

Efficient DNN Training with Knowledge-Guided Layer Freezing

Yiding Wang¹, Decang Sun¹, Kai Chen¹, Fan Lai², and Mosharaf Chowdhury²

¹Hong Kong University of Science and Technology

²University of Michigan

Abstract

Training deep neural networks (DNNs) is time-consuming. While most existing solutions try to overlap/schedule computation and communication for efficient training, this paper goes one step further by *skipping* computing and communication through DNN layer freezing. Our key insight is that the training progress of internal DNN layers differs significantly, and front layers often become well-trained much earlier than deep layers. To explore this, we first introduce the notion of *training plasticity* to quantify the training progress of internal DNN layers. Then we design KGT, a knowledge-guided DNN training system that employs semantic knowledge from a reference model to accurately evaluate individual layers' training plasticity and safely freeze the converged ones, saving their corresponding backward computation and communication. Our reference model is generated on the fly using quantization techniques and runs forward operations asynchronously on available CPUs to minimize the overhead. In addition, KGT caches the intermediate outputs of the frozen layers with prefetching to further skip the forward computation. Our implementation and testbed experiments with popular vision and language models show that KGT achieves 19%-43% training speedup w.r.t. the state-of-the-art without sacrificing accuracy.

1 Introduction

Recent advances in deep learning (DL) benefit significantly from training larger deep neural networks (DNNs) on larger datasets. Due to growing model sizes and large volumes of data, DNNs have become more computationally expensive to train, raising the cost and carbon emission of large-scale training [67]. Many recent DL research works focus on improving parallelism and pipelining via sophisticated computation-communication overlapping or scheduling to build more efficient systems and reduce training time [44, 62, 69]. Nevertheless, while approaching linear scalability can reduce the time to train a model, the total amount of computation requirement remains the same.

In this paper, we move one step further to explore: *can we reduce the total computation (and communication) in large DNN training?* We propose a *knowledge-guided* training system, KGT, to accelerate DNN training via computation-communication freezing while still maintaining accuracy. Our key insight is that the training progress of internal DNN layers differs significantly, and front layers can become well-trained much earlier than deep layers. This is because DNN features transition from being task-agnostic to task-specific from the first to the last layer [99, 100]. Thus, the front layers of a DNN often converge quickly, while the deep layers take a much longer time to train, as generally observed in both vision and language models [77, 82], experimentally validated in §2.3, and theoretically analyzed in Appendix C and D. KGT can safely freeze these converged DNN layers earlier, saving their corresponding computation and communication expenses without hurting model accuracy.

While freezing layers can reduce training cost, prematurely freezing under-trained layers may hurt the final accuracy. We observe that in transfer learning, freezing layers is mainly used for solving the overfitting problem [24]. While techniques such as static freezing [49] and cosine annealing [12] can reduce backward computation cost, accuracy loss is a common side effect. Thus, the main challenge of extending layer freezing to generic DNN training is how to maintain accuracy by only freezing the converged layers.

To address this challenge, KGT introduces the notion of *training plasticity*¹ to quantify a layer's training progress and safely detect the converged DNN layers to avoid premature freezing. To this end, KGT uses a *reference model*, which is the proxy for semantic knowledge, to evaluate DNN layers' plasticity. The reference model, in essence, is a trained lightweight DNN with the same architecture as the model being trained to understand layer-wise performance (details in §4.1). We compare the intermediate activations (internal

¹*Plasticity* quantifies a layer's training progress toward convergence, which is borrowed from *neuroplasticity* in neural science and child development [19]. Basically, a DNN layer's training plasticity will gradually decrease and become stable as it converges.

outputs) between the training model and the reference model elicited by the same data batch to measure the plasticity. When the plasticity becomes stationary, it implies that the layer is converged and can be frozen safely (§4.2). In addition, KGT can unfreeze the frozen layers to continue training with learning rate decay. Our approach is informed by recent advances in knowledge distillation research [2, 66, 86] that suggest the same input data (images and word vectors) will elicit similar *pair-wise* activation patterns in *trained models*. Compared to the naive metric of gradient’s granularity or norm (against a hard label), intermediate activation is more semantically meaningful, and thus accurate.

KGT adaptively generates the reference model by instantly compressing a snapshot of the training model via quantization [28] on CPUs after the *bootstrapping stage* [1] (early iterations during which the training model converges quickly). Large DNNs are robust to quantization, according to our evaluation and ML literature [51]. KGT also profiles in the background to make sure the CPU-efficient reference model can provide accurate plasticity evaluation. The reference model exploits available CPU cores during GPU-heavy training, running forward operations parallel to the GPU training using the same input data in a non-blocking and asynchronous fashion. Hence, the system overhead is minimal and can be well hidden. In the remaining training process, KGT updates the reference model using the latest snapshots to stabilize the plasticity curve. Essentially, we trade off small CPU resources for maintaining accuracy when freezing layers.

Freezing the front layers can save the backward computation and parameter synchronization. Nevertheless, we find that the forward pass still takes up to 35% of the time of an iteration. We observe that, in DNN training, the frozen front layers will produce the same forward output given the same input. Prior work on inference also shows that caching forward results can improve performance [9]. To take advantage of this, KGT saves the intermediate activations of the frozen layers to the disk, *prefetches* the saved activation tensors to the GPU memory, and continues training the remaining layers from the cached activations in the following epochs. As the data loader knows the future data sequence, it is possible to prefetch relevant activations without stalling. Caching and prefetching are also compatible with data augmentation. Therefore, we further save the frozen layers’ forward computation without altering the training data sequence (§4.3).

We implement KGT as a framework-independent Python library (§5). Existing code can work with KGT with minimal changes. We evaluate KGT using seven popular vision and language models on five datasets (§6). It achieves 19%-43% training speedup than the state-of-the-art frameworks and can reach the target accuracy.

To summarize, the key contributions of KGT include: (1) leveraging semantic knowledge to save the backward computation and communication via DNN layer freezing while maintaining accuracy; (2) building an efficient system

to implement the idea of knowledge-guided training; and (3) caching the intermediate results with prefetching to further save the forward pass of the frozen layers with negligible overheads. KGT will be open-source.

2 Background and Motivation

2.1 DNN Training

Modern DNNs consist of dozens or hundreds of layers that conduct mathematical operations. Each layer takes an input tensor of features and outputs corresponding activations. We train a DNN by iterating over a large dataset many times and minimizing a loss function. The dataset is partitioned into *mini-batches*, and a pass through the full dataset is called an *epoch*. A DNN training iteration includes three steps: (1) forward pass, (2) backward pass, and (3) parameter synchronization. The forward and backward passes require GPU computation. In each iteration, the *forward pass* (FP) takes a mini-batch and goes through the model layer-by-layer to calculate the loss regarding the target labels and the loss function. In the *backward pass* (BP), we calculate the parameter gradients from the last layer to the first layer based on the chain rule of derivatives regarding the loss [58]. At the end of each iteration, we update the model parameters with an optimization algorithm, such as stochastic gradient descent (SGD) [11]. In data parallel distributed training, independently computed gradients from all workers are aggregated over the network to update the shared model.

2.2 Existing Optimizations for DNN Training

Training large DNNs is computation- (and communication-) intensive due to the ever-growing data volumes and model size [41, 62, 69]. One important direction for DNN training acceleration from the system perspective that is closely related to KGT is *computation-communication overlapping and scheduling*. Baseline training frameworks (e.g., TensorFlow, PyTorch, and Poseidon [102]) optimize distributed performance by issuing the gradient transmission once a layer finishes its backward computation so that the deeper layers can overlap their communication with the front layer’s BP. Priority-based communication scheduling systems (e.g., ByteScheduler [69], P3 [41] and TicTac [29]) leverage the layer-wise structure information to prioritize the front layers in communication which try to overlap the communication with FP. Pipelining solutions [38, 62, 97] add inter-batch pipelining to intra-batch parallelism to further improve parallel training throughput. While all these solutions can optimize computation-communication efficiency, the total computational cost remains the same.

Other methods. There are other optimizations, such as gradient sparsification [54] and quantization [91], to reduce com-

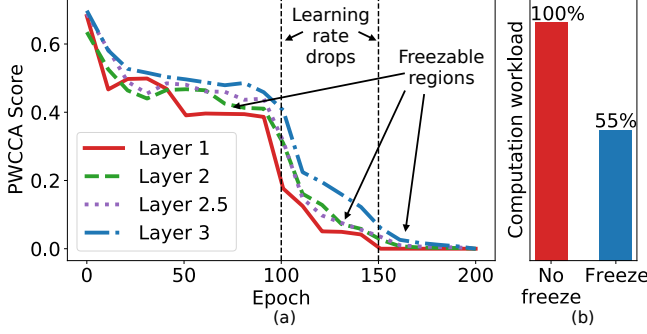


Figure 1: Post hoc layer convergence analysis with PWCCA. A lower score means the layer converges better. Layer 2.5 is the first half of Layer 3. The learning rate changes at 100th and 150th epochs and reboots training. We can freeze layers when they are stationary and unfreeze them when the learning rate decreases extremely.

munication volumes. These methods are largely orthogonal to KGT, and we will overview them in Section 7.

2.3 Opportunities of DNN Layer Freezing

In this paper, we explore the idea of *reducing computation and communication costs* through DNN layer freezing. In the following, we first show the idea and its potential, and then lay out the challenges, motivating the design of KGT.

Motivation for layer freezing. Recent efforts have shown that the front layers primarily extract general features of the raw data (e.g., the shape of objects in an image) and often become well-trained much earlier [75, 99], while deeper layers are more task-specific and capture complicated features output from front layers (details in Appendix A). Our work is also inspired by transfer learning [36, 46, 55, 73]. When fine-tuning a pre-trained model on a new task, we find that ML practitioners can freeze (i.e., fix layer’s weights) the front layers or only fine-tune them for a few iterations and focus on training the deep layers on the new dataset.

To demonstrate the potential, we use PWCCA [61], a *post hoc* layer convergence analysis tool, to track the training progress of different layers of ResNet-56 [31] as an example. ResNet-56 is a popular model for image classification on the CIFAR-10 dataset, and it consists of three layer blocks (referred to as modules or stages), where each module has 18 basic blocks of successive layers. PWCCA compares the intermediate activation (i.e., the output feature map produced by a DNN layer) with a *fully-trained* model. A low PWCCA score (0–1 range) suggests the layer is converged to the final state. We can clearly find some freezable regions in Figure 1: e.g., for layer 1, during the 50th–90th, 120th–140th, and after the 160th epoch, its score becomes stable, meaning it is temporarily converged; other layers also show relatively

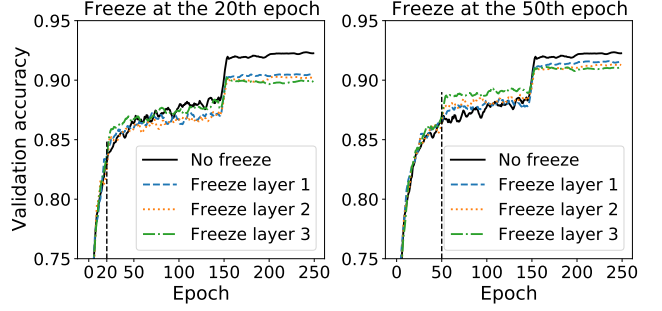


Figure 2: Prematurely freezing layers can hurt the final accuracy, especially for the deep layers.

stable regions. The scores drop at the 100th and 150th epochs because the learning rate decreases as scheduled; after that, they soon converge again. This pattern reveals the strategy during training: *freeze the layers when their performance is stable and unfreeze them when the learning rate decreases*.

We then freeze the layers during their freezable regions and find that we can reduce the computation and communication costs by 45% in theory without hurting accuracy! For natural language processing (NLP) models, this potential can be even larger because the front layers usually contain more parameters than CNNs.

Challenges in layer freezing. Although this opportunity has been pointed out by ML community, it is not feasible to run the hypothetical post hoc analysis (e.g., PWCCA) in practice. Quantifying the training progress of a layer is difficult due to the lack of prior knowledge (e.g., a trained model). Furthermore, we find that prematurely freezing DNN layers can greatly hurt a model’s accuracy. To demonstrate this, we investigate the impact of static freezing on the final accuracy when training ResNet-56. In Figure 2, we fix the parameters of a layer module at the 20th and 50th epoch in each figure and show their validation accuracies alongside the original baseline accuracy. The final degraded accuracies indicate that freezing layers prematurely can hurt the converged model’s performance. For deep layers, this accuracy degradation can be more than 2% which is huge for such models.

3 KGT Overview

We propose KGT, an efficient DNN training system that detects and freezes the converged layers in a practical manner. Lacking prior knowledge of the hypothetical fully-trained model, KGT introduces a self-generated *reference model* during training to provide semantic knowledge for evaluating a layer’s convergence with minimal system overhead (§4.1). The reference model is essentially an accompanying lightweight DNN with the same architecture as the model being trained to capture their internal layers and understand

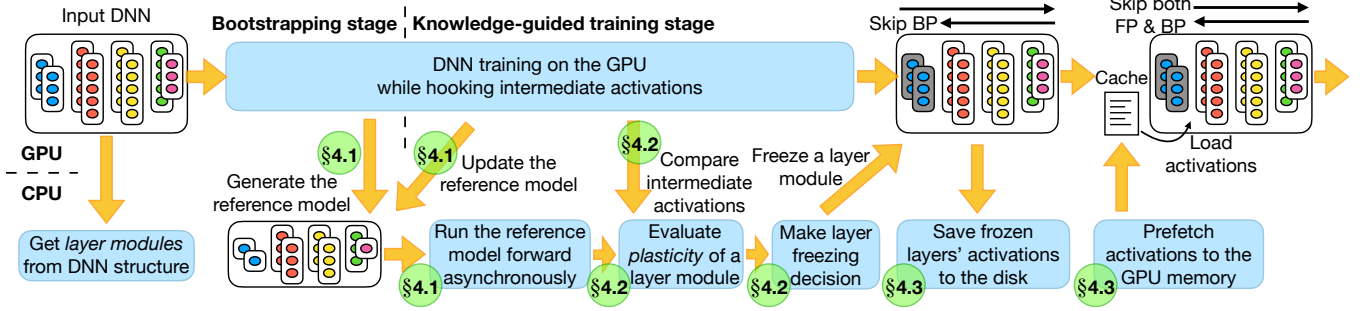


Figure 3: KGT overview. KGT covers two training stages (bootstrapping and knowledge-guided stages) and has three major design components (generating and executing the reference model, layer freezing with plasticity evaluation, and skipping FP with prefetching).

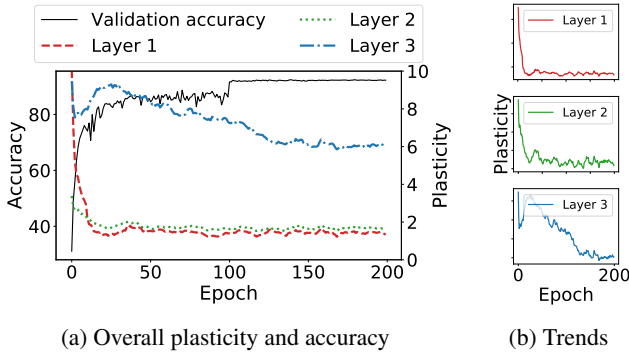


Figure 4: The plasticity of the front layers drops quickly, and they will produce semantically similar activations for most training iterations.

layer-wise performance.

To quantify the training progress, KGT defines a system metric of *plasticity*. A layer’s plasticity is formulated as the difference between the intermediate activation tensors of the training model and its reference model given the same mini-batch input.² The plasticity changes as the model evolves over training, and KGT considers the layers with stable plasticity values to be converged, whereby KGT freezes these layers without hurting the accuracy (§4.2).

To validate the effectiveness of plasticity in capturing the training progress, we use ResNet-56 and generate a reference model with the same architecture but pre-trained for only 50 epochs. We measure the plasticity of ResNet-56’s first three layer modules during training. In Figure 4a, the top black curve shows the validation accuracy, and the other three indicate the derived plasticity for each layer module. We find that, in the first ~30 epochs, the plasticity of the first two modules

²We measure the difference using the Similarity-Preserving loss (SP loss) [86], a novel loss function recently developed to compare two activation tensors for CV tasks and also echoed in NLP [66]. Unlike PWCCA [61] for post hoc convergence analysis, SP loss focuses on capturing the semantic difference for DNN training, making it a perfect fit for plasticity evaluation (detailed analysis in Appendix D).

converges quickly to a low level while the plasticity of layer module 3 is much higher and unstable. These layers show different trends of plasticity, more clearly after normalization in Figure 4b. We find similar freezable regions when using the post hoc analysis (Figure 1), e.g., layer 1 converges near the 50th epoch, while layer 3 only converges at the last epochs. Plasticity requires no prior knowledge as PWCCA, only a training-in-progress model for reference (more comparison in Appendix D), and accurately captures the trend of layer convergence. The algorithm behind plasticity is tested performant for various tasks compared to traditional gradients and other activation-based metrics [66, 86].

Training life cycle with KGT. Figure 3 describes the high-level workflow of KGT in two stages.

(1) *Bootstrapping stage*: when a job is submitted, KGT starts to monitor the job. The bootstrapping stage is a *critical period* of training, during which the DNN is sensitive and no parameter is eligible for freezing, according to recent ML research [1]. KGT monitors the changing rate of the training loss (in line with the later plasticity monitoring) and moves to the next stage as the DNN moves out of the critical period.

(2) *Knowledge-guided training stage*: KGT generates the reference model on the CPU using the latest snapshot of the training model. Afterwards, KGT collects the intermediate activation tensors of the frontmost non-frozen layer of the full model and its reference model for plasticity evaluation (§4.1), freezes the layer once it reaches the convergence criteria (§4.2), and then moves to the next active layer. KGT excludes the frozen layers during BP (and parameter synchronization in case of distributed training) to accelerate training. Meanwhile, KGT caches the frozen layer’s activations to the disk, so that we can also skip the FP computation by prefetching the intermediate results for the same input (§4.3).

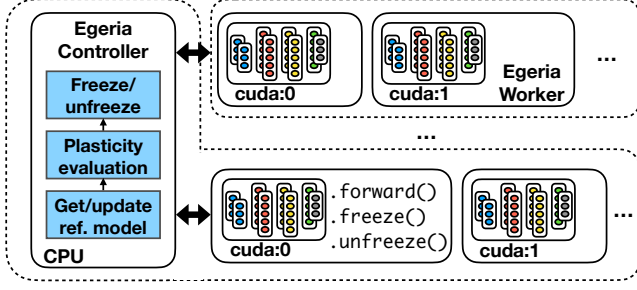


Figure 5: KGT uses a controller-worker framework. The controller co-locates on a training node. Solid lines denote device and logical boundaries; dashed lines denote machines.

4 Design

Next we dive into the details on how to capture the semantic knowledge during training (§4.1), with which KGT optimizes the computation and communication in the backward pass (§4.2), as well as the forward pass (§4.3) on the fly.

4.1 KGT Architecture

Directly running another full DNN to measure the internal layers’ plasticity can greatly slow down the training. Instead, KGT decouples the control logic and the training logic with a controller-worker abstraction (§4.1.1), and asynchronously performs plasticity evaluation (§4.1.2). Besides, KGT generates and continuously updates the reference model by fast quantization (§4.1.3).

4.1.1 Controller-Worker Framework

Figure 5 illustrates the controller-worker framework of KGT, which primarily consists of a logically centralized controller and workers:

- **Controller:** It manages the life cycle of the reference model, including its generation and execution, gathering data for plasticity evaluation, and making layer freezing/unfreezing decisions for workers. It makes plasticity evaluations at one place to reduce the overall computation overhead in case of distributed training, as multiple controllers only change the sample size.
- **Worker:** Each training worker has an KGT worker process. In addition to the original training operations, it performs KGT tasks, including transmitting data and handling controller decisions. The new `forward()` method uses hooks to obtain the intermediate activation tensors. The `freeze()` and `unfreeze()` methods will be called by the controller and apply on target layers.

KGT can wrap the model as a worker for single-GPU (e.g., `nn.Module` in PyTorch), multi-GPU-single-node

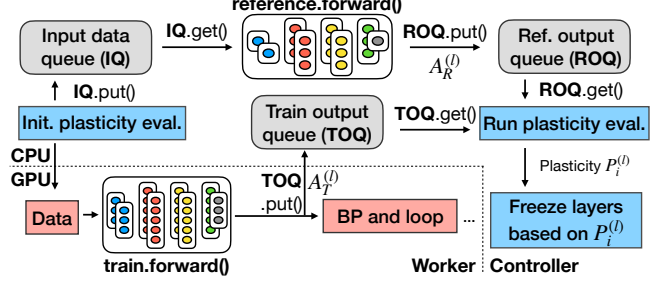


Figure 6: Red and blue blocks are worker and controller operations, respectively. They interact asynchronously through multiprocessing queues (gray blocks).

(e.g., `nn.DataParallel`), or multi-GPU-multi-node (e.g., `nn.parallel.DistributedDataParallel`) training. KGT requires minimal lines of code changes:

```
from egeria import EgeriaController, EgeriaModule

controller = EgeriaController(args, ...)
# Replace nn.Module for local training
model = EgeriaModule(arch, args, ...)
```

4.1.2 Non-Blocking Plasticity Evaluation

The controller runs the reference model on CPUs in a non-blocking and asynchronous fashion to minimize system overheads. GPU machines typically have high-end CPUs [4] which are optimized for int8 DNN inference operations [39]. During training, CPUs handle the data preparation tasks which take relatively few CPU cycles. It is widely observed that there are spare CPU cycles in most cases because GPUs are often the bottleneck in DNN training [42, 44].

To avoid slowing down the training, the controller runs an efficient reference model in parallel on CPUs, as shown in Figure 6. Note that the reference model has the same architecture as the training model, so KGT can automatically parse their building blocks and match the internal layers. We implement the asynchronous plasticity evaluation using three *single-producer/single-consumer queues*: (1) When KGT initiates a plasticity evaluation, the co-located KGT worker puts the data batch in the input queue (IQ). (2a) The controller process polls IQ, runs a forward pass on the reference model, and puts the hooked intermediate activation $A_R^{(l)}$ to the reference output queue (ROQ). The controller only executes the forward pass at low CPU load (<50%) to avoid interference to other CPU-based auxiliary operations, thanks to the non-blocking framework. (2b) The co-located KGT worker puts the hooked intermediate activation $A_T^{(l)}$ to the training output queue (TOQ) and continues the training loop without blocking. (3) The controller polls ROQ and TOQ and calculates the plasticity of the frontmost active layer modules to make freezing decisions (§4.2).

4.1.3 Generating and Updating the Reference Model

An ideal reference model should execute fast for diverse models and provide semantically meaningful activations. To this end, KGT generates the reference model for the full model being trained on the fly.

There are several techniques to generate a light-weighted version for a large model. For example, neural architecture search (NAS) [106] and knowledge distillation (KD) [33] can compress the model but have prohibitively large computation overhead. Besides, they may produce different architectures that do not match the internal layers for plasticity evaluation. Therefore, KGT adopts *post-training quantization* [28] to instantly generate a reference model with the same architecture.

Quantization is a popular model compression technique to accelerate inference on CPUs [37, 93, 96]. It reduces the precision of model’s parameters (e.g., from 32-bit floating-point to 8-bit integers). By default, KGT quantizes the reference model using 8-bit integers. This can reduce the reference memory footprint by $3\times$ to $4\times$ and accelerate the forward pass by $2\times$ on CPUs; meanwhile, a lower precision (e.g., int4 or int2) cannot further improve the performance due to the CPU instruction set [53]. In our test, int8 post-training quantization reaches a sweet spot to provide both fast and accurate semantic reference (§6.4). KGT can fall back to full-precision reference if the training DNN is extremely sensitive or running the reference on GPU in case the CPU resources are scarce, adapting to various environments. In addition, we design the freezing workflow (Algorithm 1) with robustness kept in mind.

KGT will periodically update the reference model using the latest training snapshot to keep it up-to-date for plasticity evaluation. We observe that a stale reference model can amplify the inherent fluctuations in stochastic gradient descent (SGD) [11] training, making it hard for KGT to understand plasticity trends (§4.2). To illustrate this phenomenon, we show the plasticity trace of the aforementioned ResNet model *without* reference update in Figure 7a, wherein a stale reference model introduces large fluctuations that may disrupt plasticity evaluation. Figure 7b shows that updating the reference model can better reflect the performance trend.

4.2 Freezing Layers with Plasticity

Next, we elaborate how KGT decides when to freeze the stable layers by comparing the intermediate activations of the model being trained with that of the reference model. During this stage, KGT needs to address two challenges: (1) how to quantify the training plasticity of a layer module; and (2) how to accurately make the layer freezing decision.

4.2.1 Evaluating the Plasticity of a Layer Module

The raw data we collect are the intermediate activation tensors of the training and reference models, which has been widely

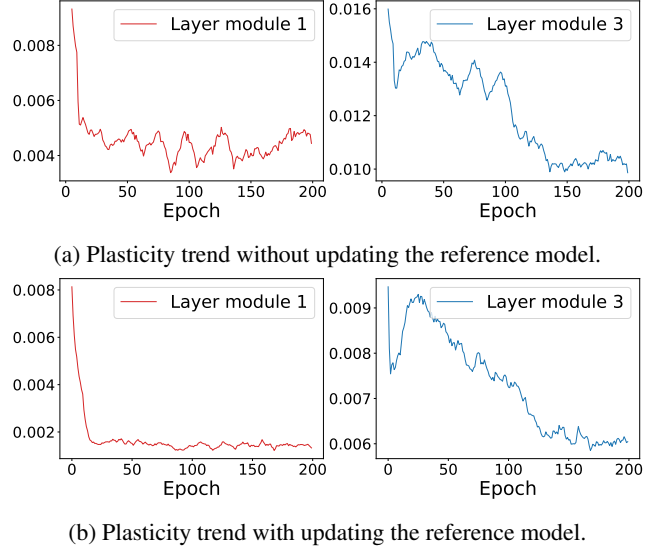


Figure 7: The controller will update the reference model to ease the fluctuation of plasticity shown above.

explored as a direct metric of layer performance in DNN training [61]. For example, knowledge distillation uses the difference of activations between the training model and a trained teacher model [78, 86] as a supervisory signal to improve accuracy, because intermediate activation plays an important role in forming the decision boundaries for the partitioning of the feature space in each hidden layer [32, 65]. Using intermediate activation as a “soft” label is proven successful compared to only using the gradients calculated against a hard label [2, 86], because it is a more semantically meaningful indicator and provides *contextual* knowledge [18, 66]. To accurately freeze layers, we quantify the training plasticity by measuring the changes of the layer’s intermediate activations.

We use the SP loss [86] that shows high efficacy in ML tasks to compare the intermediate tensors between two models. The theory behind SP loss is that the same input data will elicit similar pair-wise activation patterns in trained models for both CV and NLP tasks [66, 86]. The same idea also applies in our training plasticity evaluation. Compared to PWCCA [61] which is designed for post hoc analysis, SP loss reflects layers’ actual training progress rather than for visualization (details in Appendix D).

We focus on the semantic performance of a layer module that usually contains a consecutive of layers defined together. Layers in a layer module are closely related to perform a sequence of transformations for a certain goal [50] and have similar training progress. Meanwhile, individual layers with fewer parameters (e.g. linear layers) are less stable in SGD training. Though it’s rare, even a few individual front layers might not converge in a strict order (as observed in [101]), our module-based freezing can mitigate this (Appendix D) and revisit them in the future unfreezing stage (§4.2.2). KGT provides configuration options to customize the granularity

of layer module through regular expression, e.g., evaluating every convolutional layer.

Given the input data of batch size b , we denote the activation tensors of the training and reference models at a layer l as $A_T^{(l)}$ and $A_R^{(l)}$. For the image data, the activation tensors $A_T^{(l)}, A_R^{(l)} \in \mathbb{R}^{b \times c \times h \times w}$, where c , h , and w are channel number, height, and width; similar for the word embeddings. Then SP loss will align $A_T^{(l)}$ and $A_R^{(l)}$ to $b \times b$ -shaped matrices, which encode the pair-wise similarity in the activation tensors that are elicited by the input mini-batch (details in Appendix B). We then denote the *training plasticity* $P_i^{(l)}$ of layer l at an iteration i using the SP loss between the two matrices, representing the semantic difference compared to the reference model, as shown in Equation 1. The lower and more stable the plasticity, the DNN layers are closer to convergence.

$$P_i^{(l)} = SP_Loss(A_T^{(l)}, A_R^{(l)}) \quad (1)$$

4.2.2 How to Decide Layer Freezing

During the knowledge-guided training stage, KGT will periodically run the plasticity evaluation every n iterations and decide to freeze the layer or not. The intuition of the freezing criterion is straightforward: if a layer’s plasticity becomes stationary for *some iterations* S , KGT considers its semantic performance stable and can safely freeze it.

When obtaining the plasticity $P_i^{(l)}$, we first smooth it with the moving average of its recent values (history buffer window W or the max span if it is smaller), as shown in Equation 2.

$$\overline{P_i^{(l)}} = \begin{cases} \frac{P_{i-W}^{(l)} + \dots + P_i^{(l)}}{W}, & i \geq W \\ \frac{P_0^{(l)} + \dots + P_i^{(l)}}{i}, & i < W \end{cases} \quad (2)$$

To determine whether the curve has become stable, we fit $\overline{P_i^{(l)}}$ with linear least-squares regression to a straight line and analyze its slope (0 means no change at all). This method can filter out the drastic fluctuation in SGD training and provide a recent history context than simply evaluating delta. If the plasticity slope has been less than the tolerance T for S evaluations ($S = W$ by default), we consider the layer converged, freeze it, and move to the next layer module, as detailed in Algorithm 1. This simple yet effective method has a similar intuition to the early stopping in DNN training [45, 81].

KGT monitors the frontmost active layer module l to avoid a fragmented frozen model. According to the chain rule of automatic differentiation [58], only excluding the last link of backpropagation can reduce the workload. It is widely recognized that the front layers converge faster [77, 82, 99, 100], and KGT can handle exceptions with the aforementioned module-based freezing and unfreezing mechanism.

Unfreezing. Learning rate (LR) scheduling can largely influence the convergence [1] (e.g., the accuracy boost in Fig-

ure 2) and is an external factor to the model. LR annealing [23], the most commonly employed scheduling technique, recommends gradually lowering the LR during training with implementations like step decay and exponential decay. Given this, KGT will restart training all the frozen layers if the LR has dropped over a factor of 10 since the frontmost layers’ frozen and halve the counter and history buffer for refreezing. Another type of LR scheduling is periodically increasing and decreasing the LR, e.g., cosine annealing [56] and cyclical LR [83]. Due to its complexity, KGT lets the user customize the unfreezing and refreezing criteria, e.g., training for a few iterations in each cycle after freezing a layer.

Hyperparameters. We use four hyperparameters: n (plasticity evaluation and bootstrapping stage monitoring interval), W (history buffer to provide recent context), S (stale counter threshold, same value as W), and T (plasticity slope tolerance). They seem overwhelming but are highly related and require just little empirical knowledge to set. Our extensive empirical analyses show we can achieve consistently good performance across different parameters, similar to setting the batch size and tuning the early stopping criteria [45, 81] (sensitivity analysis of W in §6.4). The changing rate of ending the bootstrapping stage is permissively set to 10%.

4.3 Skipping Forward Pass with Caching and Prefetching

By freezing the converged front layers, we can exclude them from the backward pass and parameter update to reduce training cost. However, the forward pass is still necessary because the deep layers require the frozen layers’ activations as input [58]. Naturally, we can cache the frozen layers’ intermediate activations to save the forward pass because they output the same activation given a certain input.

There are two challenges of caching computation results for a DNN training task. First, training a large model requires a large dataset (e.g., the training set of ImageNet is over 100 GB). The size of the intermediate activation tensor depends on the output shape of the last frozen layer. In our evaluation, the storage space of ResNet-50 intermediate activations ranges from $1.5 \times$ to $5.3 \times$ of the input data. It is not technically appropriate to cache a whole epoch’s results to the GPU/CPU memory. Second, given the memory limit, caching systems usually improve their hit rate by keeping the most frequent content using replacement policies like LRU (least-recently-used). However, in DNN training, the data loader randomly samples a mini-batch, meaning there is no popular data to prioritize in the cache.

To solve these challenges, we exploit a training workflow feature: before an iteration, the *data loader* samples future mini-batches in advance, so unlike typical cache systems, we actually “know the future” (the incoming data indices)! Prefetching is an effective technique in ML applications [8,

Algorithm 1: Layer freezing algorithm.

Input:

Intermediate activations of the training and reference models $A_T^{(l)}$, $A_R^{(l)}$, layer module l , training iteration i , history buffer W , tolerance T , staleness threshold S .

Output: The (updated) frontmost active layer l .

/ Initialize global variables. */*

```

1   $pList_l \leftarrow \emptyset$ ;  $\triangleright$  Plasticity evaluation history of  $l$  across iterations.
2   $staleCounter \leftarrow 0$ ;  $\triangleright$  Number of consecutive stale  $P_i^{(l)}$ .
3  Function checkPlasticity( $A_T^{(l)}$ ,  $A_R^{(l)}$ ,  $l$ ,  $i$ ,  $T$ ,  $S$ ):
4      assert  $l$  is not the last layer
5      if  $staleCounter < S$  then
6           $P_i^{(l)} \leftarrow \text{calculateSPLoss}(A_T^{(l)}, A_R^{(l)}, l, i)$ 
7          /* Use moving average to mitigate outliers (Equation 2). */
8           $\overline{P}_i^{(l)} \leftarrow \text{smoothPlasticity}(P_i^{(l)}, W)$ 
9          /* Update the time-series plasticity list. */
10          $pList_l \leftarrow pList_l \cup \overline{P}_i^{(l)}$ 
11         /* Calculate the slope of the linear-fitted plasticity. */
12          $s \leftarrow \text{windowLinearFit}(pList_l, W).slope$ 
13         /* If the fitting line is close to horizontal. */
14         if  $s < T$  then
15              $staleCounter \leftarrow staleCounter + 1$ 
16         else
17              $staleCounter \leftarrow 0$ 
18         end if
19     else
20         freezeLayer( $l$ )
21         /* Move to the next active layer. */
22          $l \leftarrow l + 1$ 
23     end if
24     /* Learning rate-based unfreezing mechanism. */
25     if LR annealing and LR decreased by 90% then
26         unfreezeAllLayers()
27          $l \leftarrow l_0$ ;  $\triangleright$  Reset  $l$ .
28     else
29         if cyclical LR scheduling then
30             customizedUnfreeze()
31         end if
32     end if
33     return  $l$ 

```

[71]. KGT saves the forward computation results of the frozen layers to the disk and prefetches the relevant activations to the GPU memory so that the active deep layers can instantly read them as input. The cache only stores the recent five mini-batches for minimal memory usage. Users can set the storage limit for activations that are up to an epoch (see analysis in §6.5). At the early training stage, we disable prefetching if the forward pass of a few layers is faster. In this way, we can skip the forward computation of the frozen layers and efficiently

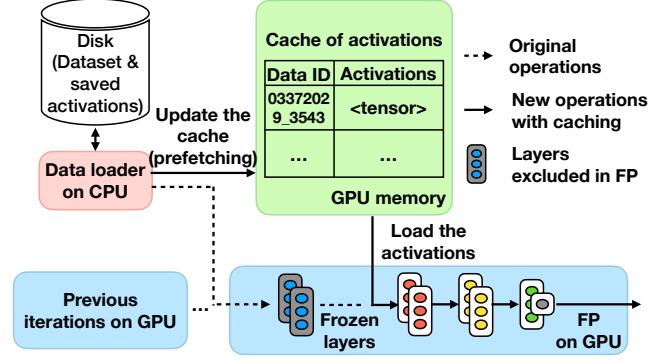


Figure 8: KGT caches the intermediate activation and prefetches the tensors into GPU memory during the FP.

overlap the slower disk access with prefetching.

Figure 8 illustrates a forward pass in KGT. We maintain a hash table in the GPU memory. The key is the sample ID, and the value is the corresponding activation tensor. During training, the data loader will sample a data batch and prefetch their activations from the storage to the GPU memory in parallel to the GPU training. KGT prefetches more than one mini-batch of the future activations, similar to the data loader [71], depending on the memory and CPU availability.

KGT can cover different training techniques and cache their outputs. KGT is compatible with stateless random operations, which are recommended [85] for random data augmentation. Thus we can deterministically keep the augmented images the same across epochs. Each worker machine maintains its own cache in distributed training to avoid extra network overhead; besides, the random seed is device-dependent. For most layers, e.g., convolutional (for CNNs) and linear (for language models), the output activation only depends on the parameters given the same data, so caching can work naturally. For batch normalization layers, the activation also depends on the specific data batch. KGT handles this case using the practice in transfer learning [46]: we set these layers to the inference mode, using the dataset statistics to normalize the input rather than the specific batch.

5 Implementation

KGT is independent of DNN training frameworks. In this paper, we implement and evaluate KGT using PyTorch and Huggingface Transformers [92]. All the technical dependencies of KGT (e.g., quantization and asynchronous computation) can work in other ML frameworks like TensorFlow and MXNet. KGT requires minimal lines of code changes, e.g., replacing `nn.Module` with our `EgeriaModule` class and adding an `EgeriaController`. KGT obtains the layer modules by parsing the model definition and adds forward hooks to obtain intermediate activations.

Reference model. When KGT controller generates or updates the reference model, it directly moves a snapshot of the training model from GPU to CPU and runs int8 quantization using PyTorch’s built-in library. We use dynamic quantization for NLP models and static quantization for convolutional networks, which add little overhead in the background. We add the same forward hooks to the reference model to match the training model.

Knowledge-guided training. Our high-level API allows us to take advantage of the framework engines to execute all the DNN computation operations. To freeze a layer, we essentially set the `requires_grad` flag of all its parameters to false to exclude the subgraph from gradient computation [70]. Distributed training requires rebuilding the communication buffer. For caching, we use the dictionary data structure, an efficient implementation of hash table [64] with $O(1)$ lookups.

6 Evaluation

In this section, we evaluate the effectiveness of KGT, namely accelerating DNN training while maintaining accuracy, for different tasks and models using single or multiple machines. The main takeaways are the following:

- KGT can work for different CV and NLP models;
- KGT accelerates training by 19%-43% without hurting accuracy; and
- KGT minimizes the system overhead while accurately freezing layers.

6.1 Methodology

Testbed setup. We evaluate KGT using two testbed configurations: a cluster of 5 machines and a single physical machine with multiple GPUs. In the 5-node cluster, each machine has 2 NVIDIA V100 GPUs, 40 CPU cores, 128 GB memory, and 2 Mellanox CX-5 NICs. The testbed has a leaf-spine topology with two core and two top-of-rack (ToR) Mellanox SN2100 switches; each ToR switch is connected to 5 servers using 40 Gbps and two core switches using 100 Gbps links. The single node has 8 NVIDIA GeForce RTX 2080 Ti and 64 CPU cores.

Tasks, models, and datasets. We use two CV and two NLP tasks in our evaluation: image classification, semantic segmentation, machine translation, and question answering; the corresponding 7 models and 5 datasets are listed in Table 1.

Training configurations. We follow the popular batch size setting of each model [63,72], and use the recommended learning rates and learning rate schedulers: 0.1 for image classification, 0.01 for semantic segmentation, 0.001 for Transformer,

and $3e-5$ for fine-tuning BERT; step decay learning rate schedule for CV training, inverse square root schedule for Transformer training, and linear schedule for fine-tuning BERT. We use the all-reduce parameter synchronization scheme for data parallel distributed training with multiple GPUs or machines and allocate one GPU per process.

Baselines and metrics. We compare KGT with the vanilla training framework (PyTorch) and a pure communication scheduling system, ByteScheduler [69], in multi-node distributed training. ByteScheduler achieves the theoretically optimal scheduling without skipping any parameter. We do not include traditional model fine-tuning techniques like static layer freezing because they would hurt the final accuracy, as presented in Section 2.3. In all evaluations, the training performance metric is the time taken to a converged validation accuracy (TTA), as listed in Table 1.

In multi-node distributed training, we also show the average training throughput, i.e., the number of samples processed per second. We use ByteScheduler’s default configuration.

6.2 End-to-End Training Performance

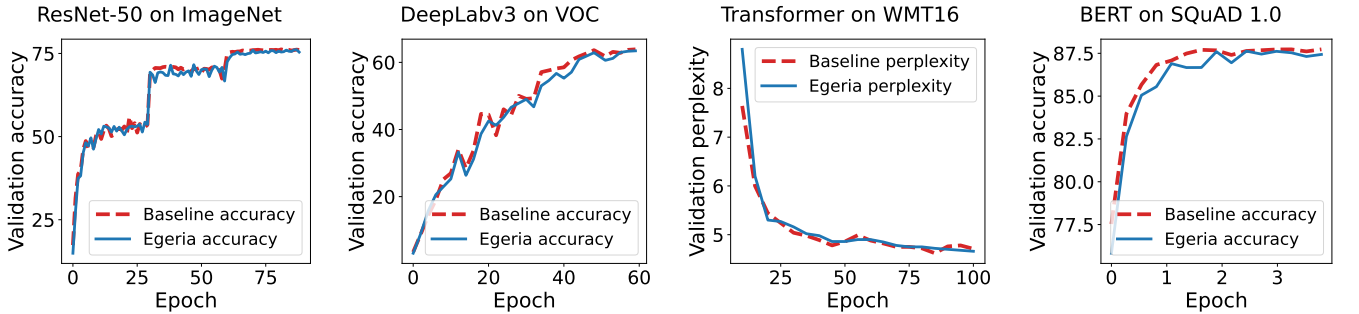
We use KGT to train different models to reach the target accuracies with largely reduced training time. Table 1 summarizes the evaluation results and the time-to-accuracy (TTA) speedups compared to the baseline DNN training system. KGT parses the model definition to extract the basic building blocks.

Image classification. ResNet-50 for ImageNet is a popular CNN benchmarking model. It consists of 48 building blocks of layers, grouped into four stages, and the deep layers have much more parameters than the front layers. Figure 9a shows the validation accuracy curves of KGT and the baseline. Within 90 epochs of training, both reach the target accuracy. With layer freezing, KGT shows a very close validation accuracy compared to that of the baseline. At the 30th and 60th epochs, the step-decay learning rate will improve the accuracy. During these critical stages, the unfreezing mechanism of KGT (§4.2.2) restarts the frozen layers and achieves the same level of accuracy boost. While KGT achieves the same accuracy, because of judicious freezing, it accelerates training time by 28% (detailed TTA plots can be found in Appendix E). The performance improvement primarily comes from later training stages when KGT freezes the deeper layer modules with more parameters. KGT can accelerate lightweight models (e.g., MobileNet V2) on smaller datasets (e.g., CIFAR-10) with 22% and 23% speedups. We also test freezing layers based on gradient norm on CIFAR-10 and find that achieving the same speedup will lose 2% of accuracy.

Semantic segmentation. We use the DeepLabv3 model with a ResNet-50 backbone for semantic segmentation train-

Task	Model	Dataset	Accuracy target	# Servers \times # GPUs/server	# Building layer modules	TTA speedup
Image classification	ResNet-50 [31]	ImageNet [20]	Top 1	1×2	48 (residual blocks)	28%
	MobileNet V2 [79]		75.9%	$2 \times 2 - 5 \times 2$	17 (inverted residual blocks)	27%-33%
	ResNet-56 [31]	CIFAR-10 [47]	71.2%	1×2	54 (residual blocks)	22%
Semantic segmentation	DeepLabv3 [16]	VOC [22]	mIoU 63.3%	1×2	49 (residual blocks and DeepLab head)	23%
Machine translation	Transformer-Base [87]	WMT16	Perplexity	4×2	12 (6 encoders & 6 decoders)	43%
	Transformer-Tiny	EN-DE [10]	4.7	$2 \times 2 - 5 \times 2$	4 (2 & 2)	33%-43%
Question answering	BERT-Base [21] (fine-tuning)	SQuAD 1.0 [76]	F1 score 87.6	1×2	12 (Transformer blocks)	19%

Table 1: Summary of evaluation tasks. Accuracy targets are the converged accuracy in baseline training. KGT accelerates different models by 19%-43% to reach the target accuracy.



(a) Epoch-to-accuracy of training ResNet-50. (b) DeepLabv3 training achieves a 21% speedup. (c) Transformer-Base for machine translation. (d) Fine-tuning BERT achieves a 41% speedup.

Figure 9: KGT can accelerate time-to-accuracy for different tasks without sacrificing training accuracy.

ing. The structure of DeepLabV3 includes a backbone module for feature computation and extraction plus a classifier module that takes the output of the backbone and returns a dense prediction. KGT will parse the backbone same as the ResNet-50 training and consider the classifier module as a whole. DeepLabv3 uses a Lambda learning rate scheduler that changes along with the training procedure, which will trigger the unfreezing mechanism of KGT at the 45th epoch. Figure 9b shows that, compared to the baseline, KGT can reach the target accuracy (mIoU of 63.3%) 21% faster and quickly improve accuracy at the later training stage when the learning rate scheduler significantly decreases.

Machine translation. KGT not only works for CV models but also for language models. A low perplexity means high accuracy for translation tasks. In Figure 9c, the model

quickly reaches a low level of perplexity then continues to improve slowly. KGT brings a 43% speedup by freezing the front encoders. Unlike CNN models that usually have heavy deep layers, Transformer has a more balanced structure, so skipping front layers can bring a considerable speedup. We also evaluate a tiny version of Transformer using an 8-GPU machine and achieve a 19% speedup.

Question answering. Training a question answering model is different from the other tasks because we fine-tune a pre-trained general-purpose language model BERT for the new task on the new dataset, rather than training from scratch [21]. Fine-tuning a pre-trained language model (e.g., BERT [21] and GPT-2 [74]) is a popular training technique for NLP tasks because it can save computation overhead and achieve state-of-the-art results for many tasks, e.g., sequence classification,

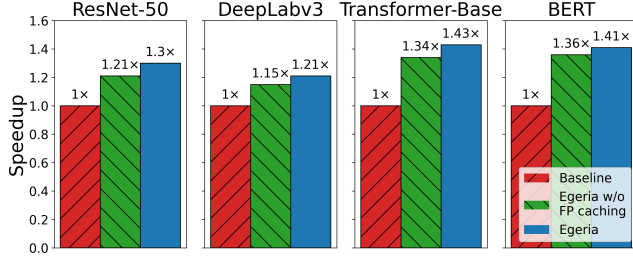


Figure 10: Performance breakdown of layer freezing and prefetching FP pass.

sentiment analysis, and text summarization.

Figure 9d shows the results of fine-tuning BERT on the SQuAD 1.0 dataset. The metric for question answering is the F1 score. KGT accelerates the baseline by 41% to reach the target accuracy. Since fine-tuning converges faster than training from scratch, KGT does not freeze many deep layers before achieving the target, but the frozen front layers can still provide a good speedup. During the training, the learning rate scheduler does not trigger the unfreezing mechanism. This task shows that KGT can not only train generic models but also perform well in model fine-tuning and transfer learning.

6.3 Performance Breakdown

FP caching benefits. For single-node training, the performance speedup comes from the BP computation of the frozen layers and prefetching the cached FP computation results. In Figure 10, we can see that caching FP generally contributes more for CNN models than language models but are all less than 10%. If there are few frozen layers or the front layers have fewer parameters, FP caching will be disabled.

Distributed training. KGT also works in multi-node distributed training and could bring an even larger speedup than single-node training, as shown in Figure 11. We also compare with ByteScheduler [69], the state-of-the-art communication scheduling framework. KGT can work together with distributed communication optimizations like ByteScheduler.

ResNet-50 and Transformer are both computation-intensive models, just like most of the recent DNN architectures, so the performance improvement of ByteScheduler is limited here. A slight throughput drop when communication is not the bottleneck is normal for ByteScheduler with the default configuration [68]. While the benefits of KGT mostly come from computation saving, since frozen layers are not required for parameter synchronization, the reduced communication traffic can speedup the training by up to 5% for ResNet-50, which can benefit the linear scalability of large scale training.

Freezing & unfreezing decisions. We take a closer look at one of our evaluations, training ResNet-56, to understand the

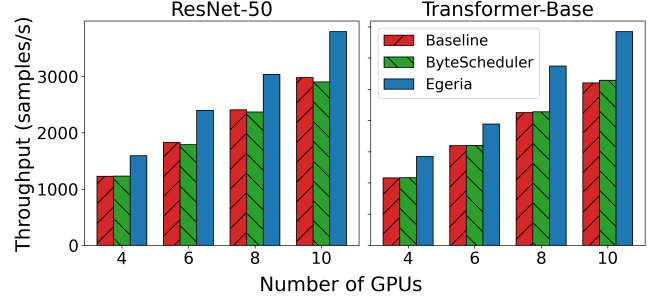


Figure 11: Distributed training performance. KGT freezes layers to exclude them from parameter synchronization.

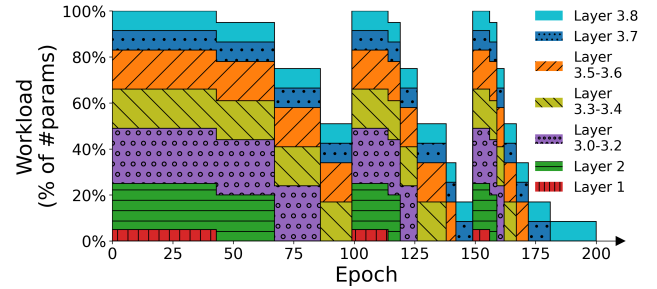


Figure 12: Freezing and unfreezing decisions breakdown through the ResNet-56's 200-epoch training. Y-axis shows the percentage of the active layers' parameters (their sizes).

decisions made by KGT in Figure 12. The bottom-up DNN consists of layer 1.0–1.8, 2.0–2.8, 3.0–3.8, and input/output layers adjacent to layer 1.0 and 3.8. KGT parses the model based on its structure and the size of each layer, so that layer 3 (75% of the total parameters), which is significantly larger than layer 2 (20%), is split finer-grained into similar-sized modules; while layer 1 (5%) and layer 2 are evaluated as a whole. Layer 3.7–3.8 (17%) is further split because it is the last module. KGT gradually freezes layers and remarkably reduces the training cost (the blanks) without hurting the accuracy. Refreezing after the 100th and 150th epochs' unfreezing takes much less time because of the relaxed criteria (§4.2.2).

6.4 Sensitivity Analysis

Impact of the reference model's precision on accuracy. KGT generates the reference model using int8 quantization by default for CPU execution efficiency (§4.1). We evaluate using higher precisions for the reference model, including float16 and float32 (full-precision), in ResNet-56 training on CIFAR-10, as shown in Table 2. We find that using the int8-quantized reference model, which averagely has a 0.6% lower accuracy (-1.2% to -0.3%), will not degrade the final accuracy and can largely improve the inference speed to obtain the intermediate activation. Other tasks show similar results.

Performance	int8	float16	float32
Final accuracy	92.1%	92.0%	92.2%
CPU inference speed	3.59×	1.69×	1×
Reference acc. gap	-0.6%	-0.2%	0

Table 2: Using difference precisions for the reference model.

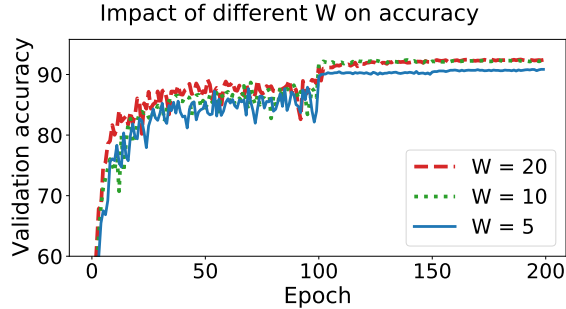


Figure 13: W is generally robust to keep the final accuracy.

Impact of the linear fitting window W on accuracy. KGT uses the slope of linear fitting on a moving window W of recent plasticity values to filter out the drastic fluctuation in SGD training and provide a recent history context than simply evaluating delta. (§4.2.2). We find that W is tolerant and robust in general if following our guideline. Only manually lowering W (e.g., to 5) may eagerly freeze unconverged layers and hurt the final accuracy, as shown in Figure 13.

6.5 System Overhead and Discussion

KGT leverages abundant CPU cores to freeze layers accurately while maintaining accuracy; it also uses disk storage to reduce forward computation overhead. Through careful system designs, we minimize the extra overhead of KGT.

The reference model. Using the reference model for plasticity evaluation involves the model generation and execution. We find that generating and updating the reference model through dynamic and static quantization on CPU take 0.5s–1.5s, thus bring no noticeable slowdown. Running the reference model on CPU could introduce up to 1.5% time overhead, which is worthwhile compared to the saving from layer freezing. If CPU resources are limited (e.g., on shared machines or using CPU-based optimizations), we support GPU execution for the reference model.

Caching and prefetching. We store the serialized intermediate tensors of the frozen layers to the disk for prefetching. The storage usage depends on the DNN architecture. For example, we need $1.5\times$ to $5.3\times$ compared to the input for ResNet-50, which is generally viable. For language models,

since the text data volume is smaller, the overall space usage is limited. Since we only keep the relevant tensors in memory, the overhead is small compared to the regular utilization. It takes hundreds of MB of GPU memory, which is a small fraction of device memory for modern GPUs.

Generalization. We design KGT as a general system. Users can adjust the usage of CPU and storage in plasticity evaluation and caching to meet their needs. Future research can study how KGT collaborates with other CPU-based (e.g., BytePS [44]) or storage-based (e.g., CoorDL [60]) optimizations and on different hardware.

7 Related Work

Efficient training systems. Accelerating DNN training is a key goal of ML systems. To optimize the computation, people may optimize the computation graph to maximize the degree of parallelism [43], or use advanced scheduling to distribute the computation across multiple machines [25, 59]. To optimize the communication, they use the layered structural information to prioritize the front layers in communication to avoid blocking layers with high priority [41, 44, 57, 69]. Furthermore, some efforts [34, 90] try to relax the synchronization requirement to speed up, while others [40, 54] measure the importance of gradient updates and filter out trivial parameters before shipping them out. Additionally, there are a wide range of networking solutions that can help in distributed DNN training, including but not limited to: performing in-network aggregation to reduce traffic volume [14, 48, 80], exploiting fast network protocols, e.g., RDMA [26, 98], designing ML-specific communication protocols [89, 94] and topology-aware parameter exchanging schemes [44, 88], reducing flow completion time by congestion control [3, 7, 15], flow scheduling [6, 52] and co-flow scheduling [84, 103, 105], etc. KGT aims to reduce the total training workload, thus should be compatible with these techniques.

Using an assistant model in training. The idea of using another DNN to assist training has been discussed on several ML topics. Co-distillation [5] trains multiple tweaked copies of the model in distributed manner and lets them share knowledge by encouraging one model to agree with others’ predictions. AutoAssistant [104] trains a lightweight assistant model to identify the hard-to-classify examples and replace the easy examples to feed the training model to accelerate training. Infer2Train [35] runs a copy of the training model on the additional hardware inference accelerator and finds the difficult examples to prioritize in the following iterations.

Freezing parameters and caching DNN results. Transfer learning practitioners usually freeze the front layers and only fine-tune the rest of the pre-trained model on another dataset

for a new task [46, 73]. FreezeOut [12] applies the freezing technique in DNN training and shows that we can trade-off accuracy for higher speed. SpotTune [27] finds which layers to fine-tune per input example to achieve better accuracy. PipeTransformer [30] applies layer freezing in Transformer training with pipeline parallelism using a gradient-based importance metric [95]. APF [13] excludes stable parameters from synchronization in federated learning; it suggests that model snapshots can best capture the performance but implements a workaround measure. Besides training, GATI [9] accelerates DNN inference by caching the intermediate results and skipping the rest of the forward pass.

8 Conclusion

We introduce a novel system KGT to accelerate DNN training while maintaining accuracy by freezing the converged layers, thus reducing the training cost. To avoid freezing layers prematurely, we employ a reference model and use semantic knowledge to evaluate the *plasticity* of internal layers efficiently during training. As a result, KGT can exclude the frozen layers from the backward pass and parameter synchronization. Furthermore, we cache the frozen layers' intermediate computation with prefetching to skip the forward pass. We evaluate KGT using several CV and language models on popular datasets and find that KGT can accelerate training by 19%-43% without hurting model accuracy.

References

- [1] Alessandro Achille, Matteo Rovere, and Stefano Soatto. Critical learning periods in deep networks. In *International Conference on Learning Representations*, 2018.
- [2] Gustavo Aguilar, Yuan Ling, Yu Zhang, Benjamin Yao, Xing Fan, and Chenlei Guo. Knowledge distillation from internal representations. In *AAAI*, pages 7350–7357, 2020.
- [3] Mohammad Alizadeh, Albert Greenberg, David A Maltz, Jitendra Padhye, Parveen Patel, Balaji Prabhakar, Sudipta Sengupta, and Murari Sridharan. Data center tcp (dctcp). In *Proceedings of the ACM SIGCOMM 2010 Conference*, pages 63–74, 2010.
- [4] Amazon. Amazon EC2 P3 instances. <https://aws.amazon.com/ec2/instance-types/p3/>.
- [5] Rohan Anil, Gabriel Pereyra, Alexandre Passos, Robert Ormandi, George E Dahl, and Geoffrey E Hinton. Large scale distributed neural network training through online distillation. In *6th International Conference on Learning Representations, ICLR*, 2018.
- [6] Wei Bai, Li Chen, Kai Chen, Dongsu Han, Chen Tian, and Hao Wang. Information-agnostic flow scheduling for commodity data centers. In *12th {USENIX} Symposium on Networked Systems Design and Implementation ({NSDI} 15)*, pages 455–468, 2015.
- [7] Wei Bai, Shuihai Hu, Kai Chen, Kun Tan, and Yongqiang Xiong. One more config is enough: Saving (dc) tcp for high-speed extremely shallow-buffered datacenters. *IEEE/ACM Transactions on Networking*, 29(2):489–502, 2020.
- [8] Zhihao Bai, Zhen Zhang, Yibo Zhu, and Xin Jin. Pipeswitch: Fast pipelined context switching for deep learning applications. In *14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20)*, pages 499–514, November 2020.
- [9] Arjun Balasubramanian, Adarsh Kumar, Yuhan Liu, Han Cao, Shivaram Venkataraman, and Aditya Akella. Accelerating deep learning inference via learned caches. *arXiv preprint arXiv:2101.07344*, 2021.
- [10] Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Varvara Logacheva, Christof Monz, et al. Findings of the 2016 conference on machine translation. In *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*, pages 131–198, 2016.
- [11] Léon Bottou. Large-scale machine learning with stochastic gradient descent. In *Proceedings of COMPSTAT'2010*, pages 177–186. Springer, 2010.
- [12] Andrew Brock, Theodore Lim, James M Ritchie, and Nick Weston. FreezeOut: Accelerate training by progressively freezing layers. *arXiv preprint arXiv:1706.04983*, 2017.
- [13] Chen Chen, Hong Xu, Wei Wang, Baochun Li, Bo Li, Li Chen, and Gong Zhang. Communication-efficient federated learning with adaptive parameter freezing. In *2021 IEEE 41st International Conference on Distributed Computing Systems (ICDCS)*, pages 1–11. IEEE, 2021.
- [14] Li Chen, Ge Chen, Justinas Lingys, and Kai Chen. Programmable switch as a parallel computing device. *arXiv preprint arXiv:1803.01491*, 2018.
- [15] Li Chen, Shuihai Hu, Kai Chen, Haitao Wu, and Danny HK Tsang. Towards minimal-delay deadline-driven data center tcp. In *Proceedings of the Twelfth ACM Workshop on Hot Topics in Networks*, pages 1–7, 2013.

- [16] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587*, 2017.
- [17] Francois Chollet et al. *Deep learning with Python*, volume 361. Manning New York, 2018.
- [18] Andy Coenen, Emily Reif, Ann Yuan, Been Kim, Adam Pearce, Fernanda Viégas, and Martin Wattenberg. Visualizing and measuring the geometry of bert. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, pages 8594–8603, 2019.
- [19] Moheb Costandi. *Neuroplasticity*. MIT Press, 2016.
- [20] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [21] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [22] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The pascal visual object classes (voc) challenge. *International Journal of Computer Vision*, 88(2):303–338, June 2010.
- [23] L. Fei-Fei, A. Karpathy, and J. Johnson. Cs231n: Convolutional neural networks for visual recognition: Annealing the learning rate.
- [24] L. Fei-Fei, A. Karpathy, and J. Johnson. Cs231n: Convolutional neural networks for visual recognition: Transfer learning.
- [25] Juncheng Gu, Mosharaf Chowdhury, Kang G. Shin, Yibo Zhu, Myeongjae Jeon, Junjie Qian, Hongqiang Liu, and Chuanxiong Guo. Tiresias: A GPU cluster manager for distributed deep learning. In *16th USENIX Symposium on Networked Systems Design and Implementation (NSDI 19)*, pages 485–500, Boston, MA, 2019.
- [26] Chuanxiong Guo, Haitao Wu, Zhong Deng, Gaurav Soni, Jianxi Ye, Jitu Padhye, and Marina Lipshteyn. Rdma over commodity ethernet at scale. In *Proceedings of the 2016 ACM SIGCOMM Conference*, pages 202–215, 2016.
- [27] Yunhui Guo, Honghui Shi, Abhishek Kumar, Kristen Grauman, Tajana Rosing, and Rogerio Feris. Spottune: Transfer learning through adaptive fine-tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [28] Song Han, Huizi Mao, and William J Dally. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding. *arXiv preprint arXiv:1510.00149*, 2015.
- [29] Sayed Hadi Hashemi, Sangeetha Abdu Jyothi, and Roy H Campbell. TicTac: Accelerating distributed deep learning with communication scheduling. In *Proceedings of Machine Learning and Systems, MLSys*, 2019.
- [30] Chaoyang He, Shen Li, Mahdi Soltanolkotabi, and Salman Avestimehr. Pipetransformer: Automated elastic pipelining for distributed training of large-scale models. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 4150–4159. PMLR, 18–24 Jul 2021.
- [31] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In *European conference on computer vision*, pages 630–645. Springer, 2016.
- [32] Byeongho Heo, Minsik Lee, Sangdoo Yun, and Jin Young Choi. Knowledge transfer via distillation of activation boundaries formed by hidden neurons. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3779–3787, 2019.
- [33] Geoffrey Hinton, Oriol Vinyals, and Jeffrey Dean. Distilling the knowledge in a neural network. In *NIPS Deep Learning and Representation Learning Workshop*, 2015.
- [34] Qirong Ho, James Cipar, Henggang Cui, Jin Kyu Kim, Seunghak Lee, Phillip B Gibbons, Garth A Gibson, Gregory R Ganger, and Eric P Xing. More effective distributed ml via a stale synchronous parallel parameter server. *Advances in neural information processing systems*, 2013:1223, 2013.
- [35] Elad Hoffer, Berry Weinstein, Itay Hubara, Sergei Gofman, and Daniel Soudry. Infer2train: leveraging inference for better training of deep networks. In *NeurIPS 2018 Workshop on Systems for ML*, 2018.
- [36] Jeremy Howard and Sebastian Ruder. Universal language model fine-tuning for text classification. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 328–339, 2018.

- [37] Pan Hu, Junha Im, Zain Asgar, and Sachin Katti. Starfish: resilient image compression for aiOT cameras. In *Proceedings of the 18th Conference on Embedded Networked Sensor Systems*, pages 395–408, 2020.
- [38] Yanping Huang, Youlong Cheng, Ankur Bapna, Orhan Firat, Mia Xu Chen, Dehao Chen, HyounJoong Lee, Jiquan Ngiam, Quoc V Le, Yonghui Wu, et al. Gpipe: Efficient training of giant neural networks using pipeline parallelism. *arXiv preprint arXiv:1811.06965*, 2018.
- [39] Intel. oneDNN library. <https://01.org/oneDNN/>.
- [40] Nikita Iykin, Daniel Rothchild, Enayat Ullah, Vladimir Braverman, Ion Stoica, and Raman Arora. Communication-efficient distributed sgd with sketching. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- [41] Anand Jayarajan, Jinliang Wei, Garth Gibson, Alexandra Fedorova, and Gennady Pekhimenko. Priority-based parameter propagation for distributed DNN training. In *Proceedings of Machine Learning and Systems 2019, MLSys 2019, Stanford, CA, USA, March 31 - April 2, 2019*, 2019.
- [42] Myeongjae Jeon, Shivaram Venkataraman, Amar Phanishayee, Junjie Qian, Wencong Xiao, and Fan Yang. Analysis of large-scale multi-tenant GPU clusters for DNN training workloads. In *2019 USENIX Annual Technical Conference (USENIX ATC 19)*, pages 947–960, 2019.
- [43] Zhihao Jia, Oded Padon, James Thomas, Todd Warszawski, Matei Zaharia, and Alex Aiken. Taso: Optimizing deep learning computation with automatic generation of graph substitutions. In *Proceedings of the 27th ACM Symposium on Operating Systems Principles, SOSP ’19*, page 47–62, 2019.
- [44] Yimin Jiang, Yibo Zhu, Chang Lan, Bairen Yi, Yong Cui, and Chuanxiong Guo. A unified architecture for accelerating distributed DNN training in heterogeneous GPU/CPU clusters. In *14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20)*, pages 463–479, 2020.
- [45] Keras. Earlystopping. https://keras.io/api/callbacks/early_stopping/.
- [46] Keras. Transfer learning & fine-tuning. https://keras.io/guides/transfer_learning/.
- [47] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009.
- [48] ChonLam Lao, Yanfang Le, Kshiteej Mahajan, Yixi Chen, Wenfei Wu, Aditya Akella, and Michael Swift. ATP: In-network aggregation for multi-tenant learning. In *18th USENIX Symposium on Networked Systems Design and Implementation (NSDI 21)*, pages 741–761. USENIX Association, April 2021.
- [49] Jaeeun Lee, Raphael Tang, and Jimmy Lin. What would elsa do? freezing layers during transformer fine-tuning. *arXiv preprint arXiv:1911.03090*, 2019.
- [50] Changlin Li, Jiefeng Peng, Liuchun Yuan, Guangrun Wang, Xiaodan Liang, Liang Lin, and Xiaojun Chang. Block-wisely supervised neural architecture search with knowledge distillation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [51] Zhuohan Li, Eric Wallace, Sheng Shen, Kevin Lin, Kurt Keutzer, Dan Klein, and Joey Gonzalez. Train large, then compress: Rethinking model size for efficient training and inference of transformers. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 5958–5968. PMLR, 13–18 Jul 2020.
- [52] Ziyang Li, Wei Bai, Kai Chen, Dongsu Han, Yiming Zhang, Dongsheng Li, and Hongfang Yu. Rate-aware flow scheduling for commodity data center networks. In *IEEE INFOCOM 2017-IEEE Conference on Computer Communications*, pages 1–9. IEEE, 2017.
- [53] Ji Lin, Wei-Ming Chen, Yujun Lin, Chuang Gan, and Song Han. Mccnet: Tiny deep learning on iot devices. *Advances in Neural Information Processing Systems*, 33, 2020.
- [54] Yujun Lin, Song Han, Huizi Mao, Yu Wang, and Bill Dally. Deep gradient compression: Reducing the communication bandwidth for distributed training. In *International Conference on Learning Representations*, 2018.
- [55] Mingsheng Long, Yue Cao, Jianmin Wang, and Michael Jordan. Learning transferable features with deep adaptation networks. In *International conference on machine learning*, pages 97–105. PMLR, 2015.
- [56] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016.
- [57] Yiqing Ma, Hao Wang, Yiming Zhang, and Kai Chen. Automatic configuration for optimal communication scheduling in dnn training. *arXiv preprint arXiv:2112.13509*, 2021.

- [58] Dougal Maclaurin, David Duvenaud, and Ryan P Adams. Autograd: Effortless gradients in numpy. In *ICML 2015 AutoML Workshop*, volume 238, page 5, 2015.
- [59] Kshiteej Mahajan, Arjun Balasubramanian, Arjun Singhvi, Shivaram Venkataraman, Aditya Akella, Amar Phanishayee, and Shuchi Chawla. Themis: Fair and efficient GPU cluster scheduling. In *17th USENIX Symposium on Networked Systems Design and Implementation (NSDI 20)*, pages 289–304, Santa Clara, CA, 2020.
- [60] Jayashree Mohan, Amar Phanishayee, Ashish Raniwala, and Vijay Chidambaram. Analyzing and mitigating data stalls in dnn training. *Proc. VLDB Endow.*, 14(5):771–784, January 2021.
- [61] Ari S Morcos, Maithra Raghu, and Samy Bengio. Insights on representational similarity in neural networks with canonical correlation. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 5732–5741, 2018.
- [62] Deepak Narayanan, Aaron Harlap, Amar Phanishayee, Vivek Seshadri, Nikhil R. Devanur, Gregory R. Ganger, Phillip B. Gibbons, and Matei Zaharia. PipeDream: Generalized pipeline parallelism for DNN training. In *Proceedings of the 27th ACM Symposium on Operating Systems Principles, SOSP ’19*, page 1–15, New York, NY, USA, 2019. Association for Computing Machinery.
- [63] Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of NAACL-HLT 2019: Demonstrations*, 2019.
- [64] Ian Ozsvald and Micha Gorelick. Dictionaries and sets. In *High Performance Python: Practical Performant Programming for Humans*, chapter 4, pages 73–88. O’Reilly Media, 2014.
- [65] Xingyuan Pan and Vivek Srikumar. Expressiveness of rectifier networks. In *International Conference on Machine Learning*, pages 2427–2435, 2016.
- [66] Geondo Park, Gyeongman Kim, and Eunho Yang. Distilling linguistic context for language model compression. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 2021.
- [67] David Patterson, Joseph Gonzalez, Quoc Le, Chen Liang, Lluís-Miquel Munguia, Daniel Rothchild, David So, Maud Texier, and Jeff Dean. Carbon emissions and large neural network training. *arXiv preprint arXiv:2104.10350*, 2021.
- [68] Yanghua Peng and ByteScheduler team. ByteScheduler issues. <https://github.com/bytedance/byteps/issues?q=label%3Abytescheduler>.
- [69] Yanghua Peng, Yibo Zhu, Yangrui Chen, Yixin Bao, Bairen Yi, Chang Lan, Chuan Wu, and Chuanxiong Guo. A generic communication scheduler for distributed DNN training acceleration. In *Proceedings of the 27th ACM Symposium on Operating Systems Principles, SOSP ’19*, page 16–29, 2019.
- [70] PyTorch. PyTorch autograd mechanics. <https://pytorch.org/docs/stable/notes/autograd.html>.
- [71] PyTorch. PyTorch DataLoader. <https://pytorch.org/docs/stable/data.html#torch.utils.data.DataLoader>.
- [72] PyTorch. PyTorch Vision. <https://github.com/pytorch/vision>.
- [73] PyTorch. Transfer learning for computer vision tutorial. https://pytorch.org/tutorials/beginner/transfer_learning_tutorial.html.
- [74] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [75] Maithra Raghu, Justin Gilmer, Jason Yosinski, and Jascha Sohl-Dickstein. Svcca: Singular vector canonical correlation analysis for deep understanding and improvement. *network*, 200(200):200, 2017.
- [76] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250*, 2016.
- [77] Anna Rogers, Olga Kovaleva, and Anna Rumshisky. A primer in bertology: What we know about how bert works. *Transactions of the Association for Computational Linguistics*, 8:842–866, 2020.
- [78] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. Fitnets: Hints for thin deep nets. *arXiv preprint arXiv:1412.6550*, 2014.
- [79] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. MobileNetV2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4510–4520, 2018.

- [80] Amedeo Sapio, Marco Canini, Chen-Yu Ho, Jacob Nelson, Panos Kalnis, Changhoon Kim, Arvind Krishnamurthy, Masoud Moshref, Dan Ports, and Peter Richtarik. Scaling distributed machine learning with in-network aggregation. In *18th {USENIX} Symposium on Networked Systems Design and Implementation ({NSDI} 21)*, pages 785–808, 2021.
- [81] Scikit-learn. Early stopping of stochastic gradient descent. https://scikit-learn.org/stable/auto_examples/linear_model/plot_sgd_early_stopping.html.
- [82] Sheng Shen, Alexei Baevski, Ari S Morcos, Kurt Keutzer, Michael Auli, and Douwe Kiela. Reservoir transformer. *arXiv preprint arXiv:2012.15045*, 2020.
- [83] Leslie N Smith. Cyclical learning rates for training neural networks. In *2017 IEEE winter conference on applications of computer vision (WACV)*, pages 464–472. IEEE, 2017.
- [84] Hengky Susanto, Hao Jin, and Kai Chen. Stream: Decentralized opportunistic inter-coflow scheduling for datacenter networks. In *2016 IEEE 24th International Conference on Network Protocols (ICNP)*, pages 1–10. IEEE, 2016.
- [85] TensorFlow. TensorFlow stateless random image transformations. https://www.tensorflow.org/tutorials/images/data_augmentation#random_transformations.
- [86] Frederick Tung and Greg Mori. Similarity-preserving knowledge distillation. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019.
- [87] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *arXiv preprint arXiv:1706.03762*, 2017.
- [88] Xinchen Wan, Hong Zhang, Hao Wang, Shuihai Hu, Junxue Zhang, and Kai Chen. Rat-resilient allreduce tree for distributed machine learning. In *4th Asia-Pacific Workshop on Networking*, pages 52–57, 2020.
- [89] Hao Wang, Jingrong Chen, Xinchen Wan, Han Tian, Jiacheng Xia, Gaoxiong Zeng, Weiyan Wang, Kai Chen, Wei Bai, and Junchen Jiang. Domain-specific communication optimization for distributed dnn training. *arXiv preprint arXiv:2008.08445*, 2020.
- [90] Weiyan Wang, Cengguang Zhang, Liu Yang, Jiacheng Xia, Kai Chen, and Kun Tan. Divide-and-shuffle synchronization for distributed machine learning. *arXiv preprint arXiv:2007.03298*, 2020.
- [91] Wei Wen, Cong Xu, Feng Yan, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. Terngrad: ternary gradients to reduce communication in distributed deep learning. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 1508–1518, 2017.
- [92] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45. Association for Computational Linguistics, October 2020.
- [93] Carole-Jean Wu, David Brooks, Kevin Chen, Douglas Chen, Sy Choudhury, Marat Dukhan, Kim Hazelwood, Eldad Isaac, Yangqing Jia, Bill Jia, et al. Machine learning at facebook: Understanding inference at the edge. In *2019 IEEE International Symposium on High Performance Computer Architecture (HPCA)*, pages 331–344. IEEE, 2019.
- [94] Jiacheng Xia, Gaoxiong Zeng, Junxue Zhang, Weiyan Wang, Wei Bai, Junchen Jiang, and Kai Chen. Rethinking transport layer design for distributed machine learning. In *Proceedings of the 3rd Asia-Pacific Workshop on Networking 2019*, pages 22–28, 2019.
- [95] Xueli Xiao, Thosini Bamunu Mudiyansele, Chunyan Ji, Jie Hu, and Yi Pan. Fast deep learning training through intelligently freezing layers. In *2019 International Conference on Internet of Things (iThings) and IEEE Green Computing and Communications (GreenCom) and IEEE Cyber, Physical and Social Computing (CPSCom) and IEEE Smart Data (SmartData)*, pages 1225–1232. IEEE, 2019.
- [96] Xiufeng Xie and Kyu-Han Kim. Source compression with bounded dnn perception loss for iot edge computer vision. In *The 25th Annual International Conference on Mobile Computing and Networking*, pages 1–16, 2019.
- [97] Bowen Yang, Jian Zhang, Jonathan Li, Christopher Ré, Christopher Aberger, and Christopher De Sa. Pipemare: Asynchronous pipeline parallel dnn training. *Proceedings of Machine Learning and Systems*, 3, 2021.
- [98] Bairen Yi, Jiacheng Xia, Li Chen, and Kai Chen. Towards zero copy dataflows using rdma. In *Proceedings*

of the *SIGCOMM Posters and Demos*, pages 28–30. 2017.

- [99] Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in deep neural networks? In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 3320–3328. Curran Associates, Inc., 2014.
- [100] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *European conference on computer vision*, pages 818–833, 2014.
- [101] Chiyuan Zhang, Samy Bengio, and Yoram Singer. Are all layers created equal? *arXiv preprint arXiv:1902.01996*, 2019.
- [102] Hao Zhang, Zeyu Zheng, Shizhen Xu, Wei Dai, Qirong Ho, Xiaodan Liang, Zhiting Hu, Jinliang Wei, Pengtao Xie, and Eric P. Xing. Poseidon: An efficient communication architecture for distributed deep learning on GPU clusters. In *2017 USENIX Annual Technical Conference (USENIX ATC 17)*, pages 181–193, 2017.
- [103] Hong Zhang, Li Chen, Bairen Yi, Kai Chen, Mosharaf Chowdhury, and Yanhui Geng. Coda: Toward automatically identifying and scheduling coflows in the dark. In *Proceedings of the 2016 ACM SIGCOMM Conference*, pages 160–173, 2016.
- [104] Jiong Zhang, Hsiang-Fu Yu, and Inderjit S Dhillon. Autoassist: A framework to accelerate training of deep neural networks. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- [105] Yangming Zhao, Kai Chen, Wei Bai, Minlan Yu, Chen Tian, Yanhui Geng, Yiming Zhang, Dan Li, and Sheng Wang. Rapier: Integrating routing and scheduling for coflow-aware data center networks. In *2015 IEEE Conference on Computer Communications (INFOCOM)*, pages 424–432. IEEE, 2015.
- [106] Barret Zoph and Quoc V Le. Neural architecture search with reinforcement learning. *arXiv preprint arXiv:1611.01578*, 2016.

Appendices

A Intermediate Activation

Figure 14 shows how intermediate activation changes across a CNN given an input image.³ The very front layers work as a collection of edge detectors. Their activations still retain the visually recognizable information from the input image. As we go deeper, the activations become more abstract and less visually interpretable, encoding higher-level concepts such as ears and eyes. The deeper layers are increasingly more “task-specific”. They show less information about the visual contents of the specific input and carry more information related to the class of the image. The activations of the deep layers are more sparse because only the patterns encoded by the filter are activated, and the specific visual appearance is filtered out. Figure 14 shows that intermediate activation can reveal layer performance and let us evaluate the training progress by comparing them with a reference.

B Pair-wise Similarity-Preserving Loss

Our knowledge-guided training system understands the semantic performance of the training DNN layers to decide if we can pause some parameters and put the training resource on those in need. Our goal matches the conceptual idea of knowledge distillation (KD): KD studies how to define the distillation loss to best capture the knowledge between the teacher and the student. Recent KD research informs that using the pair-wise similarity elicited by a mini-batch of data at a certain layer as the loss can better capture the semantic and contextual information [66, 86]. A pair of samples refers to images in a mini-batch or word vectors in a concatenated sentence.

KGT uses SP loss [86] to quantify the training plasticity of a layer (§4.2). Note that we do not consider the ML loss function behind plasticity as our contribution. We choose the similarity-preserving loss because it best matches our goal of saving the training cost. Given the input mini-batch of batch size b , we denote the activation tensors of the training and reference models at a specific layer l as $A_T^{(l)}$ and $A_R^{(l)}$. To calculate the SP Loss between them, we first reshape the tensor to 2D from the view of the batch size. For example, an image tensor $A_T^{(l)} \in \mathbb{R}^{b \times c \times h \times w}$ will become $Q_T^{(l)} \in \mathbb{R}^{b \times chw}$. We multiply $Q_T^{(l)}$ and its transpose to align the activation tensors into $b \times b$ -shaped matrices and apply a row-wise L2 normalization to get $G_T^{(l)}$, where the notation $[i, :]$ denotes the i th row in a matrix.

$$\tilde{G}_T^{(l)} = Q_T^{(l)} \cdot Q_T^{(l)T}; G_{T[i,:]}^{(l)} = \tilde{G}_{T[i,:]}^{(l)} / \left\| \tilde{G}_{T[i,:]}^{(l)} \right\|_2$$

³This example is from the book Deep learning with Python [17].

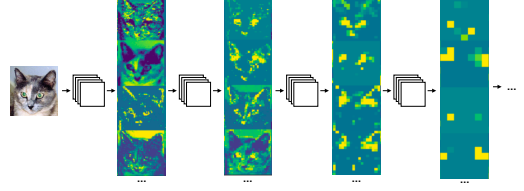


Figure 14: Intermediate activations transition from general to task-specific across layers. Intermediate activations show that how layers transform their input.

Analogously, we obtain $G_R^{(l)}$ of the reference model’s activation.

$$\tilde{G}_R^{(l)} = Q_R^{(l)} \cdot Q_R^{(l)T}; G_{R[i,:]}^{(l)} = \tilde{G}_{R[i,:]}^{(l)} / \left\| \tilde{G}_{R[i,:]}^{(l)} \right\|_2$$

Now we can define the SP loss between the training and reference activation tensor as

$$SP_loss(A_T^{(l)}, A_R^{(l)}) = \frac{1}{b^2} \left\| G_T^{(l)} - G_R^{(l)} \right\|_F^2$$

where $\|\cdot\|$ denotes Frobenius norm. For Transformer-based NLP models, the pair-wise similarity-preserving loss works between word vectors [66].

C Design Details

More theoretical background of the reference model. In KGT, the reference model is essentially an optimized training model snapshot. Recent research on understanding the progress of DNN suggests that the practically optimal method of finding stable parameters relies on maintaining all the historical update snapshots in memory [13] and compare with the latest snapshots all the time. This is not efficient nor viable due to the high cost. Our reference model closely implements this similar idea with effective system designs, so we can accurately detect the stable layers to freeze.

D Layer Convergence Analysis

Our work exploits the ML insight that the training progress of DNN layers within a model differs significantly; the front layers often converge much earlier than the deep layers, as generally observed in both vision and language models [77, 82, 99, 100]. This is because of the transitional nature of ML training: DNN features will transition from being task-agnostic to task-specific from the first to the last layer.

In this section, we illustrate the detailed analysis of the layer convergence insight with an example case and show how SP loss [86], the tool we used to quantify *plasticity*, captures the training status of DNN layers.

In recent years, ML researchers have proven that DNN layers generally converge from the input layer to the output

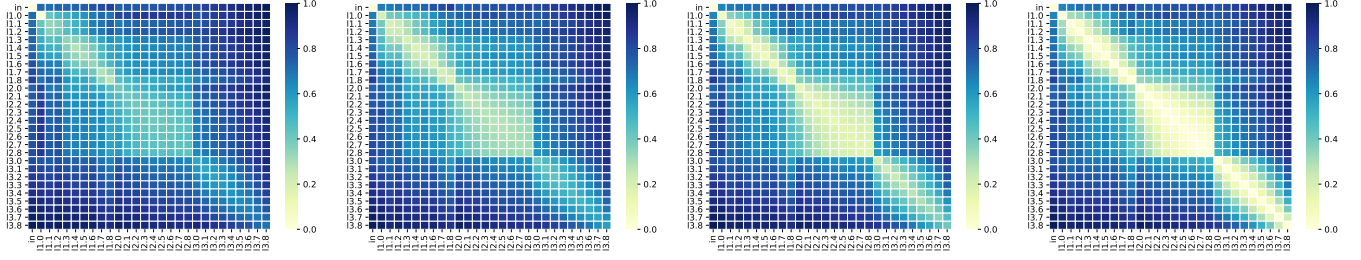


Figure 15: PWCCA heatmap of the intermediate activation across layers.

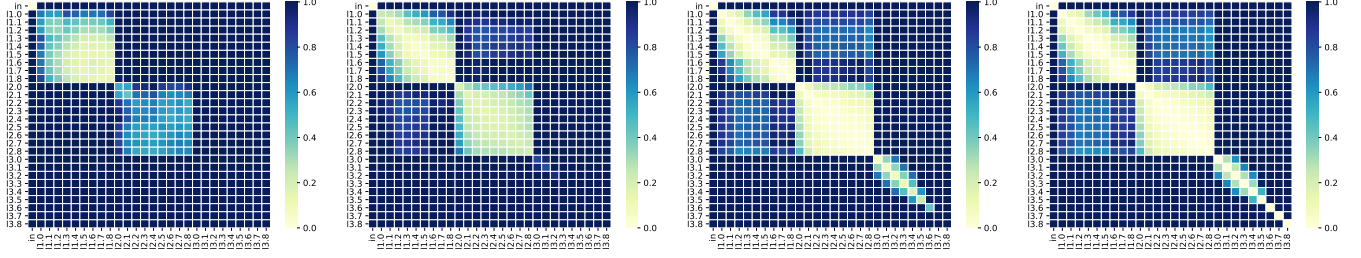


Figure 16: SP loss heatmap of the intermediate activation across layers.

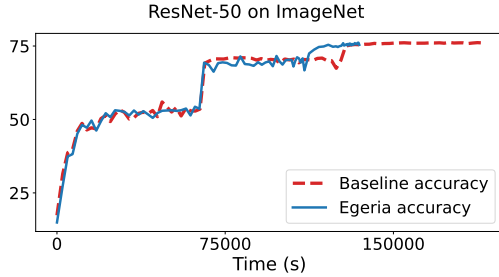


Figure 17: The time-to-accuracy of training ResNet50 for image classification.

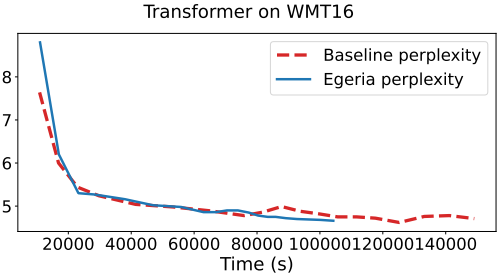


Figure 19: The time-to-accuracy of training Transformer-Base for machine translation.

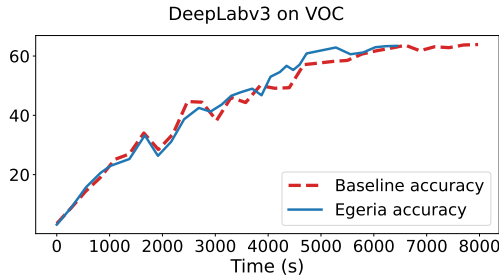


Figure 18: The time-to-accuracy of training DeepLabv3 for semantic segmentation.

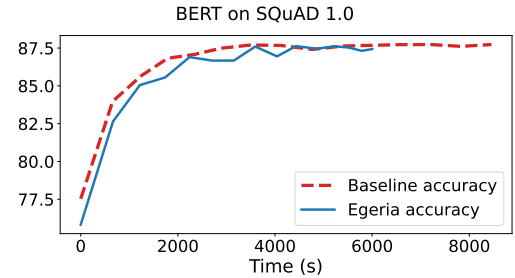


Figure 20: The time-to-accuracy of fine-tuning BERT for question answering.

layer, using post hoc analysis tools, e.g., SVCCA [75] and PWCCA [61], and suggest that some front layers do not need constant training and updating. In Figure 15, we show the PWCCA scores between the layers of a ResNet-56 model training on CIFAR-10 and the layers of its final converged

state. The PWCCA score is evaluated using the activations on the CIFAR-10 test set and normalized to 0-1; 0 means totally converged. In real-world training tasks, we focus on the cells that are on the diagonals of the heatmaps, which show layers comparing to their future selves. From top-left (front layers) to

bottom-right (deep layers), the DNN is converging in a layer-by-layer order. The transforming of a cell on the diagonals across different training stages (left to right) is equivalent to the curve of a layer in Figure 1. Figure 15 proves the correctness of the layer convergence insight we used in KGT. For further theoretical analysis, please refer to the original paper of PWCCA [61].

Similarly, we show how SP loss (plasticity) captures the training progress in a post hoc view in Figure 16. Unlike PWCCA, SP loss is not normalized into a small range: the value of actively updating layers can be $\sim 1000\times$ higher than the converged layers'. This is because they have different utility goals: SP loss is originally designed as a knowledge distillation training loss, so it needs to show the difference accurately; while PWCCA is a visualization tool and uses complex weighted projection to fit largely different values into the same range (0–1), which takes $\sim 10\times$ more computation overhead than SP loss. In Figure 16, we apply a cut-off for values above 1 to better illustrate the convergence trend (of those half-trained layers) without being distracted by the layers that are still drastically updating yet far from convergence. This is why there are more dark cells in Figure 16.

The SP loss trends in Figure 16 echo the PWCCA analysis in Figure 15, showing that the *plasticity* of KGT can accurately seize the layer freezing opportunity observed in the ML training insight. In addition, KGT uses a layer-module-based plasticity evaluation that automatically parses the model structure and divides multiple layers into modules based on their definitions and sizes. Taking ResNet-56 here as an example, all nine layers from layer1.0 to layer1.8 will form a layer module, layer2.0 to layer2.8 will form another module, while layer 3 (3.0 to 3.8) will form multiple modules. This is also explained in §6.3. Figure 16 clearly shows that close layers within a module can show high similarity (i.e., the bright boxes). This suggests that our module-based plasticity can simplify the evaluation pipeline without hurting the scheduling correctness and handle the corner case when fine-grained layers are not strictly converged one by one.

E Time-to-Accuracy Performance

We evaluate that KGT can accelerate training of multiple DNNs (§6.2). We further show the time-to-accuracy results of training ResNet-50 on ImageNet (Figure 17), DeepLabv3 on VOC (Figure 18), Transformer on WMT16 (Figure 19) and fine-tuning BERT on SQuAD 1.0 (Figure 20). KGT provides a higher speedup during the later stage of training because more layers are eligible for freezing in order.