. *arXiv preprint arXiv:2304.11277*, 2023.
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

A APPENDIX

A.1 RESULTS ON DIFFERENT GPUS

In this section, we present the results of applying the Superpipeline method on two different GPUs: the NVIDIA RTX 3090 and the H20 (shown in Table 3). The Superpipeline approach was evaluated using three different models—ViT-bigG, LLaMA2, and Stable Diffusion—under varying memory constraints and batch sizes

A.2 OPTIMIZED PARTITION TRANSFER STRATEGY

Our experiments revealed that the method of transferring partitions between GPU and CPU significantly impacts overall performance. We compared two approaches: Sequential Transfer and Batch Transfer. In Sequential Transfer, layers are transferred one-by-one to the GPU and back to the CPU. Batch Transfer, on the other hand, moves all layers to the GPU simultaneously, then back to the CPU as a batch. As illustrated in Figure 4, the batch transfer method proved significantly faster, despite

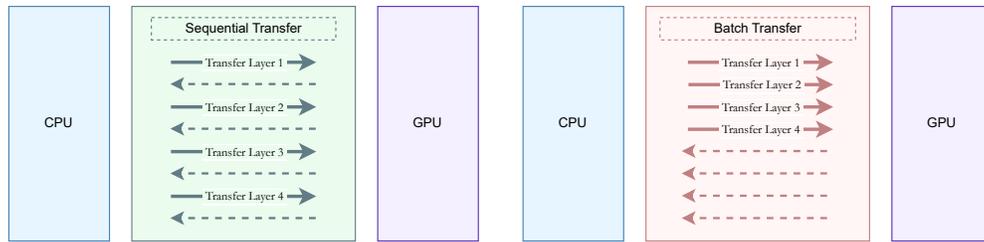


Figure 4: Comparison of Sequential and Batch Transfer Strategies

involving the same number of total transfers. Figure 5 provides empirical evidence of this performance difference across various model architectures. These findings underscore the importance of

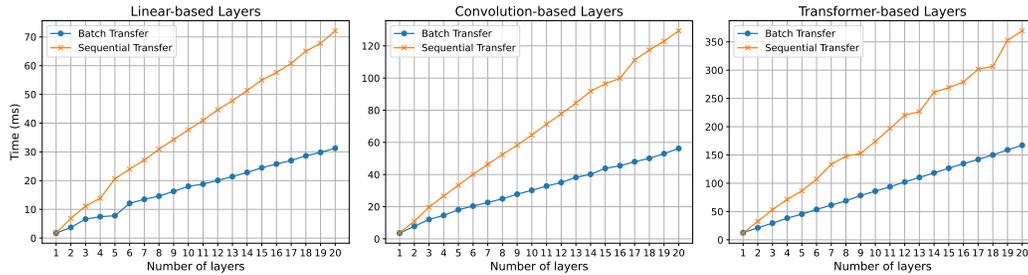


Figure 5: Performance comparison of Sequential vs. Batch Transfer strategies

optimizing not just the partitioning of the model, but also the mechanisms for data transfer between different memory hierarchies.

Table 3: Superpipeline Performance During Inference On RTX 3090

Model	Method	GPU Usage (GB)	Time (ms)	K	K'
ViT-bigG	Superpipeline	4.1	242.8 ms /embed	4	3
	Superpipeline	5.7	223.25 ms /embed	5	3
	Superpipeline	8	212.5 ms /embed	7	5
	Superpipeline	8.8	213.1 ms /embed	9	4
	Superpipeline	11.9	197.5 ms /embed	11	7
LlaMA2	Superpipeline	4.8	108.1 ms /token	4	2
	Superpipeline	5.1	104.6 ms /token	6	4
	Superpipeline	5.5	100.5 ms /token	7	6
	Superpipeline	6.4	98.3 ms /token	11	8
	Superpipeline	7.7	91.2 ms /token	16	12
	Superpipeline	8.6	86.2 ms /token	18	16
Stable Diffusion	Superpipeline	2.9	70 s /image	3	2
	Superpipeline	3.8	54 s /image	6	2
	Superpipeline	4.1	55 s /image	8	5
	Superpipeline	4.7	53 s /image	9	3
	Superpipeline	5.0	49 s /image	11	2

Table 4: Superpipeline Performance During Inference On H20

Model	Method	GPU Usage (GB)	Time (ms)	K	K'
ViT-bigG	Superpipeline	4.7	34.1 ms /embed	4	3
	Superpipeline	5.1	33.2 ms /embed	6	4
	Superpipeline	5.9	32.3 ms /embed	9	6
	Superpipeline	7.2	30.3 ms /embed	14	12
	Superpipeline	8.5	28.8 ms /embed	17	16
LlaMA2	Superpipeline	4.5	41.25 ms /token	4	2
	Superpipeline	5.7	40.6 ms /token	5	3
	Superpipeline	8.0	37.3 ms /token	7	5
	Superpipeline	7.6	36 ms /token	8	2
	Superpipeline	10.3	31.6 ms /token	11	3
	Superpipeline	11.9	30.25 ms /token	12	5
Stable Diffusion	Superpipeline	3.3	11.7 s /image	3	2
	Superpipeline	3.8	9.2 s /image	6	3
	Superpipeline	4.1	8.8 s /image	8	5
	Superpipeline	4.8	8.3 s /image	9	3
	Superpipeline	5.0	7.5 s /image	10	4