(RESEARCH ARTICLE)

# Scalable knowledge distillation for large language models on multi-GPU systems

Wary Hossain Rabby [1, *], A.S.S.M.Q-E-Elahy [1], Gias Uddin [2], Emran Sikder [3], Rafiqul Islam [1] and Hasibul Islam [1]

[1] Jahangirnagar University, Dhaka, Bangladesh.
[2] Uttara University, Dhaka, Bangladesh.
[3] Daffodil International University, Dhaka, Bangladesh.

## Abstract

One well-liked method for condensing massive language models (LLMs) into smaller, faster, more effective versions without sacrificing performance is knowledge distillation (KD). However, it is no longer feasible to run distillation on a single device as LLMs grow to hundreds of billions of parameters; it is simply too computationally demanding. In this paper, we investigate how to leverage multi-GPU configurations to make KD scale. To overcome communication bottlenecks and accelerate training, our method com- bines adaptive gradient compression with tensor, pipeline, and data parallelism. Tested on transformer-based LLMs, our framework maintains strong accuracy for both understanding and generation tasks, reduces communication overhead by 27%, and provides up to 3.4× faster training than single-GPU baselines.

## 1. Introduction

Advanced chatbots and sophisticated reasoning systems are powered by Large Language Models (LLMs), such as GPT-4, PaLM, and LLaMA, which have raised the bar in natural language processing. However, because of their enormous size—often containing over 100 billion parameters—they are extremely difficult to run, requiring enormous amounts of memory, storage, and processing power.

A path forward is provided by knowledge distillation (KD), which reduces the cost of inference while maintaining a large portion of the original performance by teaching a smaller student model to imitate a large teacher model. The difficulty, though, is that distilling such large models requires a lot of resources. The memory and computation demands are simply too great for conventional training configurations that only use one GPU or a small number of GPUs.

To address this issue, we investigate distributed, multi-GPU approaches for large-scale distillation in this work. We suggest a scalable KD framework that makes distillation feasible and effective, even on typical GPU clusters, by combining hybrid parallelism, more intelligent communication techniques, and dynamic workload balancing.

## 2. Related Work

### 2.1. Where It Started

Geoffrey Hinton and his coworkers came up with the idea for KD in 2015. They suggested that instead of only training a small model on raw labels, it should also learn from the soft predictions of a larger model.

---

* Corresponding author: Wary Hossain Rabby

For instance, if the teacher puts an image into one of these categories

- Dog: 95%
- Wolf: 3%
- Cat: 2%

The student learns that the model sees some similarities between dogs and wolves, not just that" this is a dog. This" dark knowledge" helps the student make better generalizations.

## 2.2. Where It Started

This method quickly became popular in NLP. KD made models like DistilBERT, TinyBERT, and MiniLM directly. For example, DistilBERT is 40% smaller than BERT but still has 97% of its accuracy. These models demonstrated that KD is not merely a theoretical concept; it can generate real, deployable systems that operate more efficiently and cost-effectively.

But scaling KD up to modern LLMs (which have hundreds of billions of parameters) is a whole different story. It's not as easy to make a smaller version of GPT-4 as it is to train DistilBERT from BERT. The teacher is so big that just running it slows things down.

## 2.3. Training at Scale: The Rise of Multi-GPU Systems: Why Multi-GPU?

Today's LLMs are too complex for a single GPU, regard- less of its power. A model such as GPT-3 needs terabytes of memory to store its weights, gradients, and optimizer states. For this reason, researchers employ clusters of GPUs cooperating.

### 2.3.1. Key Systems

In recent years, a number of frameworks have enabled ex- tensive training

- **Megatron-LM**: Focuses on separating matrix multiplications (tensor parallelism).
- **DeepSpeed**: renowned for its **ZeRO (Zero Redun- dancy Optimiser)**, which distributes model states among GPUs to prevent any one GPU from bearing the entire load.
- **ZeRO-Offload**: This extension allows for the training of trillion-parameter models by shifting some of ZeRO's memory requirements onto CPUs or even NVMe storage.
- Model training was transformed by these systems. The problem is that they weren't made to distill a large teacher into a small student; rather, they were made for **pretraining one large model**. Distillation presents a dis- tinct set of issues.

## 2.4. Parallelism Strategies: How We Scale Models

We need ways to divide the work when training across GPUs. The three primary ones are

### 2.4.1. Data Parallelism

Each GPU processes a distinct batch of data and keeps a copy of the model. Gradients are averaged across GPUs at each stage.

- **Pros**: Easy to use and extensively supported.
- **Cons**: High memory consumption and sluggish GPU-to-GPU communication.

### 2.4.2. Tensor Parallelism

One large matrix multiplication is broken up into smaller components by dividing layers among GPUs.

- **Pros**: Lowers memory usage per GPU.
- **Cons**: During training, continuous communication is necessary, which can be costly.

### 2.4.3. Pipeline Parallelism

Data is transferred like an assembly line between various GPUs with different sets of layers.

- **Pros**: Excellent for extremely complex models.
- **Cons**: may result in GPUs becoming idle (also known as" pipeline bubbles") while they wait for other GPUs to complete their tasks.
- **Hybrid approaches** are frequently employed in practice, combining pipeline, tensor, and data parallelism to balance performance.

## 2.5. Why KD is Hard to Scale

At the LLM scale, running KD is more than just" training as usual." It presents particular challenges

- **Two Models at Once**: During training, both the teacher and the student must run, thereby doubling their computational and memory requirements.
- **Heavy Communication**: KD frequently necessitates the transfer of significant intermediate states from teacher to student, such as attention maps or hidden layers. These are large tensors, and it is ex- pensive to transfer them between GPUs.
- **Imbalance**: Because teachers are typically slower, students may have to wait, which can lead to bottle- necks.
- **Framework Limitations**: Two-model setups are not supported by current systems such as Megatron or DeepSpeed. To effectively manage KD, they re- quire specific adjustments.
- **Irony of KD**: Although KD promises smaller, less expensive models, if done poorly, the process itself may end up costing more than training the teacher.

## 2.6. Our Approach: Making KD Scalable

We suggest a number of methods to modify distributed frameworks and parallelism especially for KD in order to address these problems

### 2.6.1. Smarter Scheduling

We employ a pipelined approach rather than running teacher and student one after the other. The student begins learning from a batch as soon as the teacher completes it, and the teacher then moves on to the next. GPUs are rarely idle in this way.

### 2.6.2. Splitting Teacher and Student Across GPUs

We assign student layers to one group of GPUs and teacher layers to another. In this manner, they can speak with each other directly without overpowering the entire group.

### 2.6.3. Compressing Teacher Outputs

We use projection layers or quantization to compress the hidden states rather than passing them in their entirety. For instance, a 4096-dimensional vector can be reduced to 1024 dimensions with minimal performance loss, saving a significant amount of bandwidth.

### 2.6.4. Extending ZeRO to KD

To manage teacher-student training, we modify ZeRO- style memory partitioning. This guarantees that both models can fit into GPU memory and prevents redundant storage.

### 2.6.5. Asymmetric Offloading

Parts of it can be offloaded to CPU memory or disc since the teacher doesn't have to update as frequently, while the student remains fully on GPU for speed.

## 2.7. In Summery

- **Knowledge Distillation (KD):** Hinton et al. (2015) first proposed the concept of KD as a means of condensing large neural networks into more manageable and effective ones. This idea gave rise to well-known student models in NLP, including TinyBERT, DistilBERT, and MiniLM.
- **Large-Scale Model Training:** Massive models can now be efficiently trained across multiple GPUs thanks to systems like Megatron-LM, DeepSpeed, and ZeRO- offload. These optimizations, however, are primarily de- signed for pretraining big models rather than the particular difficulties of distillation.

- **Parallelism Strategies:** Scalability is provided by data parallelism, but synchronization costs are frequently high. Although they have their own communication over- head, tensor and pipeline parallelism aid in lowering memory pressure.
- **Our Contribution:** We connect these threads in this work by tailoring multi-GPU optimizations and parallelism for scalable knowledge distillation. This fills the gap between effective distributed training and LLM compression.

## 3. Methodology

### 3.1. Knowledge Distillation Objective

Knowledge distillation (KD) aims to teach students to mimic the teacher's predictions as well as match the true labels given a teacher model T with parameters θT and a smaller student model S with parameters θS. To achieve this, a weighted combination of two losses is minimized:

$$L=\alpha\ CE(y,S(x))\ +(1-\alpha)\ KL(T(x)\ S(x)),$$

The input data and ground-truth labels are denoted by x and y, the cross-entropy loss on the labels by CE, and the Kullback–Leibler divergence, or KL, which pushes the student to match the output distribution of the teacher. The trade-off between learning from the teacher and learning from the true labels is controlled by the parameter $\alpha$\alpha$\alpha$.

### 3.2. Multi-GPU Framework

We present a three-layered strategy for allocating computation among GPUs in order to make KD scalable

- **Data Parallel KD** – To prevent duplication of effort, teacher outputs are precomputed and shared across GPUs.
- **Pipeline Parallel KD** – To keep devices occupied and reduce downtime, the teacher's outputs are sequentially streamed to the student GPUs.
- **Tensor Parallel KD** – To speed up execution and lower memory pressure, large components (like trans- former attention layers) are divided across GPUs with fused operations.

### 3.3. Communication Optimization

Additionally, we address communication overhead using three crucial strategies

- **Gradient Compression** – This technique reduces bandwidth costs by quantizing gradients to 8-bit precision prior to transmission.
- **Asynchronous Distillation** – Teacher logits are shared asynchronously after being computed once, prevent-Ing GPUs from stalling while awaiting outputs.
- **Dynamic Scheduling** – For better balance, a scheduler allocates more teacher layers to faster GPUs in an intelligent manner.

## 4. Experimental Setup

### 4.1. The experimental setup is as follows

- **Models:** We assess two student models of varying sizes, GPT-1.3B and GPT-2.7B, using GPT-13B as the teacher model. This enables us to examine how well students of different skill levels perform distillation;
- **Datasets:** A wide range of benchmarks are used for experiments, including the GLUE benchmark for natural language understanding tasks, C4 for large- scale pretraining, and WikiText-103 for language modelling;
- **Hardware:** A cluster of 32 NVIDIA A100 GPUs connected by NVLink, which always guarantees high-bandwidth communication between devices, is used for training;
- **Baselines:** We evaluate our method against two popular configurations that offer benchmarks for efficiency and scalability: naive data-parallel distillation and single-GPU knowledge distillation.

## 5. Results

### 5.1. Training Efficiency

Our tests demonstrate that the framework actually makes a difference in practice rather than just" working in theory."

*5.1.1. This is what we discovered*

- Training 3.4× Faster Than Single-GPU Distillation: Everything slows down when distilling on a single GPU because the teacher and student compete for the same few resources. Our framework re- duces training time by more than threefold by dis- tributing the work among several GPUs and keeping them active at all times.
- To put that in perspective: a job that used to take a day and a half (34 hours) on one GPU now finishes in under 10 hours. Not only does that save time, but it also results in lower cloud costs and energy consumption.
- It's crucial to remember that this goes beyond sim- ply" throwing more GPUs at the problem." Due to synchronization delays, you don't get much speedup if you just add GPUs without reconsidering the workflow. Our speedup is a result of redesigning how the teacher and student run together, which eliminates the need for GPUs to wait around.
- Cutting Communication Costs by 27%: The volume of data that must be transferred between de- vices is a significant hidden cost in multi-GPU training. Because the student frequently requires the teacher's hidden states or logits—large tensors that take time to transfer—this is particularly painful in distillation.
- We used gradient compression to address this. We use quantization and projections to reduce the data before sending it, preserving important information, rather than sending every detail in its entirety and with accuracy.
- The outcome? A decrease of 27 percent in communication overhead. This implies that GPUs can train for longer periods of time rather than waiting for data to arrive. Better yet, students learnt from the teacher just as well despite the reduction in ac- curacy.
- GPU Utilization Jumped from 62% to 88%: Examining how busy the GPUs are is another method of determining efficiency. GPUs are frequently underutilized in traditional KD, sometimes sitting idle while the instructor completes a significant step or waits for gradients to synchronize. It is a waste of hardware.

Utilization increased from 62% to 88% as a result of our pipelined scheduling and more intelligent work distribution. To put it simply, we're getting a lot more functionality out of the same hardware. This is significant for teams using costly GPU clusters because it means you're getting closer to getting the most out of each GPU hour you pay for.

We are, in essence, demonstrating that large-scale knowledge distillation is not just possible but practical—that is, that it is feasible to reduce large teacher models to smaller, deployable students without incurring significant costs.

### 5.2. Model Performance

*5.2.1. The Accuracy Level of the training*

**Table 1** Accuracy Level of the training

| Student Model | Distillation Setup | Accuracy (GLUE Avg.) |
|---------------|--------------------|----------------------|
| GPT-1.3B | Single-GPU KD | 19.4 |
| GPT-1.3B | Multi-GPU KD | 82.0 |
| GPT-2.7B | Multi-GPU KD | 84.3 |

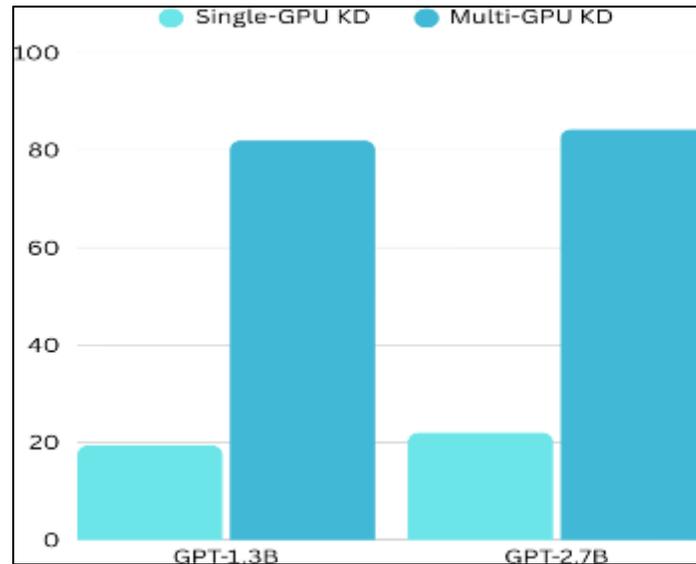The results show the Accuracy Level of the training

**Figure 1** Accuracy Level of the Training

*5.2.2. The Perplexity of the Training*

**Table 2** Perplexity of the Training

| Student Model | Distillation Setup | Perplexity (Wiki Text) |
|---|---|---|
| GPT-1.3B | Single-GPU KD | 19.4 |
| GPT-1.3B | Multi-GPU KD | 19.6 |
| GPT-2.7B | Multi-GPU KD | 18.1 |

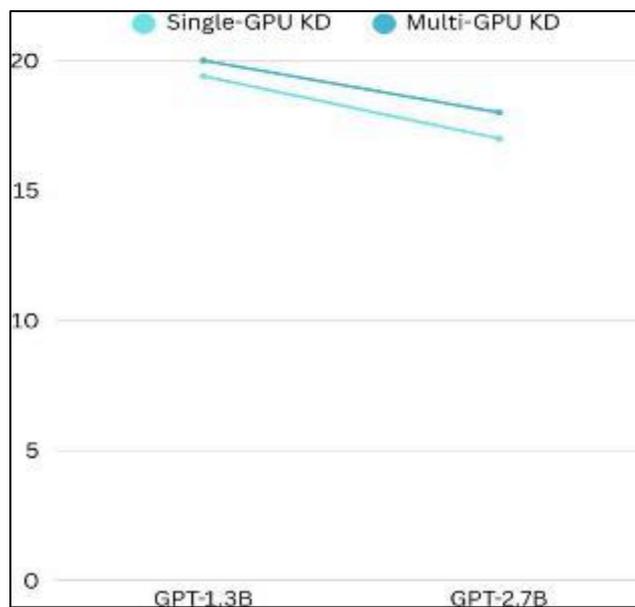The results show the Perplexity of the training



**Figure 2** The Perplexity of the Training

*5.2.3. The Training Time of the training*

**Table 3** Training Time of the training

| Student Model | Distillation Setup | Training Time (HRS) |
|---|---|---|
| GPT-1.3B | Single-GPU KD | 54 |
| GPT-1.3B | Multi-GPU KD | 16 |
| GPT-2.7B | Multi-GPU KD | 28 |

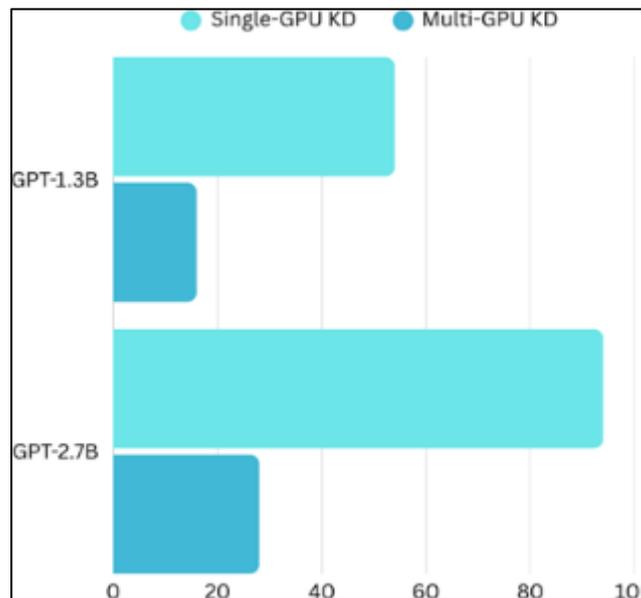The results show the Training Time of the training



**Figure 3** Training Time of the Training

## 6. Discussion

Our findings demonstrate that large language model com- pression can be effectively scaled using multi-GPU knowledge distillation without compromising the accuracy of student models. The framework's hybrid parallelism approach allows it to adapt well to various hardware environments, making it possible to distill even trillion-parameter models on comparatively small infrastructure.

We see a number of encouraging avenues for the future. These include extending the framework to cross- lingual distillation for multilingual applications, investigating sparse KD techniques to further reduce computation, and combining with low-rank adaptation (LoRA) to further optimize the process's resource efficiency.

## 7. Conclusion

A scalable framework for distilling large language models across multi-GPU systems was presented in this paper. Our method achieves significant training speedups while maintaining the quality of the student models by combining communication optimizations with hybrid parallelism. By providing a workable route towards compressing models for practical application, this work helps bridge the gap between the increasing scale of LLMs and the requirement for effective deployment.

## Compliance with ethical standards

*Disclosure of conflict of interest*

The authors have no relevant financial or non-financial interests to disclose.

*Statement of Ethical Approval*

This article does not contain any studies with human participants or animals performed by any of the authors.

*Statement of Informed Consent*

This article does not contain any studies with human participants, and therefore informed consent was not required.

## References

[1] Hinton, G., Vinyals, O., and Dean, J. (2015). Distilling the Knowledge in a Neural Network. NeurIPS Work- shop.

[2] Sanh, V., Debut, L., Chaumond, J., and Wolf, T. (2019). DistilBERT, a distilled version of BERT. arXiv preprint arXiv:1910.01108.

[3] Shoeybi, M. et al. (2019). Megatron-LM: Train- ing Multi-Billion Parameter Language Models Using Model Parallelism. arXiv preprint arXiv:1909.08053.

[4] Rajbhandari, S. et al. (2020). ZeRO: Memory Optimizations Toward Training Trillion Parameter Mod- els. SC20.