

Pipe-Déjàvu: Hardware-aware Latency Predictable, Differentiable Search for Faster Config and Convergence of Distributed ML Pipeline Parallelism

Pengcheng Xu, Kaiyang Chen, Yuanrui Zhang, Gupta Indranil

Abstract—Pipe-Déjàvu automates pipeline-parallel training of large deep learning models, emphasizing three main innovations: (1) a predictive model that considers communication cost, model computational cost, and hardware information to predict latency and resources of each parallel configurations. It saves time on pre-profiling before searching the parallel configuration; (2) a differentiable parallel configuration search space inspired by DARTS[39], can potentially reach optimal configuration faster than the original dynamic programming; and (3) parallel random initialization employed for faster training loss convergence. By using hardware-aware scheduling based on latency and resources, Pipe-Déjàvu prioritizes training larger layers on more accessible and better computational devices, shortening the longest training stage and reducing overall training time. Given the time estimation on each stage, Pipe-Déjàvu optimizes the pairing between stage and hardware to ease the longest overhead. With knowledge of hardware capacity and network latency, Pipe-Déjàvu predicts optimal scheduling, saving time and resources from pre-profiling. With differentiable search space, it accelerates searching of parallel configuration. With parallel random initialization, it improves speed of loss convergence.

I. INTRODUCTION

Deep learning has made remarkable advancements in recent years, thanks in part to the increasing computational power of hardware and the development of more efficient algorithms. Despite these advancements, training large-scale neural networks remains a time-consuming and resource-intensive task. To address this challenge, we propose Pipe-Déjàvu, a method that automates pipeline-parallel training of large deep learning models, offering significant improvements in training efficiency and resource utilization.

Pipe-Déjàvu emphasizes three main innovations: (1) a predictive model that considers communication cost, model computational cost, and hardware information to predict latency and resources of each parallel configuration, saving time on pre-profiling before searching the parallel configuration; (2) a differentiable parallel configuration search space inspired by DARTS, which can potentially reach optimal configurations faster than the original dynamic programming; and (3) parallel random initialization employed for faster training loss convergence.

By using hardware-aware scheduling based on latency and resources, Pipe-Déjàvu prioritizes training larger layers on more accessible and better computational devices, shortening the longest training stage and reducing overall training time. Given the time estimation on each stage, Pipe-Déjàvu optimizes the pairing between stage and hardware to ease the longest

overhead. With knowledge of hardware capacity and network latency, Pipe-Déjàvu predicts optimal scheduling, saving time and resources from pre-profiling.

In addition to the aforementioned innovations, Pipe-Déjàvu leverages a differentiable search space to accelerate the searching of parallel configurations, further reducing the time required for optimization. Moreover, the parallel random initialization approach improves the speed of loss convergence, contributing to the overall efficiency of the training process.

Our contributions in this paper include:

- 1) A predictive model that considers communication cost, model computational cost, and hardware information to predict latency and resources of each parallel configuration, saving pre-profiling time.
- 2) A differentiable parallel configuration search space inspired by DARTS, which accelerates the searching of optimal configurations.
- 3) Parallel random initialization for faster training loss convergence.
- 4) Hardware-aware scheduling that prioritizes training larger layers on more accessible and better computational devices, reducing overall training time.

The rest of the paper is organized as follows: Section 2 presents a background on gradient descent, data parallelism, model parallelism, and Bayesian optimization. Section 3 describes the proposed parallel random initialization approach in detail, including its integration with the Bayesian optimization framework. Section 4 discusses the algorithmic implementation of the approach, and Section 5 presents experimental results demonstrating its effectiveness, as well as a comparison to other state-of-the-art methods. Finally, Section 6 concludes the paper and discusses future research directions.

II. RELATED WORK

A. Data parallelism

1.1.1 Definition: Data parallelism splits the training data among distributed workers, but copies the model to each worker. Every worker calculates the parameter updates on its own data slice and shares them with other workers before updating the weights. This way, all workers have the same model parameters during training.

1.1.2 Related works: Horovod [18] and PyTorchDDP [20] are two popular data-parallel training systems that use all-reduce to synchronize gradients. BytePS [7, 16] combines all-reduce

and parameter servers and leverages heterogeneous resources in data center clusters. AutoDist [26] employs learning-based methods to devise a data-parallel training strategy. ZeRO [17, 23] enhances the memory usage of data parallelism by reducing duplicated tensors. MiCS [27] lowers the communication scale on top of ZeRO for better scalability on the public cloud.

B. Model/Operator parallelism

1.2.1 Definition: When the model cannot fit in one device, operator parallelism is a viable model parallelism option. Operator parallelism splits the computation of a specific operator, such as "matmul" shown in Appendix Fig. 2b, along non-batch axes, and perform each part of the operator in parallel across multiple devices. Because input tensors are jointly split, when a device performs its op part, the needed portions of input tensors may not be in its local memory. Communication is thus needed to get the input data from other devices. When the tensors are split evenly, i.e., SPMD [24], all devices will follow the same collective communication patterns such as all-reduce, all-gather, and all-to-all.

1.2.2 Related works: The main types of model parallelisms have been discussed in definition. Mesh TensorFlow [19], GSPMD [9,24] and OneFlow [25] provide annotation APIs for users to manually specify the intra-op parallel plan. ColocRL [12] places separate model parts on different devices without pipelining, so the concurrency only happens when there are parallel branches in the model.

C. Pipeline parallelism

1.3.1 Definition: Pipeline parallelism puts different sets of ops from the model graph, called stages, on different workers; at the same time, it divides the training batch into several microbatches, and pipelines the forward and backward passes across microbatches on distributed workers, as Appendix Fig. 2d illustrates. Unlike operator parallelism, pipeline parallelism sends intermediate activations at the forward and backward passes between different workers using point-to-point communication.

1.3.2 Related works: Gpipe [5] divides the input data into micro-batches and forms pipeline parallelisms. PipeDream [13,14] improves GPipe by using asynchronous training algorithms, reducing memory usage, and integrating it with data parallelism. However, PipeDream is asynchronous while Alpa[28] is a synchronous training system. TeraPipe [11] finds a new pipeline parallelism dimension for transformer-based LMs. Google's Pathway system [2] is a concurrent work of Alpa[28]. Pathway advocates a single controller runtime architecture combining "single program multiple data" (SPMD) and "multiple program multiple data" (MPMD) model. This is similar to Alpa's[28] runtime part, where SPMD is used for intra-op parallelisms and MPMD is used for inter-op parallelism.

D. Manual combination of parallelisms

1.4.1 Definition and related works: The latest developments indicate that the methods discussed above have to be combined

to scale up current large DL models [15, 24]. The state-of-the-art training systems, such as Megatron-LM [15, 20], manually create a specific execution plan that merges these parallelisms for transformer language models, which is also called 3D Parallelism. It assumes that the model has the same transformer layer repeated and assigns the same number of layers to each pipeline stage and applies a hand-crafted operator and data parallelism setup uniformly for all layers. It's not generalizable to all models.

E. Automatic combination of parallelisms

1.5.1 Definition: The individual parallelism settings, how they depend on each other and on model and cluster setups, create a complex space that makes it hard to automatically combine these parallelisms. For example, when operator parallelism is used with data parallelism, adding a dataparallel replica means allocating a new set of devices (not just one device) as the worker and finding out the best operator parallelism settings within those devices. When pipeline parallelism is included, the optimal pipelining scheme depends on the data and operator parallelism choices of each pipeline stage and how devices are allocated for each stage. With this view, previous explorations [4, 6, 22, 26] of auto-parallelization are limited to combining data parallelism with at most one model parallelism approach, which misses significant performance opportunities.

1.5.2 Related works: Some research works on how to automatically find the best way to train models in parallel. Tofu [22] uses a dynamic programming method to find the best way to split an operation within a node for linear graphs. FlexFlow [6] defines a "SOAP" problem and uses an MCMC-based random search method. But it only works for placing devices and not for pipeline parallelism. It also cannot handle large graphs or clusters well and does not guarantee optimal. TensorOpt [3] uses a dynamic programming method to find the best way to split an operation that considers both memory and computation cost. Varuna [1] works for clusters with low bandwidth and focuses on finding the best way to use pipeline and data parallelism together. Piper [21] also finds a parallel way that uses both inter- and intra-op parallelism, but it depends on manually designed ways to split operations and assumes a uniform network topology and asynchronous pipeline schedules.

III. MOTIVATION AND CONCRETE PROBLEM STATEMENT

Parallel training has become essential for training today's large models. Parallel training involves distributing the model training process across multiple processing units, such as GPUs or TPUs, which work together to process the data simultaneously. This results in faster convergence and reduces the overall training time. There are many well-studied parallel mechanisms as we describe in above related work, in which pipeline parallel plays an important role in. But pipeline parallel has some inherent shortcomings in fully utilizing all computing units due to sequential dependency of different stages. We can provide a simple example in Fig 2. Assume we divide our

training graph into N stage, stage i generally took t_i to finish, the overall onepass latency for the B batches pipeline is

$$T = \sum_{i=1}^N t_i + (B - 1) \cdot \max_{1 \leq j \leq N} \{t_j\}$$

From the above equation we can know that in order to shorten

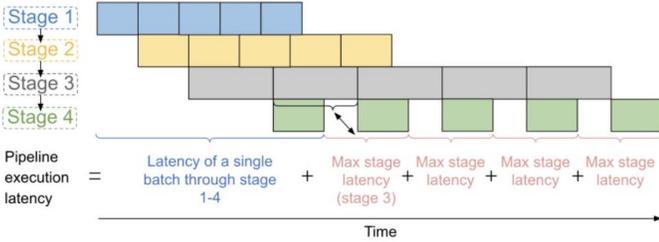


Fig. 1. Example for pipeline parallelism timeline with multiple stages

the overall execution time, our concrete problem is improving the performance (reduce execution time $\max\{t_i\}$) of the longest stage. Instinctively, those longest execution stages are more computational intensive regarding physical or hardware resources like GPU throughput/FLOPS, GPU memory, network latency and even CPU and NVMe capability depending on the user-defined optimization strategy. For example, the sharding strategy introduced by ZeRO incurs heavy communication costs between nodes which emphasize the importance of large network bandwidth. Thus, if we can have some predictor algorithm built in the scheduler, that using information like layer implementation (parameters size / structure) within each stage, hardware conditions of each processing units and user optimization strategy, to map the relatively computational intensive stages to those relatively more "powerful" machine, then we can shorten the execution time of the longest stage compared to random scheduler.

IV. METHODS

A. Hard-ware Latency Predict Model for Saving Pre-profiling Time

Fit a function which inputs the (1.Communication cost 2. Model Computation Cost 3. Hardware information), outputs the (1. Estimated execution time 2. Estimated GPU resources to consume), so that we can plan pipeline parallelism without or with less profiling which takes certain part of automatic parallel configuration time.

$$f(\text{communication}, \text{model}, \text{hardware}) = (\text{time}, \text{resources})$$

To be more specific, intuitively here is a naive example of how we can model this predicted execution latency for specific

ML task mathematically:

$$\begin{aligned} \text{pred_time}(\text{communication}, \text{model}, \text{hardware}) = \\ \alpha \times \frac{g(\#matrix_multiplication_operations)}{h(\text{hardware_computational_speed})} + \\ \beta \times \frac{p(\text{network_bandwidth})}{q(\text{data_to_transmit})} + C \end{aligned} \quad (1)$$

For the model above, we can learn the unknown parameters or function using some linear or nonlinear regression methods.

B. Differentiable Search Space for Faster Parallel Configuration Search

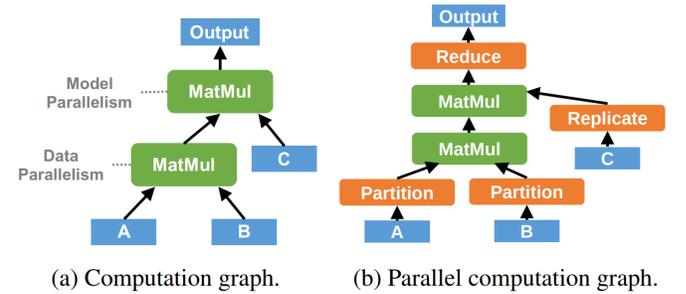
1) **Dynamic Programming:** Many existing automatic pipeline parallel use dynamic programming to minimize the training latency given the time estimated for each stage.

$$T^* = \min_{\substack{s_1, \dots, s_S; \\ (n_1, m_1), \dots, (n_S, m_S)}} \left\{ \sum_{i=1}^S t_i + (B - 1) \cdot \max_{1 \leq j \leq S} \{t_j\} \right\} \quad (2)$$

Alpa's[28] DP algorithm computes the slicing in $O(K^3NM(N + \log(M)))$ time for a fixed $t_{\max} \cdot t_{\max}$ has at most $O(K^2(N + \log(M)))$ choices: $t_{\text{intra}}((o_i, \dots, o_j), \text{Mesh}(n_s, m_s))$ for $i, j = 1, \dots, K$ and all the submesh choices. The complexity of this DP algorithm is thus $O(K^5NM(N + \log(M))^2)$. This complexity is not feasible for a large computational graph of more than ten thousand operators. Therefore, we rethink the search space and make it differentiable to optimize it.

2) **Parallel Computational Graph Representation:** To represent the parallel search space as a computational graph, we can model it using a directed acyclic graph (DAG) where each vertex represents a stage of computation and the edges denote the flow of data or model parameters between the stages. The nodes in this graph can be partitioned into different parallelization strategies, such as data parallelism, model parallelism, and pipeline parallelism.

For example, let's consider a simple graph $G(V, E)$, where V is the set of vertices, and E is the set of edges. We can represent a parallel search space for a two-layer neural network as follows:



(a) Computation graph. (b) Parallel computation graph.

Fig. 2. comparing Computation graph and Parallel Computation graph

Both graphs in the figure above describe the same parallelization of two consecutive matrix multiplications $(A \times B) \times C$ (a simplified form of attention). The green and orange boxes denote regular DNN operators and parallelization operators respectively.

- **Data Parallelism (DP):** In this strategy, each layer of the neural network is replicated across all available devices. The input data is then divided into equal partitions and processed concurrently. The parallel search space for Data Parallelism can be depicted as two distinct sets of nodes, where each set contains nodes corresponding to each device.
- **Model Parallelism (MP):** In Model Parallelism, each layer of the network is assigned to a different device, distributing the model’s layers across multiple devices. This strategy is particularly useful when dealing with large models that exceed the memory capacity of a single device. The parallel search space for this approach can be represented as a single set of nodes, where each node corresponds to a specific layer-device assignment.
- **Pipeline Parallelism (PP):** Pipeline Parallelism, on the other hand, focuses on improving the computational efficiency of model training by exploiting the pipelined execution of layers across different devices. While it also assigns each layer to a different device, similar to Model Parallelism, the key difference lies in the data processing approach. In Pipeline Parallelism, input data is processed in a pipelined manner across the devices, allowing for concurrent execution of different layers, which reduces the overall training time. The parallel search space for Pipeline Parallelism can be depicted as a sequence of stages, where each stage contains a node representing the layer-device assignment and its position in the pipeline.

To represent both parallelization strategies in the same graph, we can create a graph-like structure with multiple layers and vertices. We can then use the softmax function to assign probabilities to the edges between nodes and layers, which represent the likelihood of selecting a specific parallelization strategy. This continuous representation enables us to differentiate the parallelization search space and optimize the parallel configuration along with the model parameters.

3) **Differentiable Parallelization Search Space:** To use a computational graph representation for the differentiable parallelization search space, let’s first define a graph $G(V, E)$, where V is the set of vertices and E is the set of edges. In this graph, each vertex represents a computation stage, and the edges denote the flow of data or model parameters between the stages.

For each vertex v_i in V , we associate a set of N parallelization strategies, represented as $\alpha_{i1}, \alpha_{i2}, \dots, \alpha_{iN}$. We can then represent the search space as a matrix A of size $|V| \times N$, where A_{ij} denotes the discrete choice for strategy j at stage i . We can transform this search space into a differentiable problem

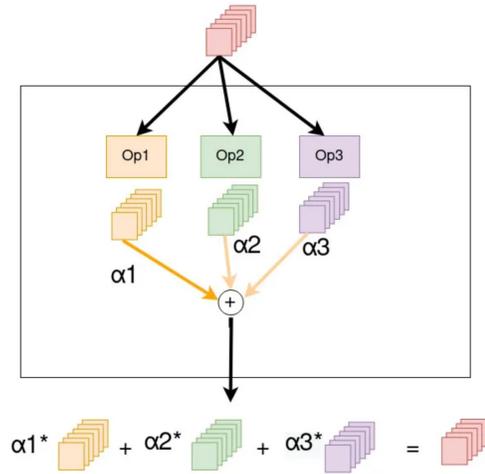


Fig. 3. Softmax Intuition: relaxing the discrete set of candidate operations

by applying the softmax function to each row of A :

$$\text{softmax}(A)_{ij} = \frac{\exp(A_{ij})}{\sum_{k=1}^N \exp(A_{ik})} \quad (3)$$

Here, the continuous approximation of the discrete choice is given by $\text{softmax}(A)_{ij}$ for each strategy j at stage i .

Now, we can use this continuous representation of the search space in the context of a computational graph to describe the algebra transformations and parallel strategies. Each vertex v_i in the graph can be associated with a continuous approximation of its parallelization strategy, given by the row $\text{softmax}(A)_i$. These continuous approximations can then be used to compute the forward and backward passes in the computational graph while considering the different parallel strategies.

The rest of the optimization process, including gradient computation and end-to-end optimization, remains the same as in the previous response, with the parallel configuration matrix A being updated during the optimization.

This approach allows us to incorporate a graph representation of the parallelization search space into the differentiable optimization process, enabling the joint optimization of model training and parallelization strategies.

4) **Bilevel Optimization:** In the DARTS[39] paper, the authors propose a bilevel optimization problem where they optimize the model’s architecture and its weights simultaneously. To adapt this approach for parallelization strategies, we’ll introduce a loss function that jointly evaluates the parallel strategy and the training loss.

Let’s denote the model’s weights by θ and the parallelization strategy represented as a matrix A . We’ll define the loss function as $L(\theta, A)$, which is a combination of the training loss and the evaluation of the parallel strategy. We can write the joint optimization problem as:

$$\text{minimize } L(\theta, A), \quad \text{with respect to } \theta, A \quad (4)$$

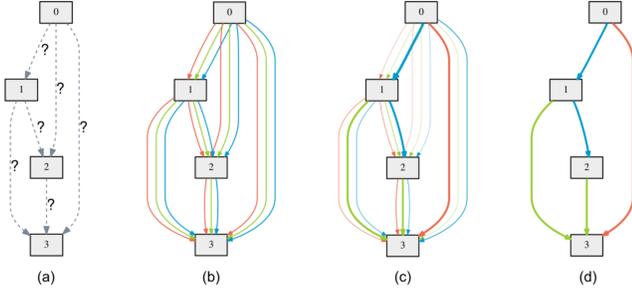


Fig. 4. An overview of Differentiable Search Space: (a) Operations on the edges are initially unknown. (b) Continuous relaxation of the search space by placing a mixture of candidate operations on each edge. (c) Joint optimization of the mixing probabilities and the network weights by solving a bilevel optimization problem. (d) Inducing the final architecture from the learned mixing probabilities.

This bilevel optimization problem can be approximated using gradient-based optimization. The gradients required for updating θ and A are as follows:

$$\begin{aligned}\nabla_{\theta}L(\theta, A) &= \frac{dL(\theta, A)}{d\theta} \\ \nabla_A L(\theta, A) &= \frac{dL(\theta, A)}{dA}\end{aligned}\quad (5)$$

Here, $\nabla_{\theta}L(\theta, A)$ is the gradient with respect to the model parameters θ , while $\nabla_A L(\theta, A)$ is the gradient with respect to the parallelization strategy matrix A .

To perform end-to-end optimization, we'll update both the model parameters and the parallelization strategy using gradient descent:

$$\begin{aligned}\theta &= \theta - \eta_{\theta}\nabla_{\theta}L(\theta, A) \\ A &= A - \eta_A\nabla_A L(\theta, A)\end{aligned}\quad (6)$$

where η_{θ} and η_A are the learning rates for the model parameters and the parallelization strategy matrix A , respectively.

This approach allows us to jointly optimize the model training and the parallelization strategies using a gradient-based optimization method similar to the one used in the DARTS[39] paper. Note that in practice, the loss function $L(\theta, A)$ should be designed to effectively balance the trade-offs between training loss and the evaluation of the parallel strategy.

5) **Complexity Analysis:** Analyzing the complexity of the differentiable search space inspired by DARTS[39] requires considering the computational costs involved in the forward and backward passes, as well as the optimization process.

Let's break down the complexity analysis:

- **Forward pass complexity:** The forward pass involves computing the softmax for each vertex in the graph. Since there are $|V|$ vertices, each with N parallelization strategies, the complexity of the forward pass is $O(|V|N)$.
- **Backward pass complexity:** The gradient computation involves computing the gradients for both the model parameters (θ) and the parallelization strategy matrix (A). Since there are $|V|$ vertices, each with N parallelization strategies, the complexity of computing the gradients with

respect to A is $O(|V|N)$. The complexity of computing the gradients with respect to the model parameters θ depends on the specific model and its architecture, but let's denote it by $O(C_{\theta})$.

- **Optimization complexity:** Updating both the model parameters and the parallelization strategy matrix involves applying gradient descent, which has a complexity of $O(|V|N)$ for the parallelization strategy matrix A and $O(C_{\theta})$ for the model parameters θ .

In total, the complexity of the differentiable search space inspired by DARTS[39] is:

$$O(|V|N) + O(C_{\theta}) + O(|V|N) + O(C_{\theta}) = O(2|V|N + 2C_{\theta})$$

However, this complexity analysis doesn't account for the number of iterations required for convergence. If the optimization process takes I iterations, the overall complexity becomes:

$$O(I(2|V|N + 2C_{\theta}))$$

This complexity analysis shows that the differentiable search space has a linear dependence on the number of vertices, the number of parallelization strategies, and the complexity of the model. It also highlights the importance of the number of optimization iterations for determining the overall complexity of the approach.

C. Parallel Random Initialization for Faster Loss Convergence

In addition to the commonly known methods such as data parallelism and model parallelism, we propose a novel approach to accelerate the convergence of the training loss during the optimization of neural networks. This approach involves utilizing parallel random initialization, which refers to the initial randomization of the network parameters.

One way to conceptualize gradient descent is to imagine a rugged hill in search of the minimum value, with gradient descent acting as a ball rolling down the steepest path. Mathematically, the goal of optimization is to minimize the objective function $L(\theta)$, where θ represents the parameters of the neural network. The update rule for gradient descent is given by:

$$\theta_{t+1} = \theta_t - \eta_t \nabla_{\theta} L(\theta_t) \quad (7)$$

where η_t is the learning rate at time t and $\nabla_{\theta} L(\theta_t)$ is the gradient of the objective function with respect to the parameters at time t . In data parallelism, n workers are assigned to look at n small directions from a randomly initialized starting point. The gradients are then combined and used for gradient descent.

However, during the initial randomization phase, there may be points where the loss is very low, and gradient descent at those points may quickly reach the minimum. Currently, standard neural network training involves randomization only once at the beginning, rather than sampling multiple times.

In the proposed parallel random initialization approach, we can sample multiple randomization points in parallel at the start of the training process. These points can be filtered based on their respective loss values and gradient magnitudes. Only

N nodes random initialization parallelly

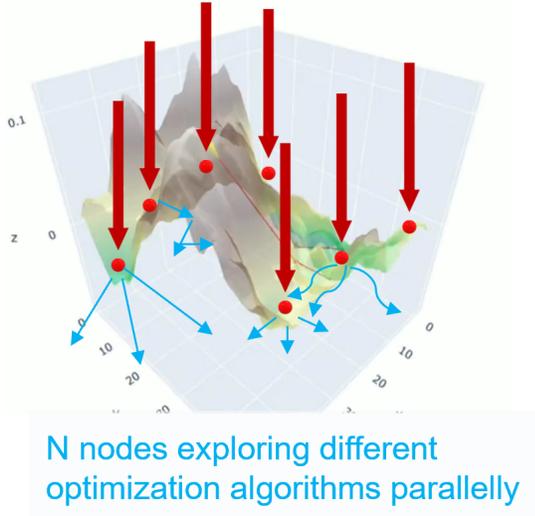


Fig. 5. Parallel Random Initialization and Parallel Optimization Intuition

the ones with small loss and large gradients are retained for further processing.

Subsequently, we can continue to sample within a smaller scope, and further refine the set of points with small loss values. This process can be repeated until a certain level of convergence is achieved, at which point we can focus on the gradient descent of a few selected points.

To implement this approach, we can start with a rough grid with a large interval, and then identify the point with the smallest loss within that grid. We can then perform finer-grained sampling in the vicinity of that point to refine the set of initialization points.

The proposed parallel random initialization approach offers a promising alternative to the standard method of initializing network parameters. By introducing randomness and parallelization during the initialization phase, we can accelerate the convergence of the training loss and potentially reduce the overall training time of neural networks.

D. Model/Pipeline Parallel Fault Tolerance Analysis

Zeno provide a fault-tolerance algorithm from the perspective of SGD for parameter server structure, which is mainly for data parallel. The analysis is mainly on the gradient they send. While after literature review, we find that there are very few research about fault tolerance in model parallel and pipeline parallel. Let w_i be the weight assigned to the i^{th} machine, where $\sum_{i=1}^N w_i = 1$. The overall gradient is computed as a weighted sum of the local gradients from each machine, where the weights are the weights assigned to each machine:

$$g = \sum_{i=1}^N w_i g_i$$

To compute the weight assigned to each machine, we use a suspicion-based approach similar to the one proposed in Zeno.

Algorithm 1 Bayesian Optimization for Parallel Random Initialization

```

1: procedure BAYESIANOPTIMIZATION
2:   function LOSSFUNCTION(params)
3:     initVal  $\leftarrow$  RESHAPEPARAMS(params, NN())
4:     return PARALLELLOSS(NN(), initVal,
5:       trainLoader, device, epoch, w_id)
6:   end function
7:   bounds  $\leftarrow$  [(-1, 1)] * SUMOFNUMEL(NN())
8:   result  $\leftarrow$  GPMINIMIZE(LossFunction,
9:     x0, bounds, numWorkers, randomState)
10:  bestParams  $\leftarrow$  result.x
11:  return RESHAPEPARAMS(bestParams, NN())
12: end procedure
13: function RESHAPEPARAMS(params, model)
14:  return {param.reshape(p.shape) | param, p  $\in$ 
15:    zip(params, model.parameters())}
16: end function
17: function SUMOFNUMEL(model)
18:  return  $\sum_{p \in \text{model.parameters}()} p.\text{numel}()$ 
19: end function

```

Specifically, we compute the suspicion level of each machine i based on the deviation between its local gradient g_i and the overall gradient g :

$$s_i = \|g_i - g\|$$

The suspicion level s_i represents the degree of deviation between the local gradient g_i and the overall gradient g . Machines with high suspicion levels are likely to be faulty and should be given lower weights during the gradient combination process.

We then update the weight assigned to each machine using the following equation:

$$w_i = \frac{\alpha}{\alpha + s_i}$$

where α is a hyperparameter that controls the weight given to the suspicion level. Higher values of α give greater weight to the suspicion level and vice versa. To identify the k machines with the least influence on the overall gradient, we sort the machines based on their weight in ascending order and select the first k machines. We then replace the gradients of the identified machines with the average of the remaining gradients:

$$\hat{g}_i = \frac{\sum_{j=1}^{N-k} g_j}{N - k}$$

The proposed fault-tolerance algorithm for model parallel/pipeline parallel with alpha-like strategy is designed to mitigate the impact of faulty machines on the overall gradient computation. By identifying and replacing the gradients of the k machines with the least influence on the overall gradient, our approach improves the fault tolerance of model parallel/pipeline parallel and can be used to improve the robustness and reliability of distributed machine learning systems.

V. EXPERIMENT RESULTS

A. Hard-ware Latency Predict model for Saving Pre-profiling Time

We use one-shot mathematical formula to replace the pre-profiling of alpa[28], having a reduce of profiling time. In the Table below, we demonstrate the pre-profiling time is significant for large model, while our one-shot prediction only took a few minutes.

Models	Time
GPT-39B	> 24hr
Stable Diffusion	> 12hr
GPT2	> 10hr
ViT	> 6hr

Table: Pre-profiling time for different models

The advantage is that it saves profiling time which can be long, up to 30% percent of whole compilation time. The disadvantage is that it can be not so accurate.

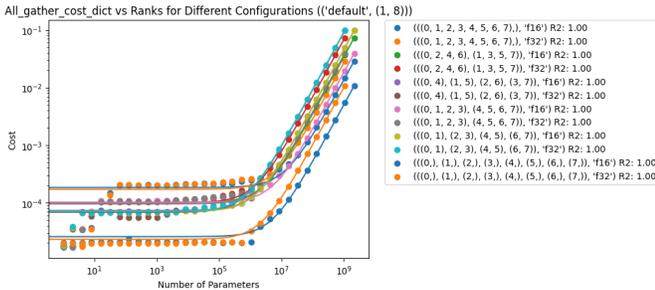


Fig. 6. all gather operation cost measured vs prediction

In figure 6 we have demonstrated the cost for all gather operation. The measured mapping cost under different configuration are represented in dot and our prediction algorithm were drew in solid lines. As number of parameters increase in logarithm, we observed a better fitting line across all configurations. With coefficient of determination reaching 1.0, we proved our regression are capable to predict the cost of placement given appropriate operation and configuration.

B. Differentiable Search Space for Faster Parallel Configuration Search

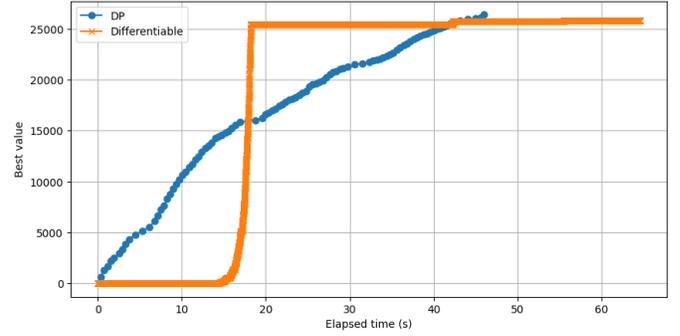


Fig. 7. Differentiable Search results vs Dynamic Programming on Knapsack problem $n = 1000$, $\text{weight_range} = (1, 1000)$, $\text{value_range} = (1, 100)$, $\text{capacity} = 100000$. Higher y axis is better.

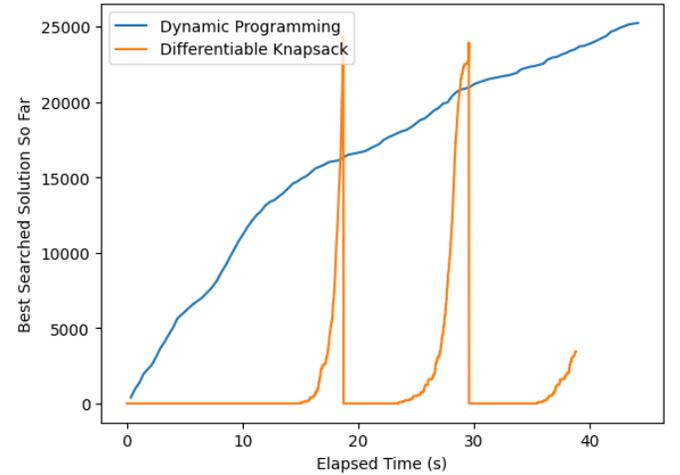


Fig. 8. DP VS Differentiable Search on Knapsack problem $n = 1000$, $\text{weight_range} = (1, 1000)$, $\text{value_range} = (1, 100)$, $\text{capacity} = 100000$. Notice that for Differentiable search, this picture plots the current searched solution for specific Time stamp, the accumulated curve is the above one. The DP curve in this plot is still accumulative so far best solution. We can notice the spike of Differentiable Search Algorithm, displaying the pattern of how it searches the optimal.

We first implemented the differentiable search space onto the classic Knapsack problem. We compete it with the classical dynamic programming algorithm. We found that if the search space is quite small, dynamic programming works faster than differentiable method. However, when the search space becomes extremely larger, as the picture shows (when $n = 1000$, $\text{weight range} = (1, 1000)$, $\text{value range} = (1, 100)$, $\text{capacity} = 100000$. Higher y axis is better.), the differentiable dynamic programming finds the suboptimal solution faster than the dynamic programming. It takes only half of the time of dynamic programming to reach the same suboptimal solution. However, due to the gradient descent is not guaranteed to get the global

optimal solution, it seems that our differentiable search stop on suboptimal and it's hard for it to reach global optimal afterwards. It's an interesting phenomenon. It gives us insight that differentiable may be faster at finding suboptimal when the search space is extremely large.

We may later develop a hybrid algorithm, so that it will reach suboptimal fast, then switch to dynamic programming to find the global optimal.

Now we will discuss the dynamic programming of finding parallel placement of distributed ML training. The complexity of Alpa's[28] DP algorithm is $O(K^5NM(N + \log(M))^2)$. This complexity is not feasible for a large computational graph of more than ten thousand operators.

Therefore, we rethink the search space and use differentiable search space inspired by DARTS[39].

We use differentiable search space instead of dynamic programming, to search optimal parallel configuration faster. We make the complexity become $O(I(2|V|N + 2C_\theta))$

I : The number of iterations required for the optimization process to converge.

$|V|$: The number of vertices in the graph representation of the differentiable parallelization search space.

N : The number of parallelization strategies associated with each vertex in the graph.

C_θ : The complexity of computing the gradients with respect to the model parameters θ .

For demonstration purpose, our experiment uses partial code from Alpa's[28] dynamic programming space, to prove that our differentiable search space could potentially outperform the DP method.

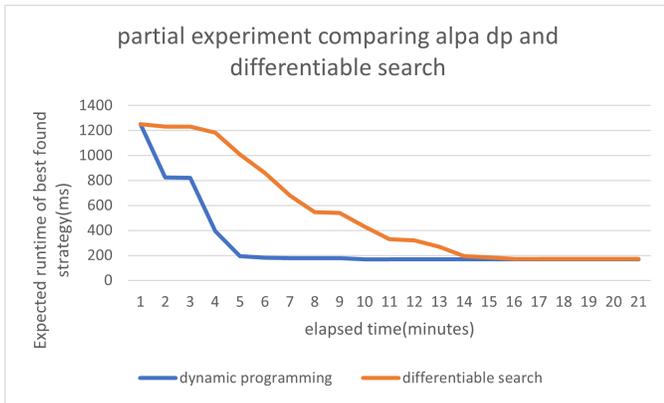


Fig. 9. Differentiable Search results vs Dynamic Programming

The advantage is that it will not like DP, have to cut the search space to make time complexity smaller. It may reach better global optimal. The disadvantage is that the gradient descent may cost more computational power. Its implementation is more complex, can be buggy when implementing the algorithm. Still, whether we will outperform dynamic programming is not clear. We will fully implement the algorithm to see the results in the future. Since the Alpa search space could be large when the number of node and the their computational

power is heterogeneous, there could be some cases when the differentiable search method could outperform the dynamic programming method in finding the sub-optimal solution. We will experiment it on more large cluster to see the results.

C. Parallel Random Initialization for Faster Loss Convergence

We use parallel random initialization to select inherently better initialization to make training loss converge faster. We compare different sampling methods like Single Random Initialization(original),Uniform Sampling. We partially implemented Latin hypercube sampling,Bayesian Optimization,Adaptive Sampling, and will test the results of them in the future.

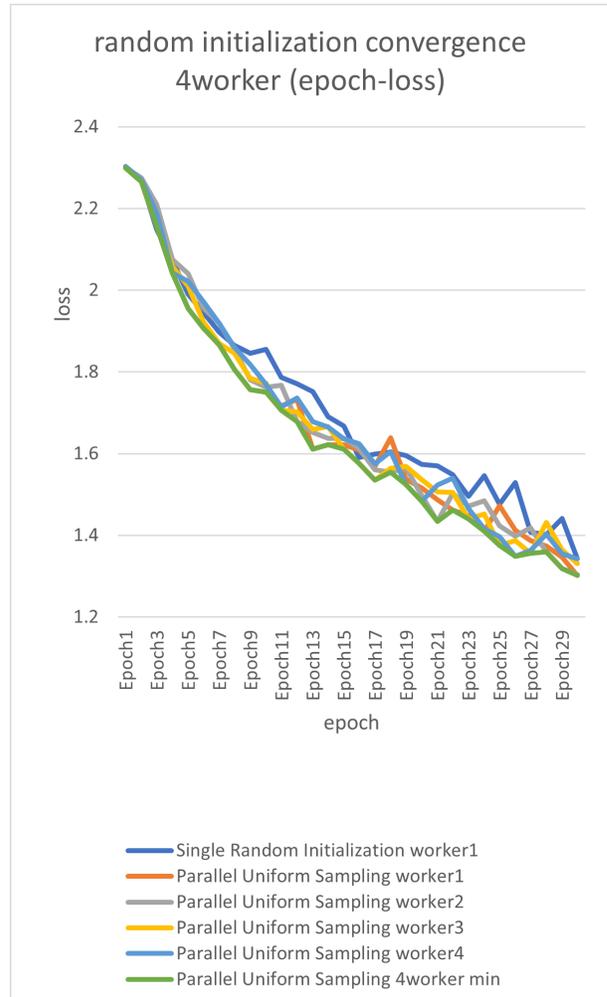


Fig. 10. Parallel Random Initialization Results 4 workers

From the figure and experimental results, when 40workers, the convergence speed can get around 1.5% improve. That's a small improve but it shows it may work if we improve our initialization or sampling method. In the future we will test different init method like kaiming_normal or xavier_normal, and different sampling methods like Latin hypercube sampling,Bayesian Optimization,Adaptive Sampling. The advantage is that if all methods is used up, this method

can drain the resource to gain faster loss convergence. The disadvantage is that it will cost many resources.

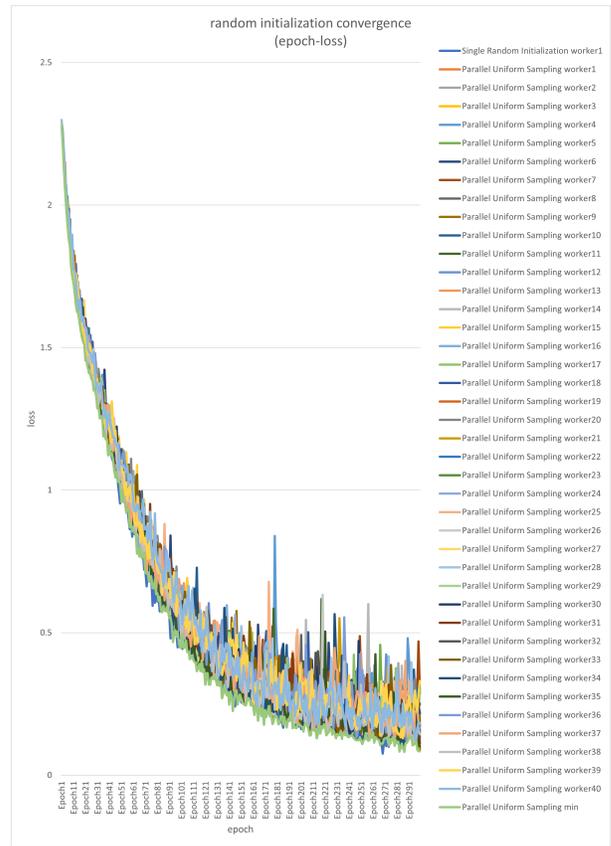


Fig. 12. Parallel Random Initialization Results 40 workers

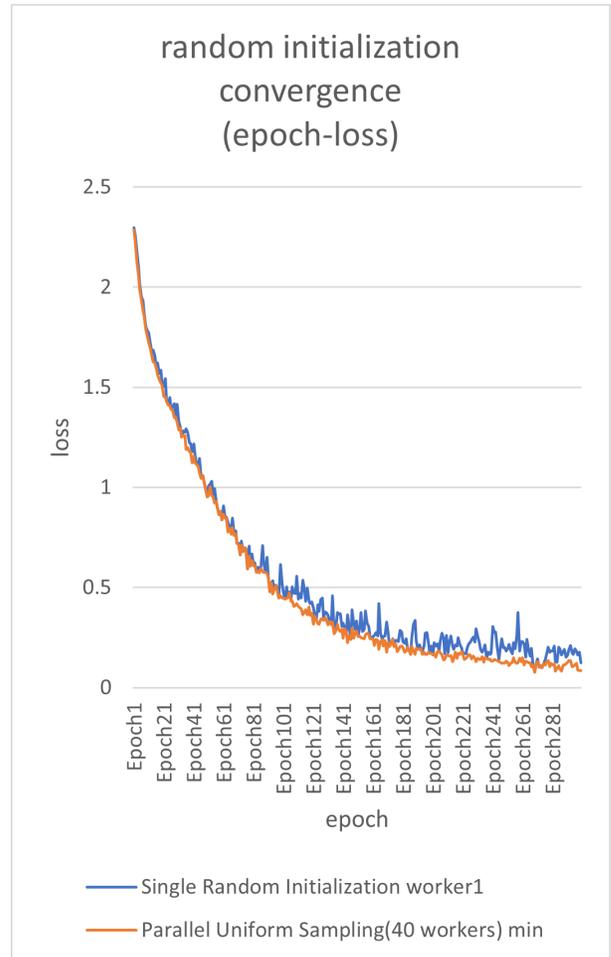


Fig. 13. Parallel Random Initialization Results 40 workers

VI. TIMELINE/EXPECTED MILESTONES FOR ACHIEVING OUR GOALS

4.1 Mar 15th: Get some initial experiment results that prove our motivations (detect the idle bubble in pipeline flow).

4.2 Mar, 30th: Figured out the concrete approaches for predictor & finish most part of mid-term report.

4.3 Apr, 15th: Implemented a simple heuristic binding scheduler (maybe only taken several simple factors like network) based on existing codebase like ColossalAI, Alpa[28]. Experiment Differentiable Search on small part of parallel configuration. Experiment Parallel Random Initialization.

4.4 Apr, 30th: Implemented the fully functional hardware aware binding algorithm between pipeline stages and computation meshes. Experiment Differentiable Search on larger part of parallel configuration. Experiment Parallel Random Initialization on different init strategy and random sampling algorithm.

4.5 May, 6th: Perform more training to evaluate our system and finish final report

APPENDIX A REFERENCES

[1] Sanjith Athlur, Nitika Saran, Muthian Sivathanu, Ramachandran Ramjee, and Nipun Kwatra. Varuna: scalable, low-cost training of massive deep learning models. In Proceedings of the Seventeenth European Conference on Computer Systems, pages 472-487, 2022.

[2] Paul Barham, Aakanksha Chowdhery, Jeff Dean, Sanjay Ghemawat, Steven Hand, Daniel Hurt, Michael Isard, Hyeontaek Lim, Ruoming Pang, Sudip Roy, et al. Pathways: Asynchronous distributed dataflow for ml. Proceedings of Machine Learning and Systems, 4, 2022.

[3] Zhenkun Cai, Xiao Yan, Kaihao Ma, Yidi Wu, Yuzhen Huang, James Cheng, Teng Su, and Fan Yu. Tensoropt: Exploring the tradeoffs in distributed dnn training with auto-parallelism. IEEE Transactions on Parallel and Distributed Systems, 33(8):1967-1981, 2021.

[4] Shiqing Fan, Yi Rong, Chen Meng, Zongyan Cao, Siyu Wang, Zhen Zheng, Chuan Wu, Guoping Long, Jun Yang, Lixue Xia, et al. Dapple: A pipelined data parallel approach for training large models. In Proceedings of the 26th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming, pages 431-445, 2021.

[5] Yanping Huang, Youlong Cheng, Ankur Bapna, Orhan Firat, Dehao Chen, Mia Chen, Hyounjoong Lee, Jiquan Ngiam, Quoc V Le, Yonghui Wu, et al. Gpipe: Efficient training of giant neural networks using pipeline parallelism. Advances in neural information processing systems, 32:103-112, 2019. [6] Zhihao Jia, Matei Zaharia, and Alex Aiken. Beyond data and model parallelism for deep neural networks. arXiv preprint arXiv:1807.05358, 2018.

[7] 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.

[8] Soojeong Kim, Gyeong-In Yu, Hojin Park, Sungwoo Cho, Eunji Jeong, Hyeonmin Ha, Sanha Lee, Joo Seong Jeong, and Byung-Gon Chun. Parallax: Sparsity-aware data parallel training of deep neural networks. In Proceedings of the Fourteenth EuroSys Conference 2019, pages 1-15, 2019.

[9] Dmitry Lepikhin, Hyounjoong Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping Huang, Maxim Krikun, Noam Shazeer, and Zhifeng Chen. Gshard: Scaling giant models with conditional computation and automatic sharding. arXiv preprint arXiv:2006.16668, 2020

[10] Shen Li, Yanli Zhao, Rohan Varma, Omkar Salpekar, Pieter Noordhuis, Teng Li, Adam Paszke, Jeff Smith, Brian Vaughan, Pritam Damania, et al. Pytorch distributed: Experiences on accelerating data parallel training. arXiv preprint arXiv:2006.15704, 2020.

[11] Zhuohan Li, Siyuan Zhuang, Shiyuan Guo, Danyang Zhuo, Hao Zhang, Dawn Song, and Ion Stoica. Terapipe: Token-level pipeline parallelism for training large-scale language models. arXiv preprint arXiv:2102.07988, 2021.

[12] Azalia Mirhoseini, Hieu Pham, Quoc V Le, Benoit Steiner, Rasmus Larsen, Yufeng Zhou, Naveen Kumar, Mohammad Norouzi, Samy Bengio, and Jeff Dean. Device placement optimization with reinforcement learning. In International Conference on Machine Learning, pages 2430-2439. PMLR, 2017.

[13] 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, pages 1 – 15, 2019.

[14] Deepak Narayanan, Amar Phanishayee, Kaiyu Shi, Xie Chen, and Matei Zaharia. Memoryefficient pipelineparallel dnn training. In International Conference on Machine Learning, pages 79377947. PMLR, 2021.

[15] Deepak Narayanan, Mohammad Shoeybi, Jared Casper, Patrick LeGresley, Mostofa Patwary, Vijay Korthikanti, Dmitri Vainbrand, Prethvi Kashinkunti, Julie Bernauer, Bryan Catanzaro, et al. Efficient large-scale language model training on gpu clusters using megatron-1m. In Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis, pages 1-15, 2021.

[16] 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, pages 16-29, 2019. [17] Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yuxiong He. Zero: Memory optimizations toward training trillion parameter models. In SC20: International Conference for High Performance Computing, Networking, Storage and Analysis, pages 1-16. IEEE, 2020. [18] Alexander Sergeev and Mike Del Balso. Horovod: fast and easy distributed deep learning in tensorflow. arXiv preprint arXiv:1802.05799, 2018.

[19] Noam Shazeer, Youlong Cheng, Niki Parmar, Dustin Tran, Ashish Vaswani, Penporn Koanantakool, Peter Hawkins,

HyoukJoong Lee, Mingsheng Hong, Cliff Young, et al. Meshtensorflow: Deep learning for supercomputers. arXiv preprint arXiv:1811.02084, 2018.

[20] Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan Catanzaro. Megatron-Im: Training multi-billion parameter language models using model parallelism. arXiv preprint arXiv:1909.08053, 2019.

[21] Jakub M Tarnawski, Deepak Narayanan, and Amar Phanishayee. Piper: Multidimensional planner for dnn parallelization. Advances in Neural Information Processing Systems, 34, 2021. [22] Minjie Wang, Chien-chin Huang, and Jinyang Li. Supporting very large models using automatic dataflow graph partitioning. In Proceedings of the Fourteenth EuroSys Conference 2019, pages 1-17, 2019.

[23] Yuanzhong Xu, HyoukJoong Lee, Dehao Chen, Hongjun Choi, Blake Hechtman, and Shibo Wang. Automatic cross-replica sharding of weight update in data-parallel training. arXiv preprint arXiv:2004.13336, 2020.

[24] Yuanzhong Xu, HyoukJoong Lee, Dehao Chen, Blake Hechtman, Yanping Huang, Rahul Joshi, Maxim Krikun, Dmitry Lepikhin, Andy Ly, Marcello Maggioni, et al. Gspmd: General and scalable parallelization for ml computation graphs. arXiv preprint arXiv:2105.04663, 2021

[25] Jinhui Yuan, Xinqi Li, Cheng Cheng, Juncheng Liu, Ran Guo, Shenghang Cai, Chi Yao, Fei Yang, Xiaodong Yi, Chuan Wu, et al. Oneflow: Redesign the distributed deep learning framework from scratch. arXiv preprint arXiv:2110.15032, 2021.

[26] Hao Zhang, Yuan Li, Zhijie Deng, Xiaodan Liang, Lawrence Carin, and Eric Xing. Autosync: Learning to synchronize for data-parallel distributed deep learning. Advances in Neural Information Processing Systems, 33, 2020.

[27] Zhen Zhang, Shuai Zheng, Yida Wang, Justin Chiu, George Karypis, Trishul Chilimbi, Mu Li, and Xin Jin. Mics: Near-linear scaling for training gigantic model on public cloud. arXiv preprint arXiv:2205.00119, 2022.

[28]Zheng, Lianmin, Zhuohan Li, Hao Zhang, Yonghao Zhuang, Zhifeng Chen, Yanping Huang, Yida Wang et al. "Alpa: Automating Inter-and Intra-Operator Parallelism for Distributed Deep Learning." In 16th USENIX Symposium on Operating Systems Design and Implementation (OSDI 22), pp. 559578. 2022.

[29]<https://www.deepspeed.ai/tutorials/pipeline/>

[30]<https://huggingface.co/docs/transformers/v4.15.0/parallelism>

[31]<https://www.usenix.org/system/files/osdi22-zheng-lianmin.pdf>

[32]Rapp, Martin, Ramin Khalili, Kilian Pfeiffer, and Jörg Henkel. "Distreal: Distributed resource-aware learning in heterogeneous systems." In Proceedings of the AAAI Conference on Artificial Intelligence, vol. 36, no. 7, pp. 8062-8071. 2022. <https://arxiv.org/abs/2112.08761v1>

[34]Amaris, Marcos, Raphael Camargo, Daniel Cordeiro, Alfredo Goldman, and Denis Trystram. "Evaluating execution time predictions on GPU kernels using an analytical model and machine learning techniques." Journal

of Parallel and Distributed Computing 171 (2023): 66-78.

https://www.sciencedirect.com/science/article/pii/S0743731522001903?casa_token=Kqi7243EJsAAAAA:P9NoXtKwokKF11fN8OgGYys5r0gaFgVweA7ZyMfYvSMkUkGQ37oMASLi49cnXp0F9Oogk3Fg

[35]Pham, Thanh-Phuong, Juan J. Durillo, and Thomas Fahringer. "Predicting workflow task execution time in the cloud using a two-stage machine learning approach." IEEE Transactions on Cloud Computing 8, no. 1 (2017): 256-268.

[36]Liu, Hanxiao, Karen Simonyan, and Yiming Yang. "Darts: Differentiable architecture search." arXiv preprint arXiv:1806.09055 (2018).

[37]T. -P. Pham, J. J. Durillo and T. Fahringer, "Predicting Workflow Task Execution Time in the Cloud Using A Two-Stage Machine Learning Approach," in IEEE Transactions on Cloud Computing, vol. 8, no. 1, pp. 256-268, 1 Jan.-March 2020, doi: 10.1109/TCC.2017.2732344.

APPENDIX B CODE

The source code of this project is at <https://github.com/explrc/pipeDejavu>

APPENDIX C DISTRIBUTED TRAINING FAULT TOLERANCE ALGORITHM

Algorithm 2 Distributed Training Fault Tolerance Algorithm

- 1: Initialize the weights of each machine $w_i = \frac{1}{N}$
 - 2: **while** not converged **do**
 - 3: Compute local gradients g_i on each machine
 - 4: Compute overall gradient $g = \sum_{i=1}^N w_i g_i$
 - 5: Compute suspicion level of each machine $s_i = |g_i - g|$
 - 6: Update the weight of each machine $w_i = \frac{\alpha}{\alpha + s_i}$
 - 7: Sort machines based on their weight in ascending order
 - 8: Identify the k machines with the least influence on the overall gradient
 - 9: Replace gradients of identified machines with the average of remaining gradients $\hat{g}_i = \frac{\sum_{j=1}^{N-k} g_j}{N-k}$
 - 10: Compute overall gradient $g = \sum_{i=1}^N w_i \hat{g}_i$
 - 11: Update model parameters using the overall gradient
 - 12: **end while**
-

APPENDIX D
IMAGES OF DIFERENT PARALLEL PARADIGM

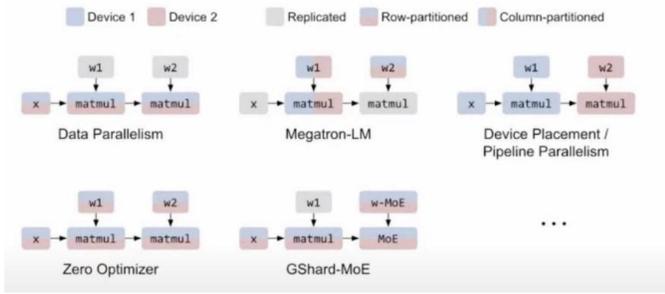


Fig. 14. Existing Parallelization Techniques

Figure 2: Common parallelization techniques for training a 2-layer Multi-layer Perceptron (MLP). Only the forward pass is shown. " x " is the input data. " w_1 " and " w_2 " are two weight matrices.

APPENDIX E
LOGIC FLOW OF HOW TO PARALLELIZE A COMPUTATIONAL GRAPH

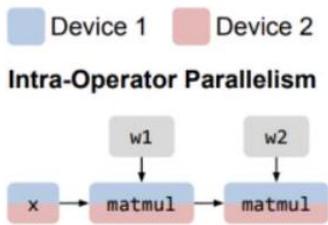


Fig. 15. (a) Data Parallelism

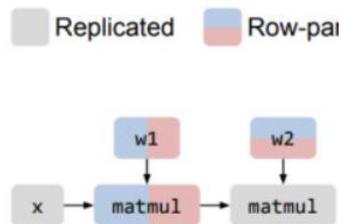


Fig. 16. (b) Operator Parallelism

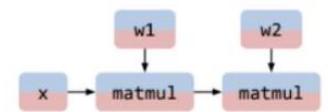


Fig. 17. (c) ZeRO Optimizer

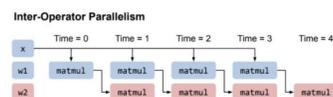


Fig. 18. (d) Pipeline Parallelism

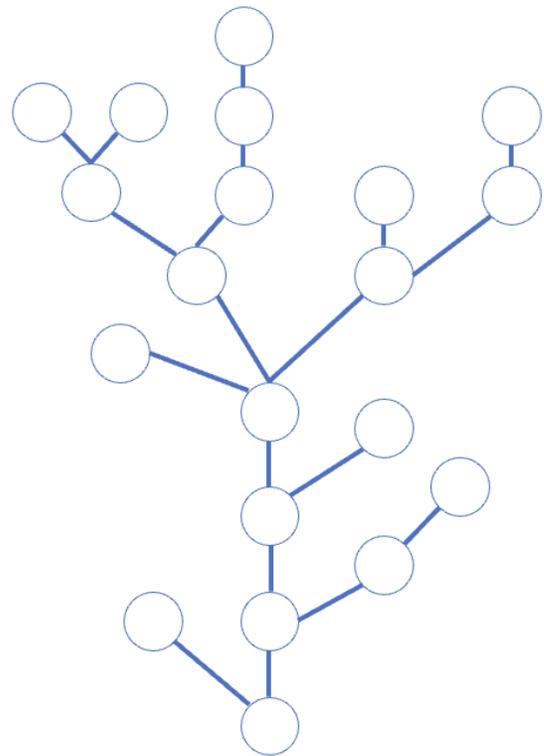


Fig. 18. Step 1: Computation Graph when node number $n = 4$

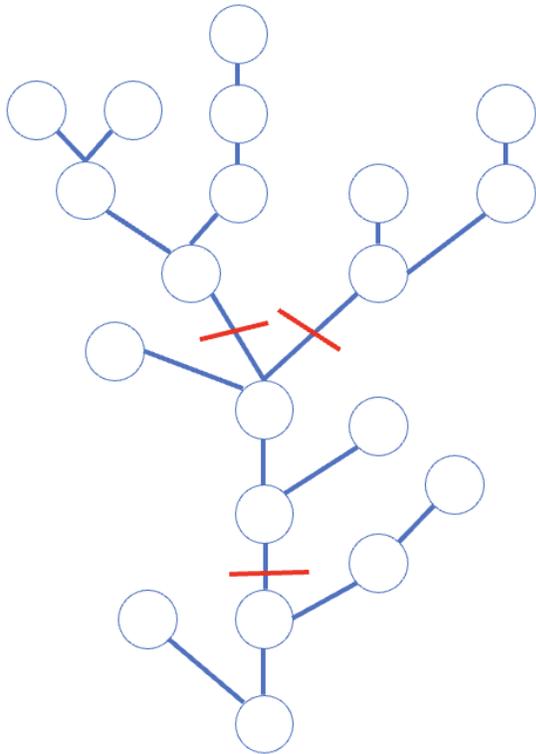


Fig. 19. Step 2: Cut $n-1$ of this computation graph

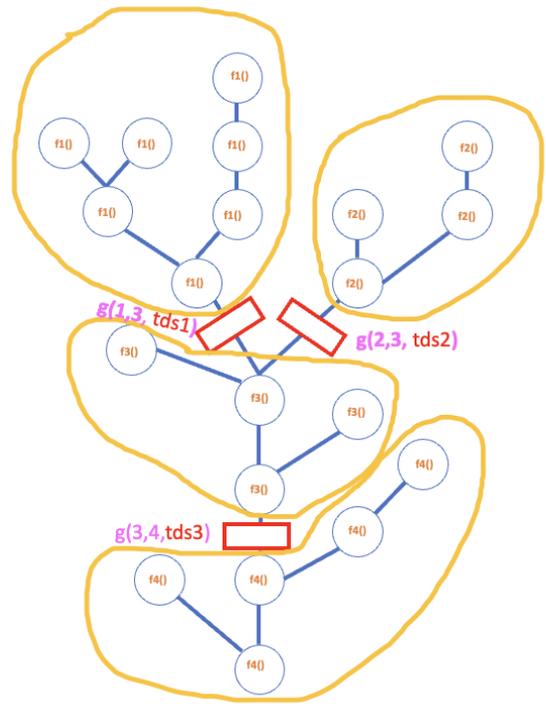


Fig. 21. Step 4: Apply hardware cost function

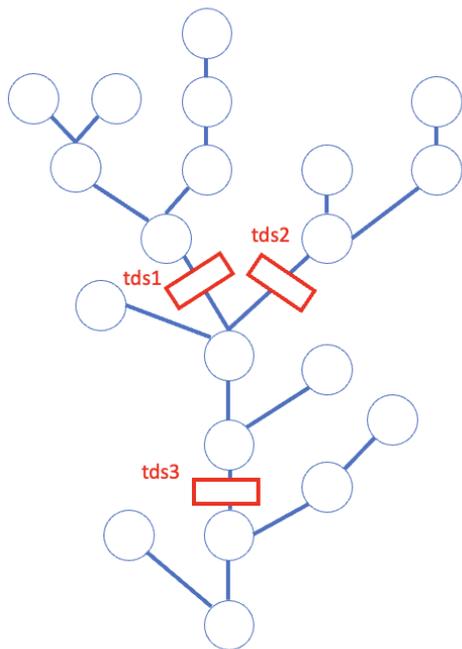


Fig. 20. Step 3: Add communication cost

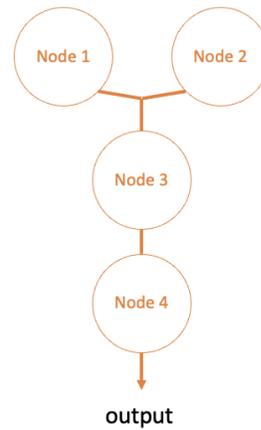


Fig. 22. Step 5: Abstract to n nodes

APPENDIX F

CODE EXAMPLE IN FITTING PROFILING DATABASE PICKLE FILE USING LINEAR REGRESSION

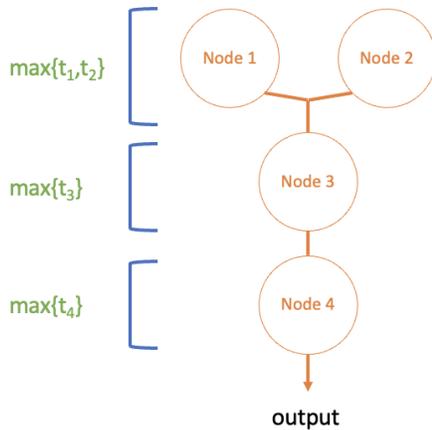


Fig. 23. Step 6: Apply MAX operator to each layer of the tree

Listing 1. Code Example in fitting profiling database pickle file using linear

```

1 regression
2 import pickle
3 import numpy as np
4 import matplotlib.pyplot as plt
5 from sklearn.linear_model import
6   LinearRegression
7 from sklearn.metrics import r2_score
8 from decimal import Decimal
9 from sklearn.preprocessing import
10   StandardScaler
11 # Load the provided data
12 with open("prof_database.pkl", "rb") as f:
13     data = pickle.load(f)
14
15 # Function to extract X and Y values from
16 # the dataset
17 def extract_data(dataset, key):
18     x = [point[0] for point in dataset[key
19         ]]
20     #x = [Decimal(point[0]) for point in
21         dataset[key]]
22     y = [point[1] for point in dataset[key
23         ]]
24     #y = [Decimal(point[1]) for point in
25         dataset[key]]
26
27     # Replace extreme values with the
28     # maximum finite representable value
29     # for float64
30     x = np.clip(x, np.finfo(np.float64).min
31         , np.finfo(np.float64).max)
32     y = np.clip(y, np.finfo(np.float64).min
33         , np.finfo(np.float64).max)
34
35     # Replace NaN values with the mean of
36     # the non-NaN elements in the array
37     y = np.where(np.isnan(y), np.nanmean(y)
38         , y)
39
40     return np.array(x).reshape(-1, 1), np.
41     array(y)
42
43 # Function to filter out infinity or large
44 # values from X and y
45 def filter_data(X, y):
46     X_flat = np.ravel(X)
47     overflow_mask = (X_flat < np.finfo(np.
48         float32).max) & (y < np.finfo(np.
49         float32).max)
50     return X[overflow_mask], y[
51         overflow_mask]
52
53 # Function to apply log transformation to y
54 # values
55 def apply_log_transform(y):
56     return np.log(y)
57
58 # List of attributes
59 attributes = [
60     'all_gather_cost_dict',
61     'all_reduce_cost_dict',
62     'all_to_all_cost_dict',
  
```

```

43     'reduce_scatter_cost_dict',
44     'available_memory_per_device',
45     'dot_cost_dict',
46     'conv_cost_dict',
47     'op_cost_dict',
48 ]
49
50 # Loop through all keys
51 for key in data.keys():
52     # Loop through all attributes
53     for attr in attributes:
54         attribute_dict = getattr(data[key],
55                                 attr)
56
57         if not isinstance(attribute_dict,
58                           dict):
59             continue
60
61         for config, cost_data in
62             attribute_dict.items():
63             X, y = extract_data(
64                 attribute_dict, config)
65
66             # Filter out infinity or large
67             # values from X and y
68             X, y = filter_data(X, y)
69             # Apply log transformation to
70             # y values
71             y = np.where(np.isnan(y), np.
72                         nanmean(y), y)
73             # Initialize the scaler
74             scaler = StandardScaler()
75             if X.size == 0 or y.size == 0:
76                 print(f"Empty arrays
77                       encountered for {key},
78                       {attr}, {config}.
79                       Skipping...")
80                 continue
81
82             y_norm = scaler.fit_transform(y
83                                         .reshape(-1, 1)).reshape
84                                         (-1)
85             # Linear Regression
86             lr = LinearRegression()
87             lr.fit(X, y_norm)
88             y_pred_norm = lr.predict(X)
89
90             # Rescale predictions back to
91             # original scale
92             y_pred = scaler.
93                 inverse_transform(
94                     y_pred_norm.reshape(-1, 1))
95                 .reshape(-1)
96             r2 = r2_score(y, y_pred)
97             '''
98             # Linear Regression
99             lr = LinearRegression()
100            lr.fit(X, y)
101            y_pred = lr.predict(X)
102            r2 = r2_score(y, y_pred)
103            '''
104            # Print accuracy results
105            print(f"Key: {key}")
106            print(f"Attribute: {attr}")
107            print(f"Configuration: {config}
108                  ")
109            print(f"R2 score: {r2:.2f}")

```

```

93     print(f"Slope: {lr.coef_[0]}")
94     print(f"Intercept: {lr.
95           intercept_}\n")
96
97     # Visualization
98     plt.scatter(X, y, label=f"{
99           config} R2: {r2:.2f}")
100     plt.plot(X, y_pred)
101
102     plt.xscale("log")
103     plt.yscale("log")
104     plt.xlabel("Number of Parameters")
105     plt.ylabel("Cost")
106
107     # Move the legend outside of the
108     # plot
109     plt.legend(bbox_to_anchor=(1.05, 1)
110               , loc='upper left',
111               borderaxespad=0.)
112
113     plt.title(f"{attr.capitalize()} vs
114              Ranks for Different
115              Configurations ({key})")
116
117     # Save the plot
118     plt.savefig(f"{attr}_vs_ranks_{key}
119               }.png", bbox_inches='tight')
120
121     # Show the plot
122     plt.show()
123
124     # Clear the plot for the next
125     # attribute
126     plt.clf()

```

APPENDIX G

PREDICTIVE MODEL RESULTS FOR PRE-PROFILING STAGE

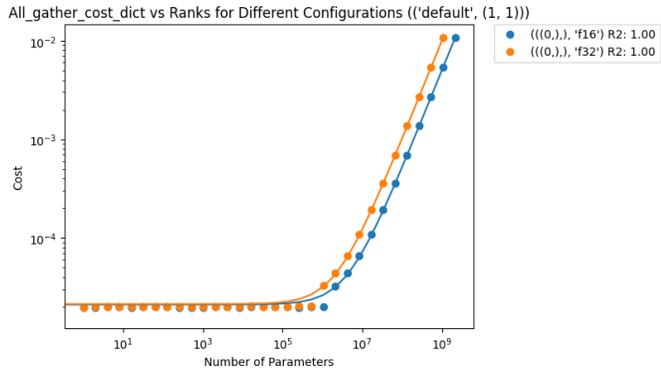


Fig. 24.

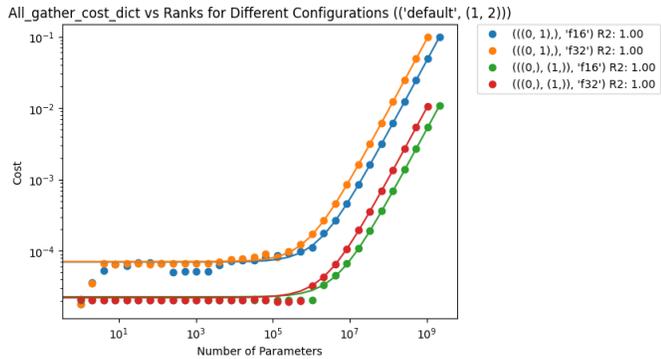


Fig. 25.

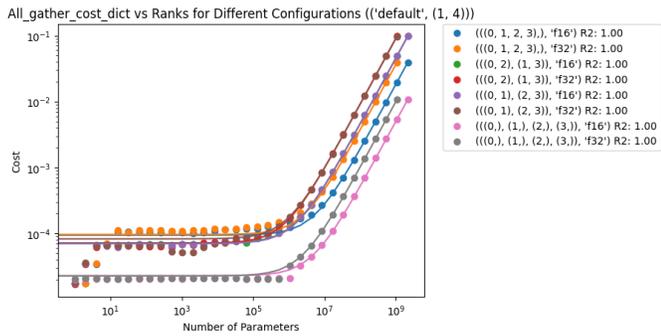


Fig. 26.

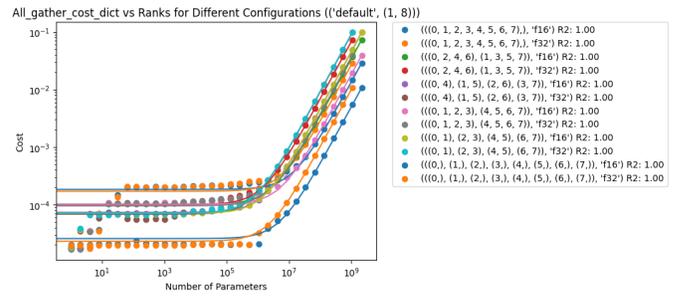


Fig. 27.

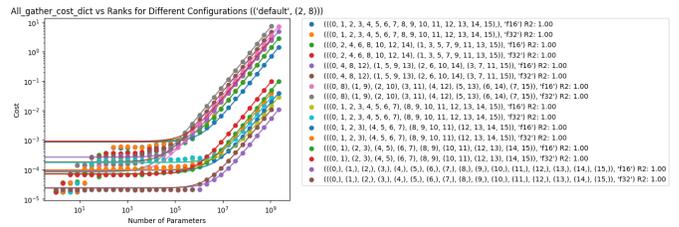


Fig. 28.



Fig. 29.

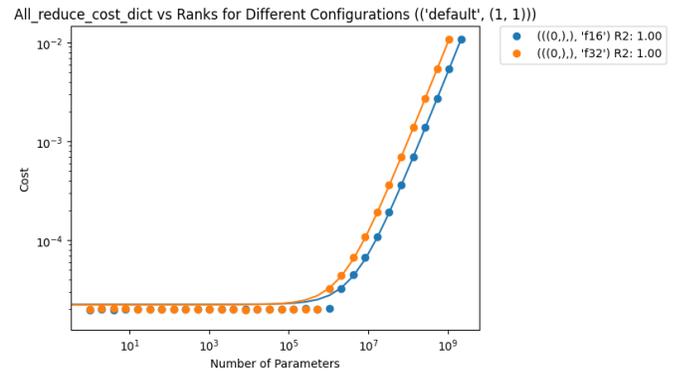


Fig. 30.

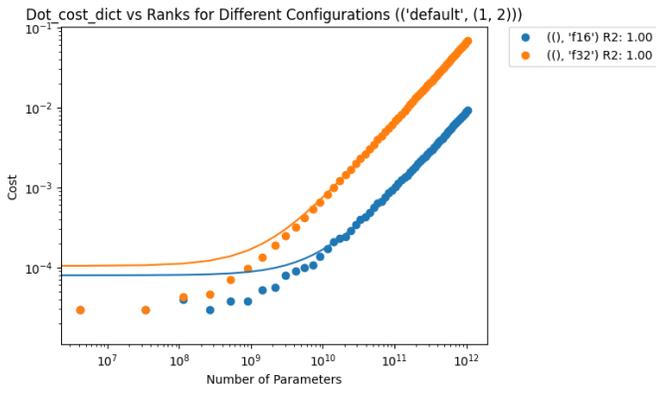


Fig. 31.



Fig. 35.

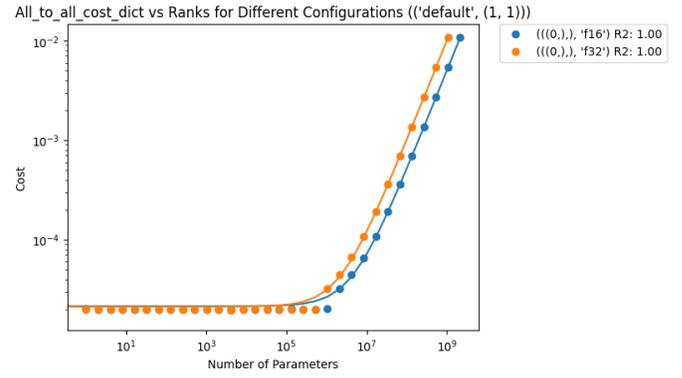


Fig. 36.

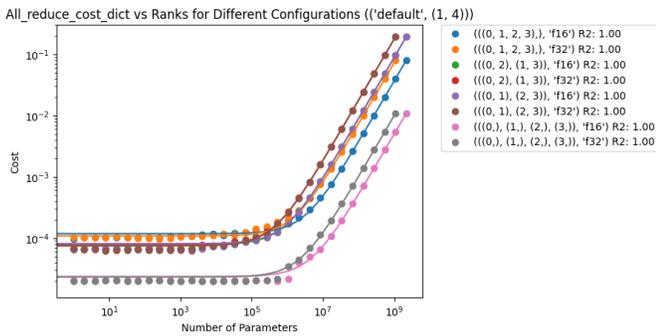


Fig. 32.

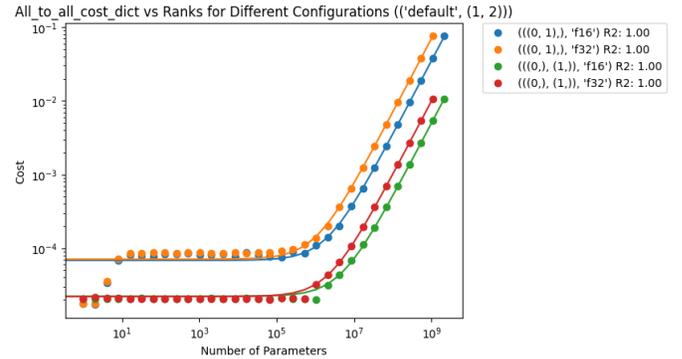


Fig. 37.

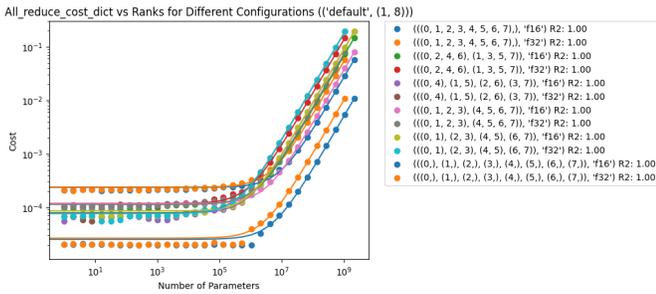


Fig. 33.

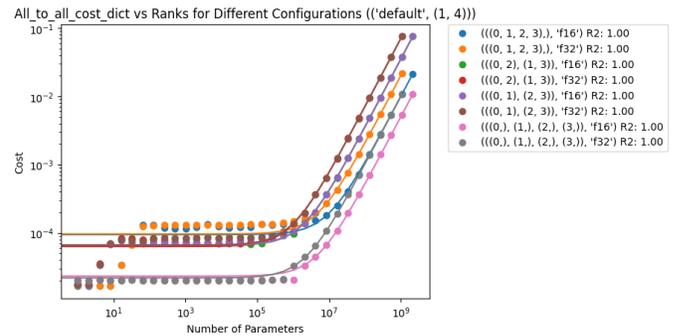


Fig. 38.

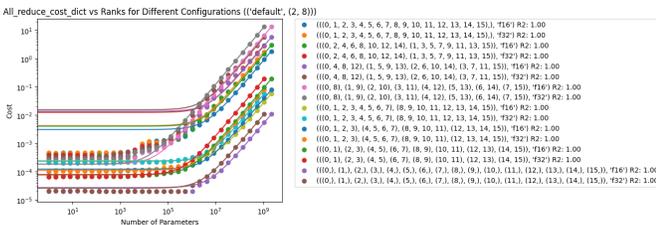


Fig. 34.

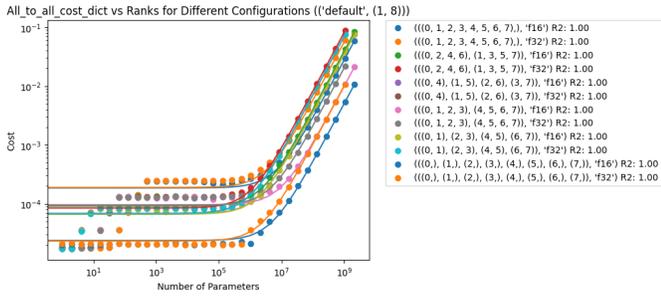


Fig. 39.

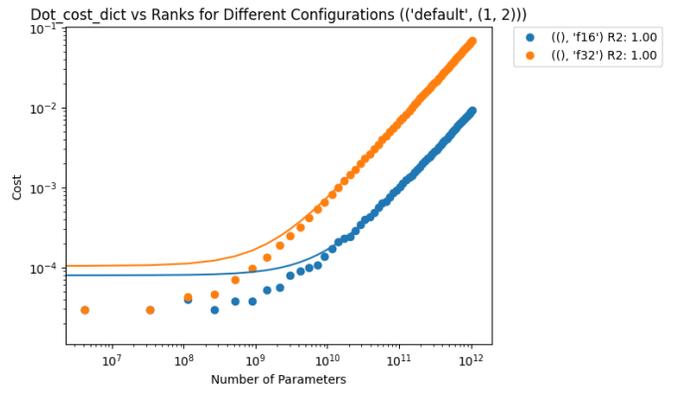


Fig. 43.

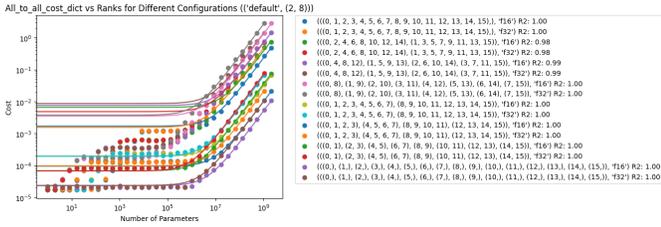


Fig. 40.

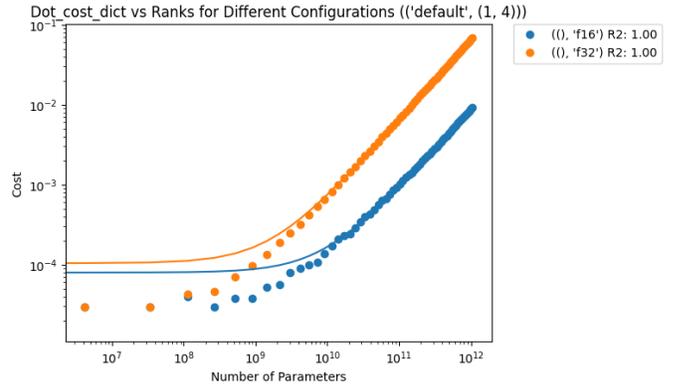


Fig. 44.

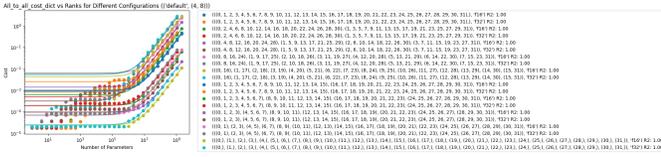


Fig. 41.

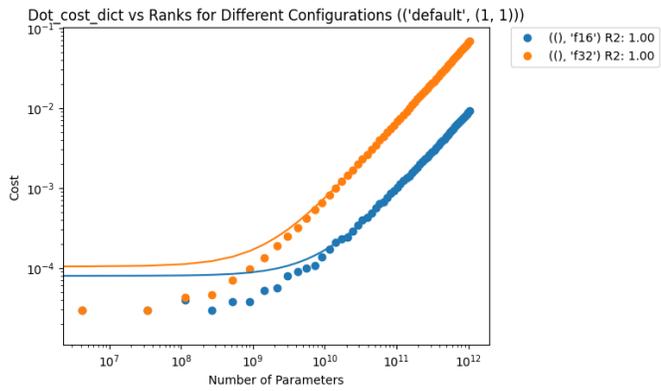


Fig. 42.

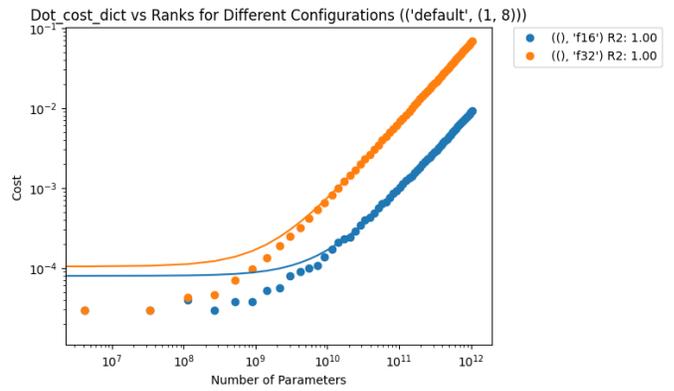


Fig. 45.

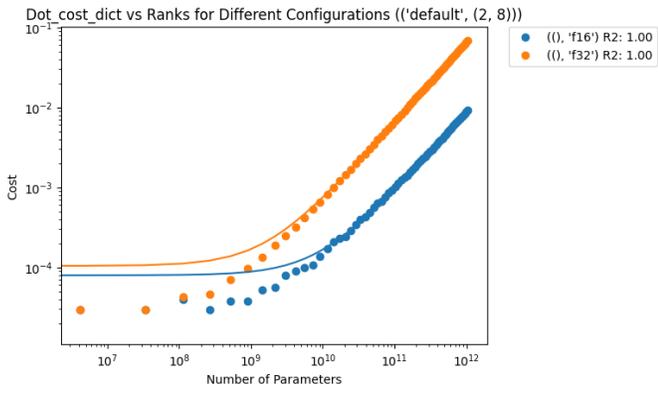


Fig. 46.

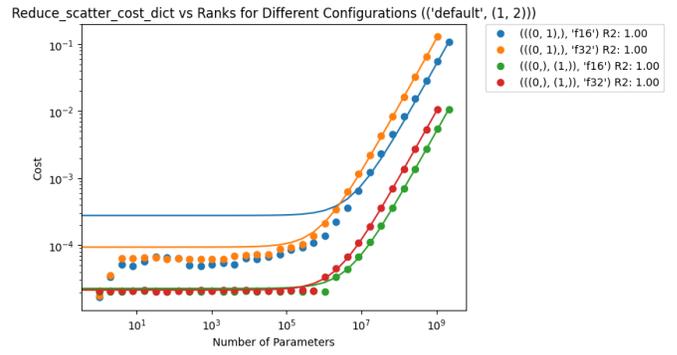


Fig. 49.

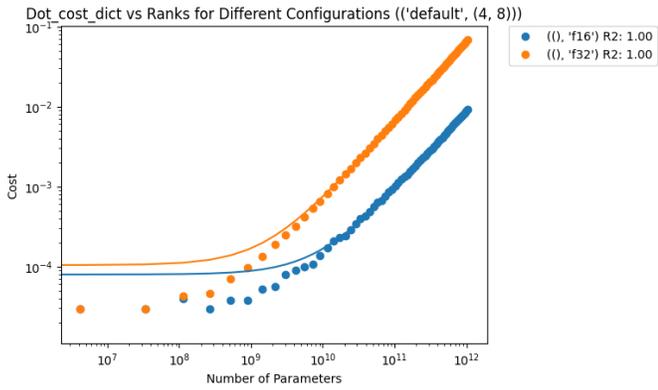


Fig. 47.

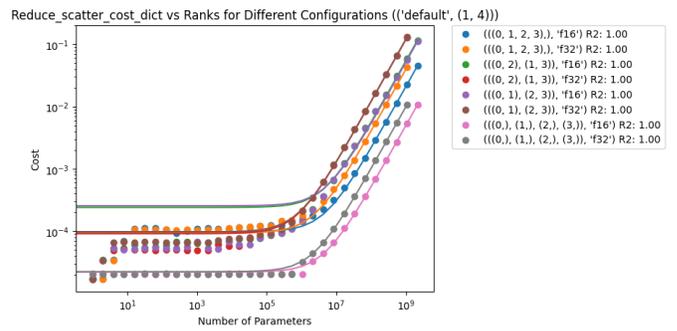


Fig. 50.

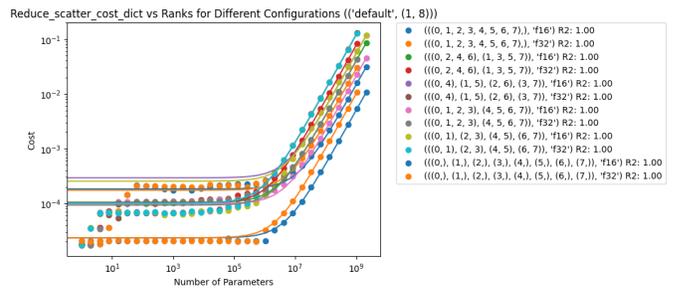


Fig. 51.

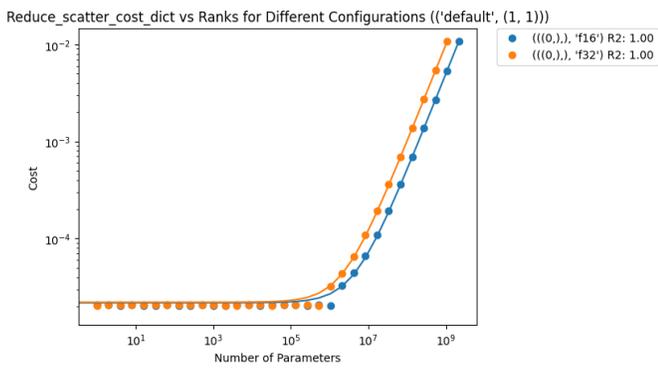


Fig. 48.

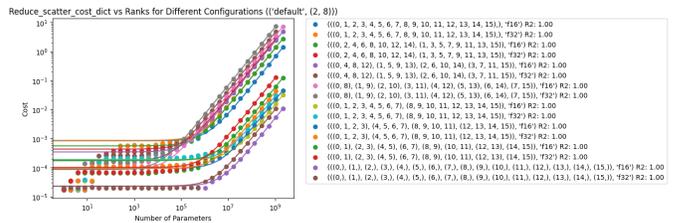


Fig. 52.



Fig. 53.

APPENDIX H EXAMPLE CODE FOR DP AND DIFFERENTIABLE SEARCH SPACE FOR CLASSIC KNAPSACK PROBLEM

Listing 2. Dynamic Programming Knapsack Algorithm

```

1 def knapsack(items, capacity, print_every
  =500, timer_every=10):
2     n = len(items)
3     dp = [[0 for _ in range(capacity + 1)]
4           for _ in range(n + 1)]
5
6     best_solutions_dp = []
7
8     start_time = time.time()
9     for i in range(1, n + 1):
10        weight, value = items[i - 1]
11        for w in range(capacity + 1):
12            if weight <= w:
13                dp[i][w] = max(dp[i - 1][w
14                               ], dp[i - 1][w - weight
15                               ] + value)
16            else:
17                dp[i][w] = dp[i - 1][w]
18
19        if i % timer_every == 0:
20            elapsed_time = time.time() -
21                start_time
22            best_solutions_dp.append((
23                elapsed_time, dp[i][-1]))
24        if i % print_every == 0:
25            print(f"Iteration {i}: Best
26                  solution so far for DP:
27                  value={dp[i][-1]},
28                  elapsed_time={elapsed_time
29                  :.2f}s")
30
31        selected_items = []
32        i, w = n, capacity
33        while i > 0 and w > 0:
34            weight, value = items[i - 1]
35            if dp[i][w] != dp[i - 1][w]:
36                selected_items.append(i - 1)
37                w -= weight
38            i -= 1
39
40    return dp[n][capacity], selected_items,
41        best_solutions_dp

```

Listing 3. Differentiable Dynamic Programming Knapsack Algorithm

```

1  from collections import namedtuple
2
3  Solution = namedtuple("Solution", ["value",
4     "items", "selection", "elapsed_time"])
5
6  def soft_knapsack(items, capacity, n_best
7     =10000, iterations=10000, learning_rate
8     =100, temperature=0.5, print_every=500,
9     penalty_factor=1e7):
10     n_best=iterations
11     n = len(items)
12     weights, values = zip(*items)
13     weights = torch.tensor(weights, dtype=
14         torch.float)
15     values = torch.tensor(values, dtype=
16         torch.float)
17
18     item_selection = torch.rand(n,
19         requires_grad=True)
20     optimizer = optim.RMSprop([
21         item_selection], lr=learning_rate)
22
23     best_solutions = []
24
25     start_time = time.time()
26     for i in range(iterations):
27         optimizer.zero_grad()
28         soft_selection = torch.sigmoid(
29             item_selection / temperature)
30         total_weight = torch.sum(
31             soft_selection * weights)
32         total_value = torch.sum(
33             soft_selection * values)
34         capacity_penalty = torch.clamp(
35             total_weight - capacity, min=0)
36             ** 2
37
38         # Multiply capacity_penalty by a
39         large constant
40         loss = -(total_value -
41             penalty_factor *
42             capacity_penalty)
43         loss.backward()
44         optimizer.step()
45
46         # Clamp item_selection values
47         between -5 and 5
48         item_selection.data.clamp_(-5, 5)
49         if total_weight <= capacity:
50             # Update best solutions
51             final_selection = torch.sigmoid(
52                 item_selection / temperature)
53                 > 0.5
54             selected_items = [i for i,
55                 selected in enumerate(
56                     final_selection) if selected]
57             max_value = torch.sum(
58                 final_selection * values).
59                 item()
60             current_solution = Solution(
61                 max_value, selected_items,
62                 final_selection, time.time() -
63                 start_time)
64             #best_solution = best_solutions
65             [0]

```

```

39     if len(best_solutions) < n_best:
40         best_solutions.append(
41             current_solution)
42         best_solutions.sort(key=
43             lambda x: x.value,
44             reverse=True)
45     elif max_value > best_solutions
46         [-1].value:
47         best_solutions.pop()
48         best_solutions.append(
49             current_solution)
50         best_solutions.sort(key=
51             lambda x: x.value,
52             reverse=True)
53         # Print intermediate results
54     if (i + 1) % print_every == 0:
55         elapsed_time = time.time() -
56             start_time
57         best_solution = best_solutions
58             [0]
59         print(f"Iteration {i+1}: loss={
60             loss.item():.2f},
61             total_value={total_value.
62                 item():.2f}, total_weight={
63                     total_weight.item():.2f},
64                     elapsed_time={elapsed_time
65                         :.2f}s")
66     print(f"Best solution so far:
67         value={best_solution.value
68             }, items={best_solution.
69                 items}")
70
71     max_value = best_solutions[0].value
72     selected_items = best_solutions[0].
73         items
74
75     best_solutions_diff = sorted([(s.
76         elapsed_time, s.value) for s in
77         best_solutions], key=lambda x: x
78         [0])
79     #best_solutions_diff = [(s.elapsed_time
80         , s.value) for s in best_solutions]
81
82     return max_value, selected_items,
83         best_solutions_diff

```

APPENDIX I
**EXAMPLE CODE FOR PARALLEL RANDOM
 INITIALIZATION SIMULATION PROGRAM**

Listing 4. Example Code for Parallel Random Initialization Simulation Program

```

1  import torch
2  import torch.nn as nn
3  import torch.optim as optim
4  from torch.utils.data import DataLoader
5  from torchvision import datasets,
   transforms
6  import numpy as np
7  import random
8  import itertools
9
10 import pyDOE2
11 from scipy.optimize import minimize
12 from skopt import gp_minimize
13
14 import torchvision
15 import torchvision.transforms as transforms
16 import torch.nn.functional as F
17 import matplotlib.pyplot as plt
18 import os
19 from tqdm import tqdm
20 import datetime
21 #from datetime import datetime
22 import csv
23 from datetime import datetime
24
25 results_dir = './results'
26 def nowTimetoString():
27     now=datetime.now()
28     return now.strftime("%m-%d-%Y_%H-%M-%S"
29 )
30 class SimpleNN(nn.Module):
31     def __init__(self, to_demo=True):
32         super(SimpleNN, self).__init__()
33         if to_demo:
34             self.conv1 = nn.Conv2d(3, 16,
35                                     3, padding=1)
36             self.conv2 = nn.Conv2d(16, 32,
37                                     3, padding=1)
38             self.fc1 = nn.Linear(32 * 8 *
39                                   8, 64)
40             self.fc2 = nn.Linear(64, 10)
41         else:
42             self.conv1 = nn.Conv2d(3, 16,
43                                     3, padding=1)
44             self.conv2 = nn.Conv2d(16, 32,
45                                     3, padding=1)
46             self.fc1 = nn.Linear(32 * 8 *
47                                   8, 64)
48             self.fc2 = nn.Linear(64, 10)
49
50         self.relu1 = nn.ReLU(inplace=True)
51         self.pool = nn.MaxPool2d(2, 2)
52         self.relu2 = nn.ReLU(inplace=True)
53         self.relu3 = nn.ReLU(inplace=True)
54
55     def forward(self, x):
56         x = self.conv1(x)
57         x = self.relu1(x)
58         x = self.pool(x)
59         x = self.relu2(x)
60         x = self.relu3(x)
61         return x
62
63 def init_weights(model, init_values):
64     for i, param in enumerate(model.
65 parameters()):
66         param.data.copy_(torch.from_numpy(
67             init_values[i]))
68
69         #param.data.copy_(torch.tensor(
70             init_values[i]))
71
72 # SimpleNN and init_weights remain the same
73 # as before
74 def simulate_parallel_loss(model,
75 init_values, train_loader, device,
76 num_epochs=10):
77     model.to(device)
78     init_weights(model, init_values)
79     criterion = nn.CrossEntropyLoss()
80     optimizer = optim.SGD(model.parameters
81                             (), lr=0.01, momentum=0.9)
82
83     loss_curve = []
84     for epoch in tqdm(range(num_epochs),
85 desc="Training"):#range(num_epochs)
86         :
87         running_loss = 0.0
88         for data, target in train_loader:
89             data, target = data.to(device),
90 target.to(device)
91             optimizer.zero_grad()
92             output = model(data)
93             loss = criterion(output, target
94                             )
95             loss.backward()
96             optimizer.step()
97             running_loss += loss.item()
98
99         loss_curve.append(running_loss /
100 len(train_loader))
101     return loss_curve
102
103 #originally work version
104 def run_simulation(num_workers,init_values,
105 train_loader, device, num_epochs=10,
106 to_demo=True):#original 2nd argument:
107     sampling_method
108     losses = []
109     for worker_id in range(num_workers):
110         #init_values = sampling_method()
111         loss_curve = simulate_parallel_loss
112 (SimpleNN(to_demo=to_demo),
113         init_values, train_loader,
114         device, num_epochs)
115         losses.append(loss_curve)
116     return losses

```

```

103 def to_csv(data,num_epochs):
104     if not os.path.exists(results_dir):
105         os.makedirs(results_dir)
106     nowTimeString=nowTimetoString()
107
108     with open(os.path.join(results_dir,
109         nowTimeString+'table1.csv'), 'w',
110         newline='') as file:
111         writer = csv.writer(file)
112         header = ['Method']#, 'Iteration
113             1', 'Iteration 2', 'Iteration
114             3', 'Iteration 4', 'Iteration
115             5', 'Iteration 6', 'Iteration
116             7', 'Iteration 8', 'Iteration
117             9', 'Iteration 10', 'Iteration
118             11', 'Iteration 12', 'Iteration
119             13', 'Iteration 14', '
120             Iteration 15', 'Iteration 16',
121             'Iteration 17', 'Iteration 18',
122             'Iteration 19', 'Iteration
123             20', 'Iteration 21', 'Iteration
124             22', 'Iteration 23', '
125             Iteration 24', 'Iteration 25',
126             'Iteration 26', 'Iteration 27',
127             'Iteration 28', 'Iteration
128             29', 'Iteration 30']
129         for i in range(num_epochs):
130             header.append("Epoch"+str(i+1))
131
132         writer.writerow(header)
133
134         for method in data:
135             for i in range(len(data[method
136                 ])):
137                 row = [method+" worker"+str
138                     (i+1)]
139                 row += [str(num) for num in
140                     data[method][i]]
141                 writer.writerow(row)
142
143 def main(to_demo=True,NUM_WORKERS_=20,
144     EPOCHS=100,DEMO_EPOCHS=50):
145     # Prepare the dataset
146     device = torch.device('cuda' if torch.
147         cuda.is_available() else 'cpu')
148     #device = torch.device('cuda:0' if
149         torch.cuda.is_available() else 'cpu
150     ')
151     #transform = transforms.Compose([
152         transforms.ToTensor(), transforms.
153         Normalize((0.5, 0.5, 0.5), (0.5,
154         0.5, 0.5))]
155     )
156     transform = transforms.Compose([
157         transforms.RandomHorizontalFlip(),
158         transforms.RandomCrop(32, padding
159         =4),
160         transforms.ToTensor(),
161         transforms.Normalize((0.5, 0.5,
162         0.5), (0.5, 0.5, 0.5))
163     ])
164     train_dataset = torchvision.datasets.
165         CIFAR10(root='./data', train=True,
166         download=True, transform=transform)
167     if to_demo:
168         train_dataset, _ = torch.utils.data
169         .random_split(train_dataset,
170             [1000, len(train_dataset) -
171             1000])
172     train_loader = DataLoader(train_dataset
173         , batch_size=32, shuffle=True,
174         num_workers=2)#, num_workers=2
175     #NUM_WORKERS_= 4
176     num_workers = NUM_WORKERS_
177
178     num_epochs = DEMO_EPOCHS if to_demo
179         else EPOCHS #30
180     # Define sampling methods
181     def single_random_initialization():
182         model = SimpleNN()
183         init_values = [p.data.clone().numpy
184             () for p in model.parameters()]
185         # originally no .numpy()
186         return init_values
187
188     def uniform_sampling():
189         #return [np.random.uniform(0, 1, p.
190             numel()).reshape(p.shape) for p
191             in SimpleNN().parameters()]
192         model = SimpleNN()
193         init_values = [p.data.clone().numpy
194             () for p in model.parameters()]
195         # originally no .numpy()
196         return init_values
197
198     def uniform_init():
199         model = SimpleNN()
200         init_values = [p.data.clone().numpy
201             () for p in model.parameters()]
202         # originally no .numpy()
203         return nn.init.uniform_(init_values
204             )
205
206     def latin_hypercube_sampling():
207         n_params = sum(p.numel() for p in
208             SimpleNN().parameters())
209         lhs_samples = pyDOE2.lhs(n_params,
210             samples=num_workers, criterion=
211             'maximin')
212         lhs_samples = lhs_samples * 2 - 1
213         # scale to [-1, 1]
214
215         init_values_list = []
216         for sample in lhs_samples:
217             init_values = []
218             start_index = 0
219             for p in SimpleNN().parameters
220                 ():
221                 end_index = start_index + p
222                     .numel()
223                 param_sample = sample[
224                     start_index:end_index].
225                     reshape(p.shape)
226                 init_values.append(torch.
227                     from_numpy(param_sample
228                     ))
229                 start_index = end_index
230             init_values_list.append(
231                 init_values)
232         return init_values_list

```

```

179
180 def adaptive_sampling():
181     def loss_function(params):
182         init_values = [param.reshape(p.
183             shape) for param, p in zip(
184                 params, SimpleNN().
185                 parameters())]
186         return simulate_parallel_loss(
187             SimpleNN(), init_values,
188             train_loader, device)
189
190     bounds = [(-1, 1)] * sum(p.numel()
191         for p in SimpleNN().parameters
192         ())
193     result = minimize(loss_function, x0
194         =np.zeros(len(bounds)), bounds=
195         bounds, method='L-BFGS-B')
196     best_params = result.x
197
198     return [best_params.reshape(p.shape
199         ) for p in SimpleNN().
200         parameters()]
201
202 def bayesian_optimization():
203     def loss_function(params):
204         init_values = [param.reshape(p.
205             shape) for param, p in zip(
206                 params, SimpleNN().
207                 parameters())]
208         return simulate_parallel_loss(
209             SimpleNN(), init_values,
210             train_loader, device)
211
212     bounds = [(-1, 1)] * sum(p.numel()
213         for p in SimpleNN().parameters
214         ())
215     result = gp_minimize(loss_function,
216         bounds, n_calls=num_workers,
217         n_random_starts=0, random_state
218         =42)
219     best_params = result.x
220
221     return [best_params.reshape(p.shape
222         ) for p in SimpleNN().
223         parameters()]
224
225 methods = [
226     ('Single Random Initialization',
227         single_random_initialization),
228     ('Uniform Sampling',
229         uniform_sampling),
230     ('Adaptive Sampling',
231         adaptive_sampling),
232     ('Bayesian Optimization',
233         bayesian_optimization),
234     ('LHS', latin_hypercube_sampling),
235 ]
236
237 losses = {}
238 for method_name, method in methods:
239     print(f"Running {method_name}...")
240     if method_name=='Single Random
241         Initialization':
242         num_workers=1
243     else:
244         num_workers=NUM_WORKERS_
245     init_values = method()

```

```

218     loss_curve = run_simulation(
219         num_workers, init_values,
220         train_loader, device,
221         num_epochs,to_demo)
222     losses[method_name] = loss_curve
223
224     print(losses)
225     to_csv(losses,num_epochs)
226     # Plot loss curves
227     plt.figure(figsize=(12, 6))
228     for method_name, loss_curve in losses.
229         items():
230         plt.plot(loss_curve, label=
231             method_name)
232
233     plt.xlabel('Epoch')
234     plt.ylabel('Loss')
235     plt.title('Loss Curves for Different
236         Initialization Methods')
237     plt.legend()
238     plt.grid()
239     plt.show()
240
241     results_dir = './results'
242     if not os.path.exists(results_dir):
243         os.makedirs(results_dir)
244
245     now=datetime.now()
246     plt.savefig(os.path.join(results_dir,
247         now.strftime("%m-%d-%Y_%H-%M-%S")+
248         'loss_curves.png'))
249     plt.show()
250
251 if __name__ == "__main__":
252     main(to_demo=True, NUM_WORKERS_=40,
253         EPOCHS=50, DEMO_EPOCHS=300)

```