



How Can We Train Deep Learning Models Across Clouds and Continents? An Experimental Study

Alexander Erben
Technical University of Munich
alex.erben@tum.de

Ruben Mayer
University of Bayreuth
ruben.mayer@uni-bayreuth.de

Hans-Arno Jacobsen
University of Toronto
jacobsen@eecg.toronto.edu

ABSTRACT

This paper aims to answer the question: Can deep learning models be cost-efficiently trained on a global market of spot VMs spanning different data centers and cloud providers? To provide guidance, we extensively evaluate the cost and throughput implications of training in different zones, continents, and clouds for representative CV, NLP and ASR models. To expand the current training options further, we compare the scalability potential for hybrid-cloud scenarios by adding cloud resources to on-premise hardware to improve training throughput. Finally, we show how leveraging spot instance pricing enables a new cost-efficient way to train models with multiple cheap VMs, trumping both more centralized and powerful hardware and even on-demand cloud offerings at competitive prices.

PVLDB Reference Format:

Alexander Erben,
Ruben Mayer, and Hans-Arno Jacobsen. . PVLDB, 17(6): 1214 - 1226, 2024.
doi:10.14778/3648160.3648165

PVLDB Artifact Availability:

The source code, data, and/or other artifacts have been made available at <https://github.com/circuit/hivemind-multi-cloud>.

1 INTRODUCTION

Deciding whether to invest in on-premise hardware or move to the cloud for deep learning (DL) is not easy. Wanting to scale existing infrastructure means paying upfront, as combining cloud and on-premise is not an option with popular DL frameworks due to needing a dedicated high-bandwidth interconnect. To enabled model- and data-parallelism, current state-of-the-art accelerators have bandwidths of 900 GB/s for intra-node [19] and 25 Gb/s for inter-node setups [1, 26]. Due to the initial investment of the cloud providers in the accelerators, they naturally want to reap profit by maximizing resource utilization. Therefore, it is common to have "spot" pricing, which offers the VMs at a strongly reduced rate, typically at a 40-90% discount (Section 1), but with the drawback that the VM can be terminated at any time if another customer is willing to pay the on-demand price [33]. Unfortunately, popular DL frameworks have not been developed with failure semantics in mind and cannot adequately deal with peers that fail [12]. While services like Amazon SageMaker [14] and projects like Skypilot [43] offer automatic job migration in case of VM termination, they are limited to single-node training due to the bandwidth requirements between accelerators.

This work is licensed under the Creative Commons BY-NC-ND 4.0 International License. Visit <https://creativecommons.org/licenses/by-nc-nd/4.0/> to view a copy of this license. For any use beyond those covered by this license, obtain permission by emailing info@vldb.org. Copyright is held by the owner/author(s). Publication rights licensed to the VLDB Endowment.
Proceedings of the VLDB Endowment, Vol. 17, No. 6 ISSN 2150-8097.
doi:10.14778/3648160.3648165

Table 1: Average us-west cloud pricing in April '23.

| Cloud Type | GC | AWS | Azure |
|------------------------------|------------|------------|------------|
| T4 Spot | 0.180 \$/h | 0.395 \$/h | 0.134 \$/h |
| T4 On-Demand | 0.572 \$/h | 0.802 \$/h | 0.489 \$/h |
| Traffic (inter-zone) | 0.01 \$/GB | 0.01 \$/GB | 0.00 \$/GB |
| Traffic (inter-region) US | 0.01 \$/GB | 0.01 \$/GB | 0.02 \$/GB |
| Traffic (inter-region) EU | 0.02 \$/GB | 0.01 \$/GB | 0.02 \$/GB |
| Traffic (inter-region) ASIA | 0.05 \$/GB | 0.01 \$/GB | 0.08 \$/GB |
| Traffic (inter-region) OCE | 0.08 \$/GB | 0.01 \$/GB | 0.08 \$/GB |
| Traffic ANY-OCE | 0.15 \$/GB | 0.02 \$/GB | 0.08 \$/GB |
| Traffic (between continents) | 0.08 \$/GB | 0.02 \$/GB | 0.02 \$/GB |

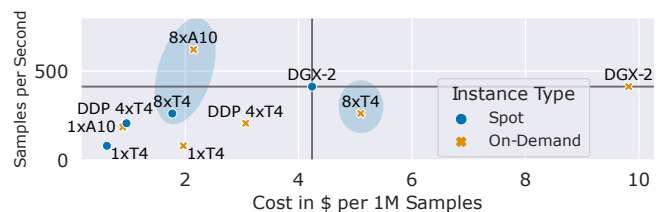


Figure 1: Cost to throughput tradeoff for ConvNextLarge at different instance types. Our training setups (circled) are cheaper (8xT4) and faster (8xA10) than centralized offerings (DGX-2).

But what if we could use spot pricing for long-running, distributed jobs and reduce bandwidth requirements to leverage multiple low-cost GPUs? This could be possible through a framework for collaborative DL training, Hivemind [39], which inherently deals with peers that can stop running at any time. While there is research on how Hivemind can be used for training on spot VMs [17, 37, 38], it does not compare the cost-throughput tradeoff for different cloud offerings or perform ablation studies on geographic distribution or model sizes.

To motivate this new possibility, we trained the ConvNextLarge model [29] on the Imagenet1K dataset [15] on different Google Cloud hardware (T4's and DGX-2), and on the very competitively priced A10 from LambdaLabs (see Section 6 for the full experimental description). Figure 1 shows the training throughput and the costs per 1 million processed samples for each setup. The single node (1xT4, 1xA10, DGX-2) experiments show the current state-of-the-art cost-throughput ratio for training on GC and LambdaLabs. The DGX-2 node is the fastest, with a throughput of 413 SPS, but it also costs \$6.30/h (\$4.24/1M samples), shown by the horizontal and vertical lines. The single-accelerator experiments (1xT4, 1xA10) have a better cost-throughput ratio (\$0.62/1M samples and \$0.9/1M samples), but have a much lower throughput of 80 and 185 SPS, respectively. However, when using our approach of distributing the training between multiple GPUs with Hivemind (circled), we make training possible that is both faster (8xA10, 621 SPS, \$2.15/1M samples) and cheaper

(8xT4, 262 SPS, \$1.77/1M samples) than using the DGX-2. Every cloud provider deals differently with how they price spot instances and network traffic (cf. Section 1) and has varying interruption rates for different accelerators [23]. Being able to choose the best option was not possible before, and having the option to combine older, more available GPUs is a net benefit for both consumers and cloud providers alike.

We aim to develop guidelines and help practitioners assess under which conditions they can cost-efficiently speed up their training tasks with spot instances. To be able to do this, they need a precise definition of the model size at which geo-distributed spot training becomes viable, what hardware can be used for it, and what the minimum bandwidth and latency are. We close this research gap by performing a comprehensive analysis of multiple DL tasks from CV and NLP, breaking down how time is spent in each epoch, and comparing them to non-distributed runs to quantify the advantages and disadvantages of distributed spot training. We determine which models scale with additional spot instances and which cannot be scaled without running into a communication bottleneck or resource inefficiencies. To quantify total training cost, we assess cost-effectiveness and evaluate a hybrid or multi-cloud approach with popular cloud providers through training on up to four continents. For comparison of the models’ scalability and to show which of them can be trained in a distributed fashion, we introduce the *granularity metric*, the ratio of calculation to communication time, and show how it can be used for predicting performance with different hardware setups. Finally, we summarize our lessons on how to design geo-distributed spot training and what to watch out for when evaluating the feasibility of such a training regime. Our contributions are:

- (1) **We analyze the impact of multi-cloud training with spot and on-demand instances from Google Cloud (GC), Microsoft Azure, Amazon Web Services (AWS), and LambdaLabs on cost-efficiency.** While we find performance penalties due to remote versus on-premise compute resources, the throughput still scales with increased computing power. By leveraging multiple spot instances with one T4 GPU each, we can be more cost-efficient than a DGX-2 node or the very competitively priced A10 offerings from LambdaLabs.
- (2) **We investigate the suitability of geo-distributed training for various CV and NLP models and hardware configurations on up to four continents.** Not surprisingly, the more parallelizable and the larger the task, the better the performance. Moreover, we verify the scalability claims of the related work and define additional constraints, such as the minimum granularity for effective training. This enables, for the first time, distributed training of smaller million-parameter models (12M-560M) over <1 Gb/s bandwidth and >150ms latency networks.
- (3) **We evaluate two different hybrid-cloud experimental setups with consumer- and server-grade on-premise hardware** and try to improve the throughput with a bandwidth of, at worst, 50 Mb/s to the cloud resources. While we show that it is possible to improve throughput even at these constraints, local cloud offerings are better suited for models that show limited suitability for distributed training.

- (4) We summarize our findings of training in a geo-distributed, multi-cloud environment. **We propose the granularity metric to compare model suitability for distributed spot training** and estimate training performance with additional spot VMs. This provides guidance on the trade-off between performance and cost when using geo-distributed spot instances. To apply our findings, we perform a case-study on a state-of-the-art model from the ASR domain and achieve speedups on low-end hardware.

2 DEEP LEARNING ON SPOT INSTANCES

In this section, we describe how the Hivemind framework works and how it can enable distributed spot training.

2.1 Hivemind

Hivemind [39] is a PyTorch-based [32] framework developed initially to enable collaborative DL training where participants could donate their heterogeneous hardware to train a single model together in a data-parallel fashion. Its main difference to other state-of-the-art distributed training frameworks, such as PyTorch DDP [26] and DeepSpeed [35], is that it runs in a decentralized fashion and can handle peers that drop out at any stage of the training. It does so with two features: a distributed hash table [31] (DHT) which spans over all participating peers for metadata storage, such as training progress and peer health, and a gradient averaging algorithm that is designed to reduce the impact of lost gradients. A key difference to other distributed training frameworks is the definition of a *hivemind epoch*, which is the number of samples that must be aggregated before an averaging step is performed. This sample count is called the *target batch size* (TBS), which corresponds to the minibatch size in standard DL training. The DHT is used for coordination, and shortly before the TBS is predicted to be reached, the peers start to form the initial groups for averaging. The time allocated for group forming is called *matchmaking time* and typically runs asynchronously to the training (cf. Section 3). The individual peer gradients are accumulated locally and sent to the other peers via an adaptive all-reduce algorithm (MoshpitSGD [38]). The next hivemind epoch starts after each peer applies the accumulated gradients to the local model. The advantage of Hivemind for geo-distributed training comes from cumulating different techniques, such as Delayed Parameter Updates [36], big-batch training [44] and aggressive communication quantization [16]. All of these combined reduce time and frequency of the communication rounds, which in turn makes training on heterogeneous devices and low-bandwidth networks possible.

2.2 Distributed Spot Training

In this paper, we focus only on models that fit into the memory of a single GPU, as we are interested in utilizing data parallelism on cheaper and more readily available hardware. However, our insights are applicable to larger models with techniques such as ZeRO offloading [36], more aggressive quantization [41] and even model parallelism [37]. The current options for data parallelism are either using multiple GPUs on the same node (e.g., a DGX system with eight GPUs) or having multiple nodes with a GPU each in the same high-bandwidth network (>25 Gb/s) to minimize communication time. The latter does not work on cheap but interruptible instances,

while the former has some use in the form of Amazon Sagemaker but is limited to a single node and is typically very pricey (spot pricing for DGX-2 is \$6.30/h versus 8xT4 at \$0.72/h on GC). However, using Hivemind, a new training scenario becomes feasible: Distributed training in a decentralized fashion on interruptable VMs with bandwidths of <1 Gb/s. Since spot instance prices change hourly depending on the time of day and zone availability [23], and can vary widely between cloud providers (cf. Section 1), training between continents and in multiple clouds could potentially be more cost-effective than using a single, more computationally powerful node at spot prices.

With the newly added training setups from Figure 1 (circled), it was not previously possible to choose the best option, and having the option to combine older, more available GPUs is a net benefit for both consumers as well as cloud providers. Our paper shows that it is possible to train on multiple clouds across multiple continents and provides guidelines on how to accomplish this cost-efficiently.

3 MODEL SUITABILITY

Selecting suitable models with a big enough parallel workload is essential to ensure successful distributed spot training. To cover a wide range of established models, we drew from MLCommons' comprehensive DL training benchmark [30]. We used models from the CV and NLP domains and gradually increased their size and TBS to increase the parallel compute amount. As discussed in Section 2, the TBS may be exclusively responsible for the success of distributed training and was chosen to cover both medium and large batches (8K, 16K and 32K). These minibatch sizes start to become more common due to the LAMB optimizer [44], which works well enough for both smaller (512) and huge batches (64K) and should be representative of state-of-the-art workloads. For a representative experimental study with a minibatch size of 256 on the automatic speech recognition model (Whisper [34]), please refer to Section 11. All experiments were run with FP16 precision, as the target T4 GPUs have a considerable improvement in FLOPs compared to FP32 (8:1).

For **CV**, we take five models from the extended ResNet family, starting with the smallest one, ResNet18 [21] (RN18), ResNet50 (RN50), ResNet152 (RN152), WideResNet101_2 [46] (WRN101) and ConvNextLarge [29] (CONV), which is almost 20 times larger than RN18. The parameter count is 11.7M, 25.6M, 60.2M, 126.9M, and 197.8M, respectively. These models were popularized due to their ability to help with the vanishing gradient problem by using residual connections between layers. Currently, they are not only used for classification, but can serve as an embedding of images by removing the classification head [18, 40]. For the dataset, we use Imagenet1K [15] and train the classification task, which tries to assign one of 1000 classes to each image.

For **NLP**, we selected three models from the BERT family: RoBERTaBase [28] (RBase), -Large (RLrg), and -XLM [13] (RXLM). The parameter count is 124.7M, 355.4M, and 560.1M, respectively. We used the same configuration as the original models and trained them on masked language modeling, a common pre-training task. RoBERTa models were a replication study of BERT but with a focus on better hyperparameter tuning, leading to state-of-the-art results and proposed using much higher minibatch sizes than previously common. The text dataset is March '22 Wikipedia [20].

When we run our experiments in a multi-cloud environment on spot instances, we cannot plug in proprietary cloud storage or wait

for the dataset to download, as the instances can be terminated anytime. To simulate a real-world deployment with a non-public dataset, we chose an independent S3 storage provider, Backblaze (B2) [4]. Backblaze has replicated data centers that can better serve requests from anywhere worldwide, guaranteeing a reasonable ingress rate from every continent. Additionally, the cost is very manageable at \$0.01/GB rate for egress and \$0.005/GB/month for storage. A detailed analysis of the costs incurred for the experiments can be found in Section 5. We access the datasets on-demand via shards in the tar format with the WebDataset library [10]. We chose WebDataset due to its features like automatic local caching, streaming decompression, streaming preprocessing, and having an easy to work with archive format that allows representing the data in its original format. Finally, for the Hivemind parameterization, we enabled delayed parameter averaging (DPU) [36] to enable simultaneous gradient communication and computation at the expense of a round of staleness. We selected FP16 compression for peer-to-peer communication.

Experimental design. First, we must verify that our models are suitable for cloud training. For this purpose, we evaluate them on the powerful Ampere GPUs first - if they scale there without facing a communication bottleneck, they should also scale on the slower T4, which is common at GC, AWS, and Azure. We use the LambdaLabs [8] for these experiments, which gives us on-demand A10 GPUs for just \$0.60/hour, but currently offer their services only in the US West region. All experiments are performed on the 515.65.01 driver, CUDA 11.6, and PyTorch 1.13.1. We profiled a network bandwidth of 3.3 Gb/s and a latency of 0.3 ms between the Lambda VMs.

To establish a fair baseline, we train all models from ?? on a single GPU that achieves large minibatch sizes through gradient accumulation. Processes logs system metrics every second and evaluates the training performance whenever a batch is processed. Finally, all multi-GPU experiments are monitored with a training monitor that scrapes the DHT every second to log the peer state and training progress synchronously.

(1) Hivemind penalty. Using Hivemind as middleware to share gradients and keep a fully decentralized architecture running harms performance compared to single-node training. We can compare the effects of Hivemind training by looking at three metrics: *baseline*, the single GPU throughput, *hivemind local*, normalized GPU throughput without the averaging step, and *hivemind global*, the actual normalized GPU throughput. When comparing the baseline and local speed in Figure 2 for a setup with two GPUs, running Hivemind reaches at best 78% (RN152) and at worst 48% (CONV) of the baseline performance. Unsurprisingly, the larger the model size, the worse the penalty gets due to the increased size of the accumulated gradients (GAC) over each step. However, the baseline also applies gradient accumulation to reach the target minibatch size without the performance drop. After isolating the respective function calls, there seems to be a slight inefficiency in how GAC is implemented in Hivemind versus the native PyTorch call. We are working with the maintainers to fix this issue [7]. On the other hand, the disadvantage of synchronization is minimal under the perfect conditions of a good interconnect. The global speed in Figures 2a and 2b only degrades at best to 97% (CONV) to at worst to 87% (RBase) compared to the local throughput, meaning that the communication under these conditions only accounts for a fraction of the total training time. This degradation is inversely correlated to the model size due to larger

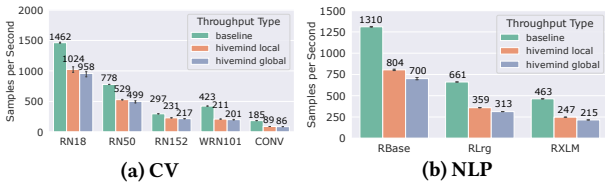


Figure 2: Hivemind penalty on normalized throughputs.

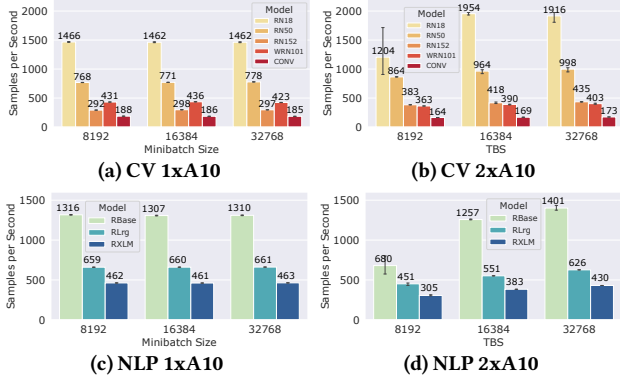


Figure 3: Throughput comparison between single GPU base-lines and the Hivemind runs with two GPUs.

models training quadratically longer per parameter, but the communication only increases linearly [37]. While an implementation issue currently affects performance, and the worst total performance drop is at 47% (CONV baseline vs. global), scaling is still possible with a ratio of roughly 2:1 of GPUs to throughput. We further refine this ratio in the following section by comparing which models are most suitable to be trained in a distributed environment.

(2) **Less suitable models for distributed spot training.** While training billion-parameter NLP models scale well due to the "square-cube" law, the minimum model size is not yet fully defined [37]. The reason is that many factors play a role in whether a model is suited for geo-distributed training. On the one hand, a small model results in small gradients exchanged between peers, so the averaging step is fast. On the other hand, a small model will also reach the TBS faster than larger models, which may lead to a low speedup if the calculation time is disproportionately lower than the communication time.

We found the granularity metric [22], typically used in high-performance computing, practical to attach a comparable value to each setup to quantify the ratio of the calculation to communication time. The higher the granularity, the more parallelizable the task, as more calculation can be distributed between peers, ensuring a good speedup. It is important to note that this metric depends on the model and the hardware being used. The communication time is affected by the parameter count, and the calculation time is affected by the layer type of the parameters (including feedforward, convolution, and transformer). Therefore, the calculation time can decrease with improved hardware, which we evaluate in Section 6. Another parameter that affects the calculation time is the TBS that all peers work to accumulate. There is a practical limit to the TBS where a model is still trainable, which is currently at 64K with the LAMB optimizer [44]. This limits the possibility of improving the speedup of small models by increasing the batch size, meaning that at some point, the speed will be limited by the communication time. It is

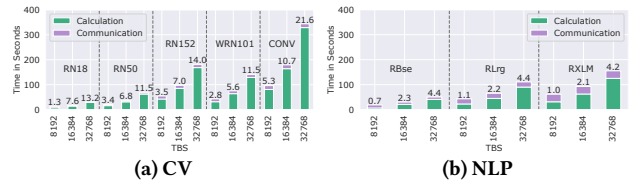


Figure 4: TBS vs. total training time on 2xA10s. Granularity is shown above each bar. Dotted lines separate different models.

important to remember that just increasing the TBS to create more calculation time can have a grave effect on training performance if the optimizer is not adequately selected and configured.

Our experimental results in Figure 3 show the practical implications of this observation. For the 2xGPU experiments in Figures 3b and 3d, we can see the effect of a TBS increase which improves the total throughput. Doubling the TBS equals cutting down the per-sample communication cost by two, which leads to the slight increase in performance visible in both CV and NLP experiments. However, the smallest models, RN18 and RBase, fluctuate significantly at a TBS of 8K due to a minimum matchmaking time of 5 seconds. Whenever all peers accumulate the TBS in less than 5 seconds, the asynchronous thread that matches the peers in groups to perform the all-reduce may still need to finish. This results in an unstable averaging time, which limits the scalability of small models with a small TBS.

To illustrate how the TBS and model size affect the individual timings, we visualize the total training time split up into the calculation and communication time in Figure 4. CV models are generally computationally more expensive and have a higher granularity than NLP models, which have slightly longer averaging rounds due to the much larger model sizes. When comparing the models at the same TBS (e.g., 32K), there is an inconclusive relation between runtime and parameter count. Some models increase their runtime with parameter count w.r.t. smaller models (RN50 to RN152, RBase to RLrg), while others decrease their runtime (RN152 to WRN101, RLrg to RXLM). This performance is due to not all layer parameters contributing similarly to computational complexity. Depending on the specific architecture, even models with more parameters can be faster to train due to a more efficient architecture, such as the WRN101 [46].

The communication time between different TBS sizes stays the same, barring the two matchmaking time exceptions (RN18, RBase), as the gradients are accumulated before being sent. For all other models, doubling the TBS leads to exactly double the amount of work and doubles the granularity. With a TBS of 32K, all models have a granularity of at least 4.2 (RXLM) and at most 21.6 (CONV), which show strong scaling potential. Therefore, we decided to use a TBS of 32K for all following experiments to ensure that the setup scales before introducing bandwidth and computational limitations.

Summarizing, whether a model is scalable without network bandwidth limitations depends on the minimum time to reach the TBS and on the granularity. Tuning the TBS is possible to a certain extent but depends on the specific training task and optimizer configuration.

(3) **Per-GPU speedup decreases with low granularity.** To evaluate the scalability with additional hardware, we profile all models on 2,3,4, and 8 GPUs with a TBS of 32K. Figure 5 shows the throughput for all models in the different hardware scenarios. Generally, all models scale well regardless of size, with the best speedup of 4.37x (RN152) and the lowest at 2.29x (RXLM) with 8 GPUs. There is a

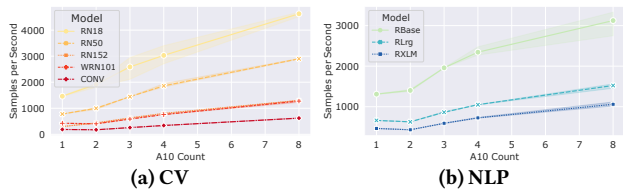


Figure 5: Throughput comparison from 1 to 8 A10 GPUs.

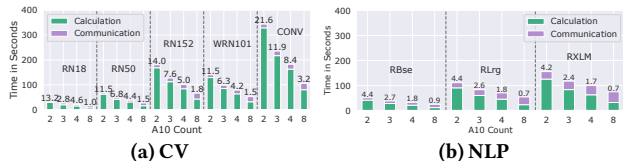


Figure 6: Multi-GPU scalability at 32K TBS. Granularity is shown above each bar. Dotted lines separate different models.

Table 2: Geo-distributed experiments on GC with T4 VMs.

| Exp. Name | Resources | Total |
|-----------------|--|-------------|
| A-{1,2,3,4,6,8} | {1, 2, 3, 4, 6, 8}xUS | 1,2,3,4,6,8 |
| B-{2,4,6,8} | {1, 2, 3, 4}xUS + {1, 2, 3, 4}xEU | 2,4,6,8 |
| C-{3,6} | {1, 2}xUS + {1, 2}xEU + {1, 2}xASIA | 3,6 |
| C-{4,8} | {1, 2}xUS + {1, 2}xEU + {1, 2}xASIA + {1, 2}xAUS | 4,8 |

visible trend in the per-GPU contribution to the speedup ($\frac{\text{speedup}}{\#\text{GPUs}}$). The more GPUs we add, the lower the contribution, e.g., RN18 goes from 0.7 to 0.4 with two to eight GPUs, respectively. This decrease is likely to continue due to a granularity of 1.0 at 8 GPUs (Figure 6a), as doubling the GPUs would, at best, increase the throughput by 33% by halving the calculation time. However, the more computationally expensive the models are, the slower the per-GPU contribution falls off and the larger the granularity is (RN152, CONV). This does not hold true for our NLP models (Figure 6b); while they have increasingly more model parameters, the only difference between the two biggest models, RLrg and RXLM, is the vocabulary size increase of 50K to 250K. Due to how embedding layers are lookups, the forward pass is not affected by the increased embedding size, but the backward pass is. This results in a smaller increase of the calculation time while communication increases linearly with the number of parameters.

Additionally, we see the drop in throughput when comparing the single GPU and dual GPU experiments for most larger models (Figure 5), which stems from observation (1) of the Hivemind penalty.

We also observe that with each subsequent doubling of GPUs, the calculation time is halved, while the communication increases sub-linearly due to the more efficient group-based all-reduce of MoshpitSGD [38]. For example, the averaging step for the RXLM on 2xA10 takes 5 seconds per GPU (10s total), while the 8xA10 averaging step takes 1.8 seconds per GPU (14.4s total).

In summary, all models show a speedup but have a decreasing per-GPU contribution due to smaller granularity with more GPUs. Therefore, the larger the model and TBS, the greater the scaling potential. High granularity is a good indicator of scalability, and since the communication time only increases linearly with additional peers (cf. Section 2.1), knowing the initial calculation time is a good indicator of future throughput. Under the optimal conditions of good compute performance and an interconnect with relatively high bandwidth, scaling was not a problem. But what happens under less favorable conditions in geo-distributed settings?

4 GEO-DISTRIBUTED PERFORMANCE

As spot prices for the same hardware differ depending on the region, zone, and time of day [23], it might be a good idea to use VMs across different data centers. However, is the connectivity between regions and continents good enough to enable distributed deep learning? To explore this question, we decided to conduct three types of experiments (Table 2):

- (A) **Intra-zone** Can we scale if the VMs are co-located in the same zone (us-central-1)?
- (B) **Transatlantic** Can we scale when we combine VMs from two regions (US and EU), and what happens when the compute is unevenly distributed across regions?
- (C) **Intercontinental** Can we scale if we combine VMs from four continents (US, EU, ASIA, AUS)?

Experimental design. Based on the insights from Section 3, we decided to use the largest models (CONV, RXLM) for all further cloud experiments in Sections 4 to 6, with the TBS of 32K as a baseline with good scaling properties. We abbreviate them with their respective domain names (CV, NLP). We used Google Cloud [5] for all experiments in this section, as they were the first to give us access to all necessary zones. The default networking solution in GC is the "Premium Tier", which tries to use a Google-owned network instead of the public internet. We measured the throughput and latency between all zones via iperf and ping and report the average of 5 consecutive runs in Table 3. Unsurprisingly, the diagonal shows that the local connectivity between zones runs at almost 7 Gb/s with a latency of 0.7ms, probably due to the hypervisors being in the same data center. While the up- and download were perfectly symmetrical in all setups, the throughput dropped to <210 Mb/s for all non-local connections. The US-based data center is located in Iowa and is best connected with at least 120 Mb/s to the remaining regions, namely Belgium in the EU (6,911km), Taiwan in ASIA (11,853km), and Sydney in Australia (AUS, 14,555km), presumably due to the physical distance. The lowest bandwidth and highest latency connections are between the EU region and ASIA and AUS, reaching around 80 Mb/s and 270ms. We decided to use the n1-standard-8 template with eight cores, 30 GB memory, and a T4 GPU, as the smaller image with 15 GB was insufficient to meet the memory requirements for gradient application on the CPU with the biggest models. The experiment naming in this section is prefixed with the type of location (A), (B) or (C) and the number of VMs, e.g., A-4 is the intra-zone experiment with 4 VMs. The full experimental description is specified in Table 2.

Table 3: Throughput and latency between GC zones.

| (a) Single stream TCP throughput in Gb/s. | | | | | (b) ICMP latency in ms. | | | | |
|---|------|------|------|------|-------------------------|--------|--------|--------|--------|
| From | To | | | | From | To | | | |
| | US | EU | ASIA | AUS | | US | EU | ASIA | AUS |
| US | 6.90 | 0.21 | 0.13 | 0.12 | US | 0.66 | 103.11 | 157.09 | 176.19 |
| EU | 0.21 | 6.81 | 0.08 | 0.07 | EU | 103.14 | 0.65 | 253.10 | 271.98 |
| ASIA | 0.13 | 0.08 | 6.79 | 0.16 | ASIA | 157.08 | 253.09 | 0.72 | 131.45 |
| AUS | 0.12 | 0.07 | 0.16 | 6.84 | AUS | 175.98 | 272.08 | 131.42 | 0.64 |

(A) **Intra-zone scalability.** Figure 7 shows the result of the intra-zone experiments, which we used as a baseline to compare geo-distributed deployments to. As the scalability of the CV and NLP models was already shown with much better hardware and slightly worse network connectivity (cf. Section 3), the scalability with the T4 GPUs is not too surprising. We do not see an improvement in throughput for two GPUs for either model due to the Hivemind penalty discussed

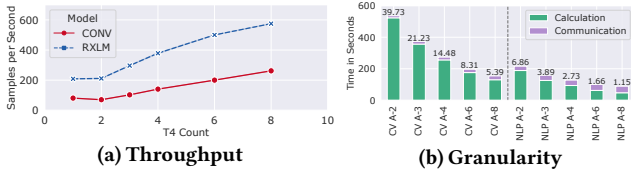


Figure 7: (A) Intra-zone performance for CV and NLP.

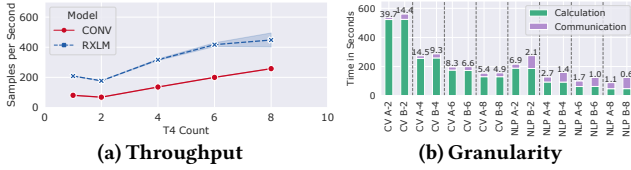


Figure 8: (B) Transatlantic performance for CV and NLP.

in Section 3. However, starting with three GPUs, we see an increase in throughput with a maximum speedup of up to 3.2x CV and 2.75x for NLP at eight GPUs. CV’s per-GPU speedup ($\frac{\text{speedup}}{\#\text{GPUs}}$) is almost linear (0.43, 0.42, 0.43, 0.41, 0.41), while NLP starts dropping off faster (0.51, 0.47, 0.45, 0.40, 0.34) for 2, 3, 4, 6 and 8 GPUs, respectively. The reason for this is the NLP granularity of 1.15 with 8 GPUs indicating an almost equal part in communication and calculation (Figure 7b) due to the much longer averaging round related to the model size (198M vs. 560M parameters). The peak network bandwidth utilization between peers was at most a symmetric 1.1 Gb/s while averaging and 33 Mb/s ingress while training due to data loading. This means that the network bandwidth of 7 Gb/s was not a limiting factor.

(B) Transatlantic scalability. We scale when computing hardware is local. However, what happens when there is cheap capacity in another region? In this case, we study the throughput of experiments with resources in the us-west and eu-central regions (B-2,4,6,8).

The B-2 experiment has one VM in the US and one in the EU, achieving a virtually identical throughput of 68.4 (US-EU) versus 70.1 (US) at CV (Figure 8a). Our maximum peak egress rate of 250 Mb/s does not affect the CV experiments, while the US experiments peaked at 1.1 Gb/s. The reduction in bandwidth penalizes NLP harder, where we are 16% slower with 177.3 SPS (US-EU) compared to the intra-zone experiment with 211.4 SPS (US). The resulting increased communication can be easily seen in the granularity analysis in Figure 8b (NLP A-2,4,6,8 vs. B-2,4,6,8). As only communication time increases in the NLP (B) experiments compared to (A), a granularity of $\gg 1$ indicates good scalability: Adding two more GPUs to the B-6 experiment with a granularity of 1.03 results in a throughput increase of 15% (B-8) relative to the baseline. Meanwhile, adding two more GPUs to the B-2 experiment with a granularity of 2.21 results in a throughput increase of 77% (B-4) relative to the baseline.

In the B-4 experiment, we look at what happens when we increase the number of VMs to four, with two in the US and two in the EU. Nothing surprising happens with CV, as the workload continues to be mostly computation, with a throughput of 135.8 (B-4), only 3% slower than the intra-zone experiment with 140.4 SPS (A-4). However, at NLP, things get more interesting as we now have more overall communication with four peers, but they can average locally first and only later transmit across the Atlantic. However, compared to their A-counterparts, we do not see a difference in relative scalability with either B-4, B-6, or B-8. This means that training across regions

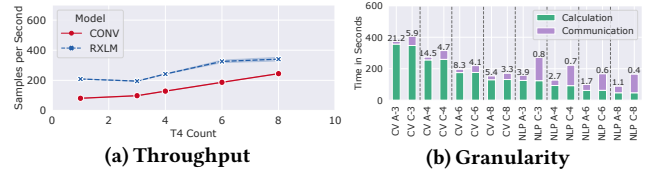


Figure 9: (C) Intercontinental performance for CV and NLP.

(B) is slower, but the contribution per GPU decreases at the same rate as in training within a zone (A). The per-GPU speedup with additional hardware reduces at the same rate for either setup (between 0.05 and 0.06). This results in two observations: First, communication overhead scales linearly with the number of peers. Second, we only have to pay the penalty for transatlantic training once. However, we cannot expect a significant improvement in communication efficiency when we increase the amount of available local resources.

Summarizing, with an transatlantic setup, CV achieves a virtually identical maximum speedup of 3.2x with 8 GPUs compared to A-1 (B-8 is 2% slower than A-8), while NLP is more affected by lower network bandwidth and only achieves a speedup of 2.15x (B-8 is 22% slower than A-8). The transatlantic training penalty is applied once; however, it does not affect the relative scaling with additional compute resources.

(C) Intercontinental scalability. To take geo-distribution to the extreme, we spawn VMs in up to 4 regions: USA, EU, ASIA, and AUS, to see how much worse bandwidth affects the training throughput (C-3,4,6,8 in Table 2).

How does the intercontinental penalty investigated in (B) affect deployments with a single GPU on each continent? Comparing the A-3 and C-3 experiments with three local versus three fully remote GPUs, CV is only 5% slower, while NLP suffers a 34% drop in throughput (Figure 9a) and does not even reach the baseline single GPU performance (A-1). The peak egress for each region was 318, 258, and 237 Mb/s for the US, EU, and ASIA, respectively. Since our bandwidth measurements were 210 and 130 Mb/s from the US to the EU and ASIA, respectively (Table 3), this suggests that the averaging was done over the US node and not an N-to-N all-reduce (a detailed analysis of how averaging affects bandwidths is discussed in Section 6). Thus, the limiting factor was the US-ASIA connection at 130 Mb/s rather than the 80 Mb/s from EU-ASIA. The same trend continues with the C-4 run, which adds AUS as a continent with one additional VM. As we know from the transatlantic experiments (B) that an additional continent has a detrimental effect on throughput, which, for the four continents experiment, C-4, results in a 9% slower throughput for CV and 36% slower for NLP compared to the A-4 runs (Figure 7a). Again, the US VM is used as an averaging intermediary with a peak egress of 365 Mb/s, while the other continents are between 318 and 330 Mb/s. When comparing the two continents (B-4) versus four continents (C-4) experiments, one GPU on each continent (C-4) is slower by 6% for CV and 20% for NLP compared to two GPUs on two continents (B-4). This reinforces that local hardware should be preferred whenever possible. However, we are always faster than the baseline (A-1), starting from 4 GPUs in both the transatlantic and intercontinental settings. While these experiments were specifically designed to be a worst-case scenario, what about a more balanced GPU distribution with at least two GPUs in each region?

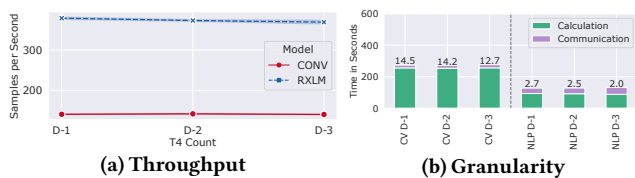


Figure 10: Multi-cloud performance for CV and NLP.

When comparing the C-6 experiment with two GPUs in three continents to the local A-6 experiments, the throughput slowdown is almost identical (CV 7%, NLP 35%) as with C-4 (CV 9%, NLP 36%) to A-4. Scaling further to two GPUs in four continents, C-8 is slightly slower at NLP (41%) compared to C-4 (36%) to their respective local runs (A-8 and A-4), due to the decreasing granularity of 0.4 (Figure 9b). The small granularity removes the additional gain of four more GPUs since the task is no longer suitable for distributed training. However, as the CV task is still at a granularity of 3.33 on C-8, it reaches a speedup of 3.02x, only 7% slower than the fully local A-8 experiment. The peak egress of 678 Mb/s was also reached on one US VM, while the remaining VMs were between 450 and 550 Mb/s. These observations show that adding another continent does not significantly reduce throughput when training on three continents with at least two VMs.

In summary, while local compute is the best choice for maximum throughput, for high granularity tasks like CV, even distributing VMs over four continents only slows down performance by 7%. However, intercontinental training leads to a significant penalty on a task with lower granularity, like NLP, resulting in a performance drop of 41% (C-8) compared to the fully local experiment (A-8). Finally, each additional region introduces a constant penalty that is not amortized by adding local hardware, which should be considered when running geo-distributed training setups.

5 MULTI-CLOUD PERFORMANCE

Using multiple cloud providers makes sense if we want to use resources cost-effectively and have additional reliability. In our scenario, we are interested in what throughput per \$ can be expected and if any barriers prevent multi-cloud training. However, one can also consider the data center’s carbon footprint, which can change depending on the season and time of day [6].

Table 4: Average multi-cloud throughput and latency.

| (a) Single stream TCP throughput in Gb/s. | | | | (b) ICMP Latency in ms. | | | | | |
|---|----|------|------|-------------------------|-------|----|-------|------|-------|
| From | To | GC | AWS | Azure | From | To | GC | AWS | Azure |
| GC | GC | 6.35 | 1.52 | 0.45 | GC | GC | 0.71 | 15.3 | 51.22 |
| AWS | GC | 1.81 | 4.87 | | AWS | GC | 13.85 | 0.15 | |
| Azure | GC | 0.47 | | 7.63 | Azure | GC | 49.80 | | 1.56 |

We have compiled the current prices for spot and on-demand instances for T4 GPUs with 8 CPU cores and the egress costs for three well-known cloud providers, GC [5], AWS [2], and Azure [9] (Section 1). There are two different pricing concepts. On the one hand, there are GC and Azure, which offer relatively cheap instances, with 69% and 73% savings over on-demand pricing, respectively, and relatively expensive egress charges between continents of up to \$0.15/GB. On the other hand, there is AWS, where the spot instance is only 51% cheaper than the on-demand instance and more than twice as expensive as GC or Azure. However, the egress fees here are much

cheaper at only \$0.02/GB. Because of the additional offerings around compute, such as networking, identity and cost management, and tooling, it is not easy to fairly compare cloud providers. Therefore, we will limit ourselves to network and VM costs.

With the multi-cloud experiments from this section, we want to evaluate the following scenarios: First, partially switching from one provider to another without stopping the training. Second, scaling resources in the same region when one of the cloud providers is already at capacity for spot-priced VMs or the current price is too high [24]. We know from Section 4 that scaling resources in the same location can significantly improve performance, which may only be possible using additional cloud providers.

Experimental design. To enable a fair comparison between the cloud providers, we rented hardware most similar to each other in the same region. We used each provider’s default settings and only changed hardware specs. For GC, it is the same instance as in Section 4. At AWS, it is a g4dn. 2xlarge with eight cores and 32 GB in the us-west-2c region. Unfortunately, we had to make two compromises with Azure. There was only the combination of four cores and 30 GB RAM (NC4as_T4_v3), and there were no T4 GPU resources available in the us-west, so we had to fall back to us-south-2.

The network profiling between all cloud providers in Table 4 shows that their intra-cloud connectivity is comparably fast with 6.4, 4.9, and 7.6 Gb/s for GC, AWS, and Azure, respectively. All connections are mostly symmetric, with inter-cloud connectivity between GC and AWS providing up to 1.8 Gb/s and a ping of 15.3ms, indicating that while they are likely not in the same data center, they are close to each other and connected to the same Internet exchange point. However, connectivity to Azure could be better since it operates in a different zone, with a bandwidth of 0.5 Gb/s and a ping of 51ms.

Our experimental setup consists of four GPUs with equal contributions from each cloud provider. D-1 is the baseline with four GPUs at GC, D-2 with two GPUs each at GC and AWS, and D-3 with two GPUs at GC and Azure. We compare moving two VMs to a different cloud provider to see the impact on cost and throughput.

(1) No inter-cloud throughput penalty. Figure 10 shows the throughput and granularity of each multi-cloud experiment. CV and NLP runs have essentially identical throughput regardless of the combination of cloud providers. Only the D-3 experiments show a very slight slowdown in communication time, reflected in the lower granularity score (Figure 10b) of 12.72 in CV and 1.99 in NLP compared to the D-1 baseline scores of 14.48 and 2.73, respectively. Actual throughput was between 1-2% slower than the baseline, which is negligible and only related to the slightly worse connection to the Azure data center. These results confirm our observation from Section 4 that network connectivity determines scalability, and one can easily train in a multi-cloud scenario.

(2) External egress costs can overshadow VM costs. One drawback to training in multiple regions or zones is that egress traffic can incur additional costs depending on the cloud provider. We have summarized the cost of egress traffic within a zone (intra-zone), between zones in each region (inter-zone), and between continents in Section 1. Notably, any traffic to Oceania (Australia, New Zealand, and others, abbreviated as OCE) generates the highest cost of \$0.15/GB for GC. We have broken down the costs for the multi-cloud experiment in Figure 11a on an hourly per-VM basis. With only four peers in the D-1/2/3 experiments, we have an N-to-N communication, i.e.,

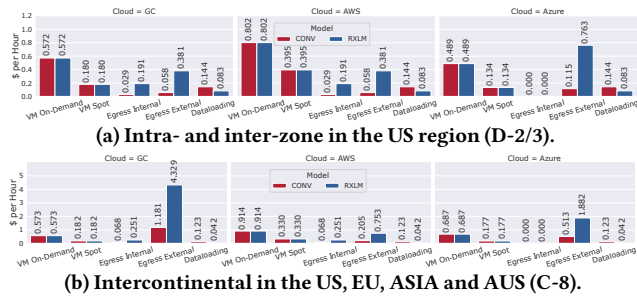


Figure 11: Costs breakdown for D-2/3 and C-8 experiments.

each peer sends its gradients to every other peer. This means that $\frac{1}{3}$ of the egress was internal to the partner VM in the same cloud, and the remaining $\frac{2}{3}$ went to the remaining two peers in the other cloud.

First, loading data from Backblaze costs \$0.01/GB from anywhere in the world, which gives us a rate of \$0.144/h for the CV and \$0.083/h for the NLP experiments. Even when CV throughput is less than half of the NLP model (Figure 10a), images are much larger than text, resulting in a higher data rate. While this is close to the spot instance costs of GC (\$0.18/h) and Azure (\$0.134/h), these are one-time costs until the entire dataset is downloaded and retrieved from the disk cache, assuming large enough local storage. A more detailed comparison of cloud provider storage offerings is beyond our scope, but current prices range from \$0.02/GB to \$0.14/GB in various GC regions, making our setting (B2) competitive.

Second, the external egress costs for the NLP experiments are very high compared to the other costs. They are 2.2x higher than the spot instance for GC and 5.7x higher for Azure, as the traffic costs in the US zone are \$0.01/GB and \$0.02/GB, respectively. The Azure cost is even higher (\$0.763/h) than the on-demand instance price of \$0.489/h. The CV experiments are much less affected due to the smaller model size, but Azure still manages to almost match its spot instance price of \$0.134/h with the external egress cost of \$0.115/h.

Finally, the total compute cost, including egress and data loading in this multi-cloud constellation, is the sum of all the cloud providers' prices times the number of VMs used. For the CV experiments, GC, AWS, and Azure cost \$0.762/h, \$1.192/h, and \$0.363/h, respectively, making the combination of GC with Azure 42% cheaper than GC with AWS. For the NLP experiments, GC, AWS, and Azure cost \$0.835/h, \$1.05/h, and \$0.973/h, respectively, and GC combined with Azure is better than GC with AWS by a smaller margin of 3.9%. However, the intercontinental network egress prices for both GC and Azure are up to 15 times higher than the inter-zone prices, so what about the cost-effectiveness compared to geo-distributed experiments?

(3) Geo-distributed egress can incur most of the cost. To illustrate the cost of intercontinental training, we use our C-8 experiment with two VMs in four continents from Section 4 to plug the cost for each cloud provider. The egress costs are calculated slightly differently than in the D-2 and D-3 experiments because four groups of two VMs average locally and then distribute the gradients across the other groups. This results in $\frac{8}{20}$ internal egress calls (two calls between each group), $\frac{6}{20}$ intercontinental egress calls (two calls between three regions), and $\frac{6}{20}$ AUS egress calls (three regions share their gradients with AUS and vice versa).

Figure 11b shows the resulting egress traffic cost per VM. The high cost between continents scales to a multiple of the remaining

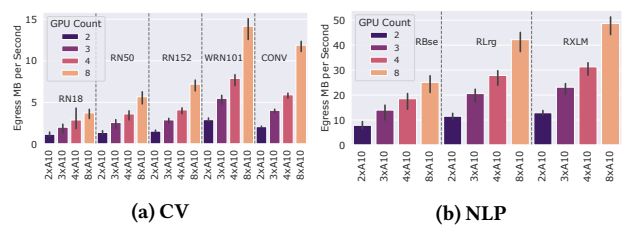


Figure 12: Baseline egress rate on 2-8 A10 GPUs.

cost for CV and NLP with GC and Azure. For NLP, the external egress cost for GC is \$4.329/h, more than 90% of the total cost per VM (\$4.804/h). Even with Azure having a more moderate rate of \$0.02/GB for intercontinental communication and only \$0.08/GB for OCE traffic, it still results in \$1.882/h external egress cost (\$2.101/h total). This is in contrast to AWS, which has a cap of \$0.02/GB to any location, resulting in the best total cost of \$1.376/h per VM. The relatively high AWS instance cost compares favorably to the other cloud providers regarding geo-distributed training. Keeping egress traffic in mind when deciding to scale to other continents is essential, as it can be the most significant part of the total cost. This raises another question: If egress traffic matters so much, how does model size affect it?

(4) Small models have lower egress rates than larger models.

Model size affects two parts of the distributed training time. First, larger models tend to have slower averaging rates, but more data movement costs due to their size. However, larger models are also averaged less frequently because they take longer to perform a step. To analyze this, we review the experiments in Section 3, where we evaluate different model sizes and GPU counts. Figure 12 shows the average egress rate over each experiment's runtime for both CV and NLP from two to eight A10 GPUs. The trend is clear: the smaller the model, the lower the egress rate for all GPUs (e.g., RN18 vs. RN50). This is surprising, as the "square-cube" law [37] states that with a decrease in parameters, the calculation time will decrease quadratically while the communication time decreases linearly. This means that with a sufficiently small model, most of the training will consist of communication time, and the egress rate would increase, as it is defined through $\frac{\text{parameter count}}{\text{calculation time}}$. However, we find that even with our smallest model, RN18, with 11.7M parameters and eight A10 GPUs, we are still not at the point where the communication time takes up most of the time.

In summary, multi-cloud training is generally possible and can be cost-effective when keeping the egress costs and granularity in mind. Regardless of the cloud provider, staying in the same region is preferred, with the US having the most favorable egress price offers. A significant portion of the cost may be hidden in egress costs, accounting for more than 90% of the total cost in our NLP experiments in GC and Azure. Based on the additional egress costs alone, renting on-demand hardware may be more advantageous than using spot instances between different regions. CV training is generally more calculation- than communication-heavy, resulting in slightly higher data-loading but fewer egress costs. However, from our experiments, this is a favorable trade-off because data-loading is much cheaper than egress costs.

6 HYBRID-CLOUD PERFORMANCE

Can augmenting on-premise hardware with cloud resources be worthwhile to speed up DL training? In this section, we examine two

Table 5: Average hybrid-cloud throughput and latency.

| (a) Single stream TCP throughput in Gb/s. | | | | (b) ICMP Latency in ms. | | | | | |
|---|----|-------|-------|-------------------------|----------------|----|-------|--------|--------|
| From | To | EU T4 | US T4 | US A10 | From | To | EU T4 | US T4 | US A10 |
| RTX8000 | | 0.45 | 0.06 | 0.05 | RTX8000 | | 16.73 | 150.80 | 159.05 |
| DGX-2 (8xV100) | | 0.55 | 0.08 | 0.07 | DGX-2 (8xV100) | | 16.19 | 150.27 | 158.54 |

Table 6: Hybrid- vs. cloud-only throughput for the (E) setting.

| Model | Setup | RTX8000 | E-A-8 | E-B-8 | E-C-8 | 8xT4 | 8xA10 |
|-------|-------|---------|-------|-------|-------|-------|--------|
| | | CONV | 194.8 | 316.8 | 283.5 | 429.3 | 261.9 |
| RXLM | | 431.8 | 556.7 | 330.6 | 223.7 | 575.1 | 1059.9 |

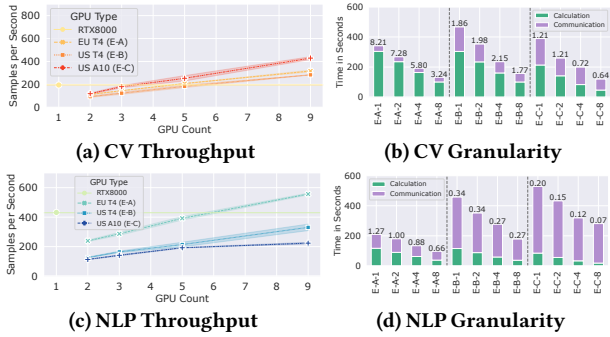


Figure 13: Hybrid-cloud experiments for the (E) setting.

settings: (E), where a consumer-grade GPU, the RTX8000, is deployed on-site, and (F), where a server-grade node, the DGX-2 (8xV100), is deployed on-site. We vary the extra resources, between one to eight T4 EU ({E,F}-A), T4 US ({E,F}-B) and A10 US ({E,F}-C) GPUs.

Experimental design. In both settings, we want to investigate how to extend local hardware with cloud resources and when this leads to better throughput. The cloud resources, in this case, are the same US/EU GC T4 instances as in Section 4 and the US LambdaLabs A10 GPUs from Section 3. We double the number of cloud VMs with each increment, starting with one additional GPU (i.e., E-A-1) until we have eight additional cloud VMs (i.e., E-A-8). This allows us to compare the same hardware in the EU and the US, and slightly weaker, local hardware (EU T4) and better, but more distant hardware (US A10).

Both the (E) and (F) setups share the network uplink between 450 and 550 Mb/s to the EU datacenter in Belgium, as they are located in the same building in Europe (Table 5). However, as this is not a Google-owned datacenter, the traffic is partly going over the public internet, which results in a lower bandwidth of 50 and 80 Mb/s to the US-based VMs compared to 210 Mb/s between the US and EU GC datacenters (Table 3a).

(E) Consumer-grade setting. The results follow the same trend as in Section 4. The CV task has a higher granularity of 8.21 with 2 GPUs at E-A-1 than NLP (1.27) (Figures 13b and 13d), and scales regardless of the location of the cloud resources (Figure 13a). We almost match the baseline throughput of 195 SPS at 5 GPUs in all settings for CV (E-A-4, E-B-4, E-C-4). The best throughput was reached at E-C-8 with the US A10 GPUs with 429 SPS. For NLP, only the E-A-8 experiment beats the baseline with a speedup of 1.29x and 556 SPS due to the low granularity and the intercontinental base penalty for the US experiments.

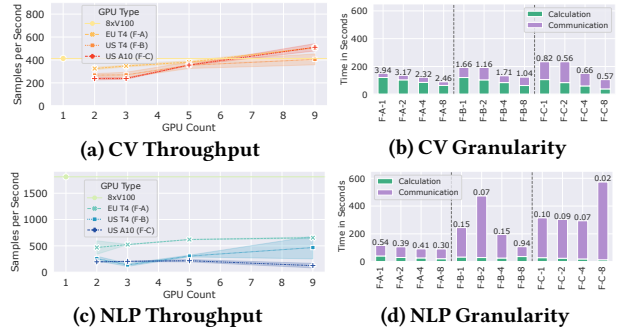


Figure 14: Hybrid-cloud experiments for the (F) setting.

However, is combining on-premise and remote cloud resources better than using the cloud without paying the intercontinental bandwidth tax? To analyze this, we compare the (E) experiments with the 8xA10 experiment from Section 3 and 8xT4 experiment from Section 4 in Section 6. First, the 8xA10 experiments are the fastest for both CV and NLP, which removes the respective hybrid-cloud combination from contention (E-C-8). Second, the 8xT4 experiments for NLP are faster than any other hybrid-cloud setup, making the cloud-only solution favorable. Finally, while we always beat the baseline 8xT4 CV throughput (261.9 SPS), but in the case of E-B-8 (283.5 SPS), just barely. The throughput of E-A-8 (316.8 SPS) makes the hybrid-cloud setup the most favorable in terms of relative GPU scaling (32.5 SPS per GPU), but it does not come close to the best cloud-only throughput of 8xA10 with 620.6 SPS.

Summarizing, the cloud-only experiments are the fastest overall due to their single-GPU throughput and locality. Adding cloud resources to on-premise hardware leads to a high communication time, which is not compensated by the additional processing speed of the GPUs. Proximity to the on-premise hardware is essential, as the more local cloud resources (E-A-8) consistently resulted in a better throughput than the same remote cloud resources (E-B-8).

(F) Server-grade setting. The baseline throughput is significantly higher compared to the RTX8000, with a much more powerful 8xV100 DGX node to 413 SPS for CV and 1811 SPS for NLP (Figures 14a and 14c) via PyTorch data parallelism [26]. This increases the penalties from Section 3, leading to the only speedup from baseline for CV in experiments F-A-8 (507 SPS) and F-C-8 (510 SPS). This is surprising, as the older T4 GPUs in the EU perform similarly to the much newer A10 GPUs in the US, showcasing the trade-off between slower, local compute and faster, remote compute. The granularity of 2.46 for F-A-8 shows that there is enough calculation time to distribute, while the F-C-8 experiments spend $\approx 62\%$ of the total training time on communication with a granularity of 0.57 (Figure 14b). The NLP experiments never reach the baseline throughput of the 8xV100 due to using most of the time for communication. The NLP F-B and F-C experiments mainly consist of communication (Figure 14d) with a granularity of up to 0.02, which results in a nonlinear, unstable training time due to the minimum matchmaking time issue (2) from Section 3.

In summary, the hybrid-cloud experiments conclude that while on-premise hardware can be augmented with cloud resources, it will likely be cost-efficient if all resources are on the same continent. Using only cloud resources is more advantageous if the on-premises hardware is not co-located.

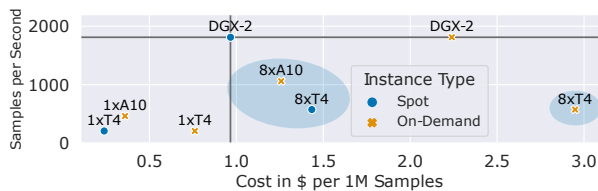


Figure 15: Cost to throughput tradeoff for RoBERTaXLM at different instance types. Our training setups (circled), that are due the low granularity of the NLP model, neither cheaper, nor faster than the centralized offering (DGX-2).

7 FURTHER INSIGHTS

Communication time can decrease with more peers. Let us compare the granularity of the experiments for E-B (Figure 13b), which uses T4 GPUs in the US as an additional cloud resource. Both the computation and communication time decrease with the number of GPUs, even increasing the granularity from 1.98 at E-B-2 to 2.15 at E-B-4. This is surprising since, usually, with more peers, the communication time should increase, and the US-EU communication bottleneck should slow us down to the same extent as the E-B-1 experiment. This reduction is a Hivemind-specific anomaly, as it uses a single TCP stream per peer. With TCP, there needs to be an acknowledgment (ACK) of each packet by the receiving peer, which is impacted by the connection’s latency. In our high latency network between continents, the round trip time (RTT) of 300-318ms limits the maximum bandwidth a single TCP stream to 50-80 Mb/s. However, a way to improve link utilization is to use multiple streams, one for each peer, which we encounter in experiments E-(B|C)-2,4,8. To verify the potential gains, we perform a microbenchmark of the multi-stream bandwidth from the RTX8000 to the EU and US data centers. Although there is wide variation, likely due to network utilization, with 80 clients, we achieve a maximum bandwidth of 6 Gb/s within the EU and up to 4 Gb/s to the US. While larger peer groups and, consequently, larger models benefit from multi-peer communication by default and do not see significant changes in communication time, small models in unevenly distributed VMs setups can be disproportionately affected. The same trend can be observed in all high latency experiments (i.e., between the EU and the US), e.g., E-B, E-C for CV and NLP (Figures 13b and 13d, and F-B and F-C for CV (Figure 14b)). In summary, uneven distribution of computational resources in high-latency networks (e.g., intercontinental) can reduce communication time with Hivemind due to more parallelism, lessening the impact of low bandwidth for a single data stream.

Cost analysis. The DGX-2 (8xV100) node from Section 6 represents server-grade hardware that could be used to train models. However, how does it compare in throughput per \$ to all of our distributed cloud experiments? The Figure 1 (CV) and Figure 15 (NLP) show the complete cost analysis of the DGX-2, the 8xT4 experiments, and the 8xA10 experiments for spot and on-demand pricing. We use the internal egress costs from Figure 11a as a reference for the 8xT4 setup. For simplicity, we compare the spot pricing without interruptions, as we assume that a new VM can be spun up fast enough not to affect the training throughput in the long run. We mark the centralized baseline (DGX-2) cost per 1M samples and the throughput in samples per second with a horizontal and vertical line. This means that we are cheaper to the left to the vertical line, and above

the horizontal line, we are faster (and vice versa). We circle the new value propositions that we enable in both figures. Our hardware setups have additional key characteristics: They are resilient by default to interruptions due to running in a decentralized fashion and they enable the combination of more GPUs than cloud providers offer in a single node. Currently, common hardware configurations (DGX) allow up to eight GPUs connected via NVLink, and with older hardware, only up to 4xT4s connected via PCIe at 10 GB/s between GPUs (with GC). We were able to combine eight single GPU nodes from GC and LambdaLabs to create competing performance and price setups without dedicated GPU interconnects.

A spot DGX-2 costs at the time of writing \$6.30/h (\$14.60/h on-demand) in GC US, which makes it the best value proposition for the low granularity NLP task. It is followed by the 8xA10, which are 41% slower and 30% more expensive than the DGX-2 (Figure 15). The 8xT4 experiments are even more expensive, as the internal egress costs take up more than half of the costs, making them the worst value proposition. However, for CV, we manage to provide two new offerings: First, the 8xA10, which is both 50% faster and 49% cheaper than the DGX-2, and 8xT4, which is 58% cheaper than DGX-2, while being 37% slower (Figure 1). The CV model can be scaled more easily due to its initially high granularity, which makes the very competitive offering of \$0.6/h per A10 from LambdaLabs an excellent value proposition. However, while we only evaluated eight T4 GPUs for our GC-based experiments, with a granularity of 5.19 (CV A-8 in Figure 7b), there is ample space to scale even further. It is important to note that LambdaLabs does not charge for any data egress, but GC does with \$0.01/GB, and the 8xT4 experiment is still cheaper. While LambdaLabs is often at capacity, Google Cloud positions itself as a hyperscaler with the advantage of rarely being at max occupancy.

We also evaluated the performance of the 4xT4 PyTorch DDP [26] for CV with the best available multi-T4 node on GC (4xT4). The NLP experiments ran OOM. Since the DDP 4xT4 runs on a single node, it causes no interconnect costs and is priced at \$0.96 per 1M samples at spot pricing, while our 8xT4 setup costs \$1.77 per 1M samples (84% more expensive). However, the 8xT4 setup has a higher throughput of 262 SPS (26% faster) compared to the 4xT4 node (207 SPS). This higher speed is not available at the price point of the 4xT4 node. Moreover, the 8xT4 setup has the potential for further scaling, which we discussed in detail in Section 4.

In summary, the lower spot prices for older GPUs allow us to train models more cost-efficiently when task granularity allows it and get more value per \$ when training on the 8xT4 or 8xA10 compared to an DGX-2 node. Combining multiple nodes with single GPUs with lower bandwidths enables scaling that was previously impossible to achieve without resorting to much more powerful GPUs. Distributed spot instance pricing opens up a new value proposition compared to on-demand offerings that can even compete with the competitive pricing of smaller cloud providers.

Spot VM Interruption Frequency. While we used low spot prices as a cost-saving argument in our experiments, we did not elaborate on the most significant drawback - the possibility of being terminated by the cloud provider at any time. There is already some research on how different cloud providers track the interruption frequency and can be used for varying workloads to achieve a positive \$-per-throughput effect [24, 42, 43].

Interruption affects three aspects: First, the interruption frequency is defined by AWS as the number of VMs terminated in the last 30 days, which is between 5 and 20% [3]. This value was not representative during our experiments with any cloud provider, as we noticed that it is highly dependent on the time of day of the zone.

Second, the time needed to setup a VM until training starts. The startup time of a VM depends on the cloud provider (e.g., a preconfigured image) and the one's technology stack (e.g., Docker, Kubernetes, Ansible). In our experience, VM startup time ranges between seconds to minutes with manual deployment taking up to 10 minutes. Although startup time can be improved, model training typically takes multiple hours or days, making it a less impactful optimization.

Third, the time required for the new peer to synchronize the training state with other peers. In our experience, this took at worst two hivemind epochs due to the averaging starting before synchronization is finished. While it is possible to create a hivemind epoch that is short enough to prevent new peers from joining, this only happens with a low enough granularity where scaling is not beneficial anymore as we are mostly communication bound.

Finally, while the VM setup and synchronization of the training state take time, the interruption frequency significantly affects the final throughput. We faced difficulties acquiring even a single spot VM during our GC experiments during daylight hours. This highlights the need for systems like SkyPilot [43], which utilizes automation to deploy spot instances across various clouds and zones. In our case, the interruption frequency can be used as a penalty on the training throughput, i.e., a 5% interruption frequency over the entire training time means roughly a 5% slower training.

8 LESSONS LEARNED

We find it important to summarize our findings more generically to provide guidance for DL practitioners that want to perform distributed spot training. These lessons are based on the Sections 3 to 6.

Small model training still scales. We have shown that models between 12M-560M parameters can be trained in a decentralized, distributed fashion achieving a speedup of up to 4.37x on eight Ampere-GPUs. The limiting factor as to when a model is suitable for (geo)-distributed training is the target batch size which all peers need to accumulate until synchronization happens. We found a TBS of 32K suitable to not only train in a single zone, but even see a speedup when using VMs in four different continents. As long as the optimizer can handle big-batch training and the dataset is big enough to accommodate large batches, the remaining issue to find the base granularity of the model to decide how to scale it cost-effectively. Finally, we found that small models induce less traffic over larger models over time, even at a much higher averaging rate, making them better suited for cost-efficient training than large models.

Egress costs can take up most of the total cost. Egress pricing for the NLP experiments overtook the spot and the on-demand costs of T4 GPUs when training on four continents or even in two zones. For example, RoBERTaXLM's high throughput and parameter count require more data to be sent between peers during averaging due to smaller granularity. Under the current pricing models, AWS has the best value for geo-distributed training, while GC and Azure are best at training in a single zone. The biggest cost-saving potential lies in cloud providers that do not charge for egress at all, like LambdaLabs.

Granularity is important to evaluate scalability. We found that the ratio between calculation and communication time, granularity, is the most important metric to track when deciding on distributed training suitability. It enables us to compare the scalability potential between different models on the same hardware due to summarizing their model size and throughput ratio. Additionally, it gives a value to the cost-efficiency: With a granularity of exactly 1, the potential speedup when doubling the number of VMs is, at best, 1.33x due to halving the calculation time. However, with a granularity of 10, the speedup with double the VMs is, at best, 1.83x due to the communication time playing a less significant role. With this, we can estimate training performance with additional resources.

Geo-distributed multi-cloud training is possible and is cost-efficient. Even with the current teething pains of Hivemind, we got a speedup in all of our experimental setups of intra-zone, transatlantic, and intercontinental settings as long as the granularity of the task permitted it. Using older and cheaper Tesla GPUs at spot pricing is not only more cost-efficient than the DGX-2 offering, but even trumps the competitive pricing model of LambdaLabs, all while including egress costs. Our network profiling showed that the current training limitations are not primarily the bandwidth but rather the intercontinental latency and the task's granularity. If the granularity is already low at high bandwidth, it can only worsen when used in a high latency, low bandwidth network. When considering both, estimating the potential cost-savings of investing in a multi-/hybrid-cloud scenario is possible.

9 RELATED WORK

Decentralized deep learning. Training with unreliable peers has been studied in a collaborative setting, resulting in the Distributed Deep Learning in Open Collaborations (DeDLOC) [17] algorithm, on which the Hivemind framework [39] is based. It can interpolate between traditional distributed DL algorithms like parameter servers [25], decentralized SGD [27], or All-Reduce SGD [1]. We used the Hivemind framework for all of our experiments, as it provided the base for training on spot instances in high latency, low bandwidth networks.

SWARM [37] applies both previous techniques and adds model parallelism to the mix by creating pipelines between nodes and rebalancing them in case of failures. The authors find a crucial insight in the "square-cube" law, which argues for better training scalability with larger model sizes; as the size increases linearly, so does the communication time, while the calculation time increases quadratically. We add to that by analyzing distributed training for smaller model sizes that pose different trade-offs. We show that while the square-cube law still holds for increasing model sizes, under consideration of granularity, we can still train small models.

Decentralized deep learning on heterogeneous hardware with slow interconnects can benefit the training of foundation models. To achieve this, model and pipeline parallelism can be used in addition to data-parallel training [45]. This is a complementary work to ours, since we target smaller models and weaker hardware.

Deep learning on spot instances. DeepSpotCloud [23] is a system that uses the AWS API to automatically migrate a DL task with checkpointing whenever the spot instance is terminated. The authors note that the volatility of GPU instance pricing and interruptions have a unique pattern compared to non-accelerated VMs,

and solve this by using intercontinental provisioning. We noticed the same trends of high interruption ratios in our experiments. However, we have shown that geo-distributed training is possible until granularity permits it, which poses a possibility for ever-migrating training between continents without checkpointing.

Amazon Sagemaker [14] is an AWS service that allows to perform ML under budget constraints. For training, it supports spot VM migration until a cost threshold is reached by checkpointing the progress. However, it lacks the option of training on multiple spot VMs. It can do either spot instance training on DGX-like nodes or combine multiple on-demand nodes with PyTorch DDP (or similar), but not both. This eliminates the potential of accelerating the training process with more GPUs that do not fit a single spot-provisioned hypervisor.

The analysis by Yang et al. [42] investigates maximizing a target accuracy from a spot pricing versus time perspective. Linear programming was used to decide how to provision the VMs with different cost-utility trade-offs. While this shows the potential of utilizing multiple clouds and continents for non-distributed tasks, we evaluated the distributed spot training problem from the throughput, cost, and model size perspective on different hardware setups. By including our insights, their technique for scheduling on spot instances could be adapted to optimize the total throughput of all peers.

Skypilot [43] is a broker system where users can submit their hardware requirements, and it tries to provision the necessary resources on any supported cloud. It features a preemption analysis that counts the number of interruptions in a zone and can decide to migrate whenever they cross a certain threshold. We have shown that multi-, hybrid-cloud, and geo-distributed training is possible, and by combining our insights, it would open up auto-migrated, decentralized DL training for the best spot prices in the world.

10 CONCLUSION

This paper analyzes multi- and hybrid-cloud training in a decentralized fashion on spot instances. We define the lower bounds of model sizes that can be scaled cost-efficiently using the granularity metric to estimate their suitability for distributed training in low-bandwidth, high-latency situations. We show that training on multiple cloud providers and four continents still scales with additional compute resources. Alternatively to the current use of spot instances in DL, we show the potential of using spot instances in a distributed, decentralized way by being more cost-efficient with eight T4 instances over a DGX-2 from the same cloud provider while paying additional egress costs. Finally, we provide an intuition about where costs in such a training scenario come from and how different model sizes from CV and NLP affect throughput and costs. Our work empowers practitioners to utilize spot-priced instances for distributed deep learning with relatively small models. Our insights show some potential that can further improve distributed training performance, such as optimizers with higher minibatch sizes and improvements regarding the communication time with, e.g., better compression.

11 APPENDIX: ASR CASE STUDY

We perform a case study on Automatic Speech Recognition (ASR) to showcase spot training on weaker GPUs. Whisper [34] is a state-of-the-art ASR model trained on 680,000 hours of labeled data to

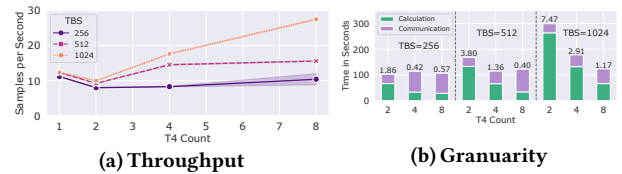


Figure 16: WhisperSmall performance with varying TBS.

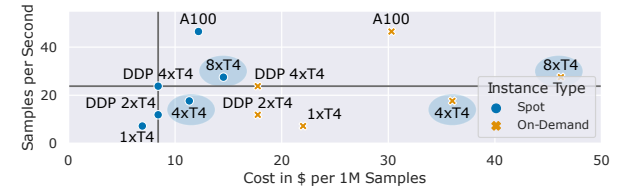


Figure 17: Cost to throughput tradeoff for WhisperSmall at TBS=1024 with different instance types. Our training setups (circled) provide mixed result of being slightly faster and more expensive than comparable, centralized DDP offering.

transcribe audio. It features different sizes, from 37.8M to 1.5B parameters, and was trained with a minibatch size of 256. We use the Commonvoice [11] dataset, preprocessed to Log-Mel spectrograms. In our distributed experiments, we start with a TBS of 256 and increase to 512 and 1024 to combat potential granularity issues. Due to memory constraints, only three model sizes (Tiny, Base, Small) were trainable on the T4 GPU. Unfortunately, the original TBS of 256 was not large enough to train the relatively small models due to their small granularity (0.04, 0.14 and 0.57 at 8xT4, respectively) with no performance benefits. The only model showing scaling potential is WhisperSmall, with a granularity of 1.8 with 2xT4. However, when scaling the target batch size to 512 and 1024, we see some benefit over the single GPU runs for the WhisperSmall model (Figure 16). By effectively increasing the amount of computation by the factors of 2 and 4, we can generate a speedup of 1.27x and 2.2x with 8xT4's for the TBS 512 and 1024, respectively. When compared to other hardware setups, our A100 80GB GPU and the best multi-T4 GPU on GC (4xT4) with Pytorch DDP (Figure 17) have almost double the throughput at 46 SPS and are slightly slower at 24 SPS, respectively, compared to our 8xT4 setup which runs at 28 SPS. This outcome is not surprising due to the generational leap in architecture for the A100 and the slower interconnect with our 8xT4 experiments compared to a single 4xT4 node (see Section 3 for a detailed throughput analysis). The proposed cost-throughput ratio is mixed: the A100 is at \$12.19/1M samples, the DDP 4xT4 is at \$8.41/1M, and our 8xT4 is at \$14.53/1M. Our proposed setup is slightly more expensive than the A100, and it will not scale beyond eight T4 GPUs due a granularity at 1.17, leaving the A100 as the fastest and the DDP 4xT4 setup as the cheaper but slower alternative. Despite these results, our proposed setup has several benefits, including resilience for spot interruptions, interruption-free migration to the lowest cloud prices, and the possibility to scale the GPU count up as long as granularity permits it.

ACKNOWLEDGMENTS

This work is funded in part by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) - 392214008.

REFERENCES

- [1] [n.d.]. Horovod: fast and easy distributed deep learning in TensorFlow, author=Sergeev, Alexander and Del Balso, Mike, journal=arXiv preprint arXiv:1802.05799, year=2018. [n. d.].
- [2] 2023. Amazon AWS. Accessed: 19 May 2023, aws.amazon.com.
- [3] 2023. Amazon AWS Spot Pricing. <https://aws.amazon.com/blogs/compute/new-amazon-ec2-spot-pricing/>. Accessed: 2023-09-27.
- [4] 2023. Backblaze. <https://backblaze.com/>. Accessed: 2023-10-05.
- [5] 2023. Google Cloud. Accessed: 19 May 2023, cloud.google.com.
- [6] 2023. Google Cloud Region Picker. <https://cloud.withgoogle.com/region-picker/>. Accessed: 2023-10-05.
- [7] 2023. Hivemind GAC Issue. <https://github.com/learning-at-home/hivemind/issues/566>. Accessed: 2023-10-05.
- [8] 2023. LambdaLabs. Accessed: 19 May 2023, lambdalabs.com.
- [9] 2023. Microsoft Azure. Accessed: 19 May 2023, portal.azure.com.
- [10] Alex Aizman, Gavin Maltby, and Thomas Breuel. 2019. High Performance I/O For Large Scale Deep Learning. In *2019 IEEE International Conference on Big Data (Big Data)*. 5965–5967. <https://doi.org/10.1109/BigData47090.2019.9005703>
- [11] Rosana Ardila, Megan Branson, Kelly Davis, Michael Henretty, Michael Kohler, Josh Meyer, Reuben Morais, Lindsay Saunders, Francis M Tyers, and Gregor Weber. 2019. Common voice: A massively-multilingual speech corpus. *arXiv preprint arXiv:1912.06670* (2019).
- [12] Alexander Borzunov, Max Ryabinin, Tim Dettmers, Quentin Lhoest, Lucile Saulnier, Michael Diskin, and Yacine Jernite. 2022. Training Transformers Together. In *NeurIPS 2021 Competitions and Demonstrations Track*. PMLR, 335–342.
- [13] Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised Cross-lingual Representation Learning at Scale. arXiv:1911.02116 [cs.CL]
- [14] Piali Das, Nikita Ivkin, Tanya Bansal, Laurence Rouesnel, Philip Gautier, Zohar Karnin, Leo Dirac, Lakshmi Ramakrishnan, Andre Perunicic, Iaroslav Shcherbatyi, Wilton Wu, Aida Zolic, Huibin Shen, Amr Ahmed, Fela Winkelmolen, Miroslav Miladinovic, Cedric Archembeau, Alex Tang, Bhaskar Dutt, Patricia Grao, and Kumar Venkateswar. 2020. Amazon SageMaker Autopilot: A White Box AutoML Solution at Scale. In *Proceedings of the Fourth International Workshop on Data Management for End-to-End Machine Learning (Portland, OR, USA) (DEEM'20)*. Association for Computing Machinery, New York, NY, USA, Article 2, 7 pages. <https://doi.org/10.1145/3399579.3399870>
- [15] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 248–255.
- [16] Tim Dettmers. 2016. 8-Bit Approximations for Parallelism in Deep Learning. arXiv:1511.04561 [cs.NE]
- [17] Michael Diskin, Alexey Bukhtiyarov, Max Ryabinin, Lucile Saulnier, Anton Sinitin, Dmitry Popov, Dmitry V Pyrkin, Maxim Kashirin, Alexander Borzunov, Albert Villanova del Moral, et al. 2021. Distributed Deep Learning In Open Collaborations. *Advances in Neural Information Processing Systems* 34 (2021), 7879–7897.
- [18] O Elharrouss, Y Akbari, N Almaadeed, and S Al-Maadeed. [n.d.]. Backbone-review: Feature extraction network for deep learning and deep reinforcement learning approaches. arXiv 2022. *arXiv preprint arXiv:2206.08016* (n. d.).
- [19] Anne C Elster and Tor A Haugdahl. 2022. NVIDIA Hopper GPU and Grace CPU Highlights. *Computing in Science & Engineering* 24, 2 (2022), 95–100.
- [20] Wikimedia Foundation. 2023. "Wikimedia Downloads". <https://dumps.wikimedia.org>
- [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- [22] Kai Hwang. 1992. *Advanced Computer Architecture: Parallelism, Scalability, Programmability* (1st ed.). McGraw-Hill Higher Education.
- [23] Kyuhyong Lee and Myungjun Son. 2017. DeepSpotCloud: Leveraging Cross-Region GPU Spot Instances for Deep Learning. In *2017 IEEE 10th International Conference on Cloud Computing (CLOUD)*. 98–105. <https://doi.org/10.1109/CLOUD.2017.21>
- [24] Sungjae Lee, Jaeh Hwang, and Kyuhyong Lee. 2022. SpotLake: Diverse Spot Instance Dataset Archive Service. In *2022 IEEE International Symposium on Workload Characterization (IISWC)*. 242–255. <https://doi.org/10.1109/IISWC55918.2022.00029>
- [25] Mu Li, David G Andersen, Jun Woo Park, Alexander J Smola, Amr Ahmed, Vanja Josifovski, James Long, Eugene J Shekita, and Bor-Yiing Su. 2014. Scaling distributed machine learning with the parameter server. In *11th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 14)*. 583–598.
- [26] Shen Li, Yanli Zhao, Rohan Varma, Omkar Salpekar, Pieter Noordhuis, Teng Li, Adam Paszke, Jeff Smith, Brian Vaughan, Pritam Damania, et al. 2020. Pytorch distributed: Experiences on accelerating data parallel training. *arXiv preprint arXiv:2006.15704* (2020).
- [27] Xiangru Lian, Ce Zhang, Huan Zhang, Cho-Jui Hsieh, Wei Zhang, and Ji Liu. 2017. Can decentralized algorithms outperform centralized algorithms? a case study for decentralized parallel stochastic gradient descent. *Advances in neural information processing systems* 30 (2017).
- [28] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach. arXiv:1907.11692 [cs.CL]
- [29] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. 2022. A convnet for the 2020s. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 11976–11986.
- [30] Peter Mattson, Christine Cheng, Gregory Diamos, Cody Coleman, Paulius Micekevicius, David Patterson, Hanlin Tang, Gu-Yeon Wei, Peter Bailis, Victor Bittorf, et al. 2020. MLPerf Training Benchmark. *Proceedings of Machine Learning and Systems* 2 (2020), 336–349.
- [31] Petar Maymounkov and David Mazieres. 2002. Kademlia: A peer-to-peer information system based on the xor metric. In *Peer-to-Peer Systems: First International Workshop, IPTPS 2002 Cambridge, MA, USA, March 7–8, 2002 Revised Papers*. Springer, 53–65.
- [32] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems* 32 (2019).
- [33] Gustavo Portella, Genaina N Rodrigues, Eduardo Nakano, and Alba CMA Melo. 2019. Statistical analysis of Amazon EC2 cloud pricing models. *Concurrency and Computation: Practice and Experience* 31, 18 (2019), e4451.
- [34] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine Mcleavey, and Ilya Sutskever. 2023. Robust Speech Recognition via Large-Scale Weak Supervision. In *Proceedings of the 40th International Conference on Machine Learning (Proceedings of Machine Learning Research)*. Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (Eds.), Vol. 202. PMLR, 28492–28518. <https://proceedings.mlr.press/v202/radford23a.html>
- [35] Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. 2020. Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 3505–3506.
- [36] Jie Ren, Samyam Rajbhandari, Reza Yazdani Aminabadi, Olatunji Ruwase, Shuangyan Yang, Minjia Zhang, Dong Li, and Yuxiong He. 2021. ZeRO-Offload: Democratizing Billion-Scale Model Training. arXiv:2101.06840 [cs.DC]
- [37] Max Ryabinin, Tim Dettmers, Michael Diskin, and Alexander Borzunov. 2023. SWARM Parallelism: Training Large Models Can Be Surprisingly Communication-Efficient. *arXiv preprint arXiv:2301.11913* (2023).
- [38] Max Ryabinin, Eduard Gorbunov, Vsevolod Plokhhotnyuk, and Gennady Pekhimenko. 2021. Moshpit SGD: Communication-Efficient Decentralized Training on Heterogeneous Unreliable Devices. In *Advances in Neural Information Processing Systems*, Vol. 34. <https://proceedings.neurips.cc/paper/2021/file/97275a23ca44226c9964043c8462be96-Paper.pdf>
- [39] Learning@home team. 2020. Hivemind: a Library for Decentralized Deep Learning. <https://github.com/learning-at-home/hivemind>.
- [40] Chathurika S. Wickramasinghe, Daniel L. Marino, and Milos Manic. 2021. ResNet Autoencoders for Unsupervised Feature Learning From High-Dimensional Data: Deep Models Resistant to Performance Degradation. *IEEE Access* 9 (2021), 40511–40520. <https://doi.org/10.1109/ACCESS.2021.3064819>
- [41] Mitchell Wortsman, Tim Dettmers, Luke Zettlemoyer, Ari Morcos, Ali Farhadi, and Ludwig Schmidt. 2023. Stable and low-precision training for large-scale vision-language models. *arXiv preprint arXiv:2304.13013* (2023).
- [42] Sheng Yang, Samir Khuller, Sunav Choudhary, Subrata Mitra, and Kanak Mahadik. 2022. Scheduling ML Training on Unreliable Spot Instances. In *Proceedings of the 14th IEEE/ACM International Conference on Utility and Cloud Computing Companion (Leicester, United Kingdom) (UCC '21)*. Association for Computing Machinery, New York, NY, USA, Article 29, 8 pages. <https://doi.org/10.1145/3492323.3495594>
- [43] Zongheng Yang, Zhanghao Wu, Michael Luo, Wei-Lin Chiang, Romil Bhardwaj, Woosuk Kwon, Siyuan Zhuang, Frank Sifei Luan, Gautam Mittal, Scott Shenker, and Ion Stoica. 2023. SkyPilot: An Intercloud Broker for Sky Computing. In *20th USENIX Symposium on Networked Systems Design and Implementation (NSDI 23)*. USENIX Association, Boston, MA, 437–455. <https://www.usenix.org/conference/nsdi23/presentation/yang-zongheng>
- [44] Yang You, Jing Li, Sashank Reddi, Jonathan Hseu, Sanjiv Kumar, Srinadh Bhojanapalli, Xiaodan Song, James Demmel, Kurt Keutzer, and Cho-Jui Hsieh. 2019. Large batch optimization for deep learning: Training bert in 76 minutes. *arXiv preprint arXiv:1904.00962* (2019).
- [45] Binhang Yuan, Yongjun He, Jared Davis, Tianyi Zhang, Tri Dao, Beidi Chen, Percy S Liang, Christopher Re, and Ce Zhang. 2022. Decentralized training of foundation models in heterogeneous environments. *Advances in Neural Information Processing Systems* 35 (2022), 25464–25477.
- [46] Sergey Zagoruyko and Nikos Komodakis. 2016. Wide residual networks. *arXiv preprint arXiv:1605.07146* (2016).