

FLEXMoE: Scaling Large-scale Sparse Pre-trained Model Training via Dynamic Device Placement

XIAONAN NIE, Peking University, China
 XUPENG MIAO, Carnegie Mellon University, USA
 ZILONG WANG, Microsoft, China
 ZICHAO YANG, Carnegie Mellon University, USA
 JILONG XUE, Microsoft Research, China
 LINGXIAO MA, Microsoft Research, China
 GANG CAO, Beijing Academy of Artificial Intelligence, China
 BIN CUI*, Peking University, China

With the increasing data volume, there is a trend of using large-scale pre-trained models to store the knowledge into an enormous number of model parameters. The training of these models is composed of lots of dense algebras, requiring a huge amount of hardware resources. Recently, sparsely-gated Mixture-of-Experts (MoEs) are becoming more popular and have demonstrated impressive pretraining scalability in various downstream tasks. However, such a sparse conditional computation may not be effective as expected in practical systems due to the routing imbalance and fluctuation problems. Generally, MoEs are becoming a new data analytics paradigm in the data life cycle and suffering from unique challenges at scales, complexities, and granularities never before possible.

In this paper, we propose a novel DNN training framework, FLEXMoE, which systematically and transparently address the inefficiency caused by dynamic dataflow. We first present an empirical analysis on the problems and opportunities of training MoE models, which motivates us to overcome the routing imbalance and fluctuation problems by a dynamic expert management and device placement mechanism. Then we introduce a novel scheduling module over the existing DNN runtime to monitor the data flow, make the scheduling plans, and dynamically adjust the model-to-hardware mapping guided by the real-time data traffic. A simple but efficient heuristic algorithm is exploited to dynamically optimize the device placement during training. We have conducted experiments on both NLP models (e.g., BERT and GPT) and vision models (e.g., Swin). And results show FLEXMoE can achieve superior performance compared with existing systems on real-world workloads — FLEXMoE outperforms DeepSpeed by 1.70× on average and up to 2.10×, and outperforms FasterMoE by 1.30× on average and up to 1.45×.

CCS Concepts: • **Computing methodologies** → **Self-organization**; *Massively parallel algorithms*.

Additional Key Words and Phrases: Deep Learning System, Distributed Computing, Sparse Model

*Bin Cui is the corresponding author.

Authors' addresses: Xiaonan Nie, xiaonan.nie@pku.edu.cn, School of CS & Key Lab of High Confidence Software Technologies (MOE), Peking University, Beijing, China; Xupeng Miao, xupeng@cmu.edu, Carnegie Mellon University, USA; Zilong Wang, zilongwang@microsoft.com, Microsoft, China; Zichao Yang, yangtze2301@gmail.com, Carnegie Mellon University, USA; Jilong Xue, jxue@microsoft.com, Microsoft Research, China; Lingxiao Ma, lingxiao.ma@microsoft.com, Microsoft Research, China; Gang Cao, caogang@baai.ac.cn, Beijing Academy of Artificial Intelligence, China; Bin Cui, bin.cui@pku.edu.cn, School of CS & Key Lab of High Confidence Software Technologies (MOE), Institute of Computational Social Science, Peking University (Qingdao), Peking University, Beijing, China.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.

2836-6573/2023/5-ART110 \$15.00

<https://doi.org/10.1145/3588964>

ACM Reference Format:

Xiaonan Nie, Xupeng Miao, Zilong Wang, Zichao Yang, Jilong Xue, Lingxiao Ma, Gang Cao, and Bin Cui. 2023. FLEXMoE: Scaling Large-scale Sparse Pre-trained Model Training via Dynamic Device Placement. *Proc. ACM Manag. Data* 1, 1, Article 110 (May 2023), 19 pages. <https://doi.org/10.1145/3588964>

1 INTRODUCTION

Large-scale pre-trained models (PTMs), such as Transformer models, have promoted the deep learning (DL) development on various complicated tasks, including natural language processing, e.g., BERT [9], GPT [6], T5 [41], computer vision, e.g., ViT [10], Swin [25], advertising recommendation, e.g., M6 [24], and so on. These models are also known as foundation models since they are trained on hundreds of gigabytes of data and can be adapted, e.g., task-specific fine-tuning, to a wide range of downstream tasks. However, recent studies [6, 20] have demonstrated the model quality scales as the power law with data size, parameter size, and computation budgets. Current state-of-the-art foundation models have excessively grown to trillions of parameters and therefore become extremely expensive and time-consuming to be trained.

To avoid hitting the train-ability wall, one well-known line of studies utilizes the sparsely-gated Mixture-of-Experts (MoE) [45] structure to make pre-trained models cost-effective. This sparse MoE architecture enlarges the model parameter size by expanding expert networks while keeping the computational budget stable. Specifically, each MoE layer consists of a gate network and a series of experts (e.g., tens of), where the gate network routes the input tokens (i.e., the inputs of the layer) to only a small number of experts rather than all experts to be computation-efficient. MoE models have been successfully used to boost the performance of various PTMs, such as language models (e.g., GLaM [11] achieve the state-of-the-art performance in the SuperGLUE NLP benchmark [47]) and vision models (e.g., V-MoE [43] matches the most powerful vision transformers with only half computation time).

Due to the superior behaviors, large-scale MoE models have attracted wide interests in industrial companies recently when building real-world big data applications, including Google [11], Meta [44], Microsoft [42] and Alibaba [24]. Such kind of data-intensive workloads is becoming a new data analytics paradigm for data science research communities. With the rapid expansion of MoE models, unprecedented opportunities are brought into the data life cycle (especially in model development, deployment, and management procedures), as well as challenges. For example, Google created the ST-MoE [50] model to help improve the question answering quality of their search engines, which achieves state-of-the-art result in the AIRC leadboard [5], e.g., increasing from 81.40% to 86.52% by using MoE. However, such a large MoE model not only needs a lot of data to train on, e.g., 1.5T tokens for ST-MoE, but also consumes large amounts of expensive computing resources for training. Therefore, it is necessary to address data management problems hidden in existing systems and conduct optimizations to overcome these challenges.

Since a single MoE layer can easily exceed the limited GPU memory, previous distributed training techniques (e.g., data parallel [23] and model parallel [30, 32, 33, 36]) are becoming unsatisfactory. To train large MoE models, Lepikhin et al. [21] proposed the *expert parallelism* technique, which distributes experts into multiple GPUs to fit within the limited GPU memory. For each training token, after the gate network has determined the routing, it will be sent to the GPUs where its target experts reside. Since it adapts to the special characteristic of MoE structures, expert parallelism has become one of the most widely-used techniques to train MoE-based models in several well-known frameworks, including DeepSpeed-MoE [42], Tutel [19] and HetuMoE [37].

The Challenge of Workload Imbalance. Despite the widespread adoption of expert parallelism [21], it suffers from the severe workload imbalance brought by the sparse and conditional computing nature of MoE layers. Specifically, the input tokens are organized at a fine-grained level

and routed to the individual experts by the gating network. We found that the input tokens show highly uneven preferences on the experts and the routing results keep varying over training [34]. As reported by [17, 22, 26] and verified by our empirical studies, such workload imbalance could lead to significantly efficiency slowdown and severe resource under-utilization.

The Current Landscape. To tackle the workload imbalance problem, most existing works leverage an auxiliary balance loss on each MoE layer to enforce balanced routing. To control the trade-off between system efficiency (workloads balance) and statistical efficiency (model quality), the balance loss is further associated with a tunable penalty weight. In addition, Fedus et al. [12] introduces the expert capacity to avoid excessive input token assignments to the experts. In other words, tokens beyond the expert capacity will be dropped by the experts. Nevertheless, these methods can only mitigate the workload imbalance problem rather than fundamentally guarantee balance. Moreover, they make users stuck in a dilemma: *whether and to what degree shall we sacrifice model quality for system efficiency?* For instance, applying a large penalty weight to the balance loss (or applying a low capacity to experts) could help load balancing among experts but harm the model quality since a substantial amount of tokens would be routed to some less-desired experts (or dropped), and vice versa.

Intuitively speaking, the aforementioned methods try to adjust the data-flow of tokens to be as even as possible, i.e., enforcing all experts to process almost the same amount of tokens. They are friendly adapted to existing expert parallelism systems, because users can easily change the definitions of their model (e.g., balance loss and capacity) and without any modifications on the implementation of underlying DL systems. In other words, such approaches are friendly to people who are not familiar with DL systems so they have to sacrifice the model quality for training efficiency. Unlike previous approaches making model modifications, which may cause model quality degradation, we try to optimize the workload imbalance from a system perspective:

How to design a distributed MoE training system that could achieve the maximum system efficiency without affecting the model quality in the meanwhile?

Summary of Our Approach. In this work, we propose FLEXMoE, a novel distributed training system for large-scale sparse MoE models. We overcome the **routing imbalance** problem from the brand new perspective of the expert placement — instead of placing each expert on a single GPU device as in traditional expert parallelism nor duplicating all experts to all GPU devices as in traditional data parallelism, we introduce a **fine-grained replicated expert parallelism** that selects specific heavy experts, duplicates them over multiple devices, and spreads the input tokens across the replicas. It is non-trivial in sparse MoE models because the device placement of these experts changes the original input data assignment (i.e., all-to-all communication across different experts) and involves additional synchronization overheads (i.e., all-reduce communication among replicated experts). By precisely modeling the actual execution procedure of the MoE layer, FLEXMoE estimates the costs that determine the concrete expert replication scheme.

We further investigate the uneven distribution of the sparsely-activated experts in various MoE training workloads. Our key finding is that MoE models exhibit **routing fluctuation** during the iterative training process. We observe that the imbalanced expert preference distribution changes continuously and smoothly during the training process (will be described in-depth in Section 2.4). The training of gate networks is highly unstable [7] because gradient-based optimization algorithms will reinforce the previous routing determination until it reaches a certain point and escapes from the wrong learning direction. To handle this obstacle, we design a **dynamic expert management** mechanism in FLEXMoE to adaptively adjust the expert-to-device mapping and the assignment of tokens during training. Specifically, we adopt a *data-driven approach* to adaptively change the expert placement during training by monitoring the traffic trends. For instance, FLEXMoE gradually

expands resources for experts with increasing workloads for faster calculations and shrinks those with decreasing workloads to reduce replica synchronization overheads.

FLEXMoE is designed as a novel scheduling module on top of existing DNN frameworks, which monitors data traffic, makes scheduling plans, and dynamically adjusts the mapping between the data-flow graph and distributed devices. Particularly, three placement adjustment primitives (i.e., expand, shrink, migrate) are provided to flexibly govern the expert placement and a gate flow-control mechanism is introduced to enable autonomous global traffic optimization. Based on these building blocks, a simple yet effective algorithm is designed to estimate the benefits or overheads of each adjustment to determine the optimal scheduling plan.

Last but not least, we implement FLEXMoE on top of PyTorch [38] and verify the superiority of our work through comprehensive experiments. When training several popular MoE-based PTMs (e.g., BERT, GPT, Swin) of various tasks on a 64-GPU A100 cluster, FLEXMoE achieves up to $2.1\times$ of speedup compared with existing state-of-the-art MoE training systems. Furthermore, FLEXMoE consistently outperforms them in terms of model quality since we do not accommodate the expert placement with compromised token-routing.

Paper Organization. In section 2, we discuss the inefficiency of training sparse MoE models with expert parallelism and explore opportunities for optimizing such problems through system-level optimizations. In section 3, we introduce the system design of FLEXMoE, which leverages the vExpert abstraction to facilitate dynamic expert management and device placement. We provide details on the system implementation in Section 4 and evaluate the effectiveness of FLEXMoE in Section 5.

2 BACKGROUND AND MOTIVATIONS

We first present the notations used in our paper.

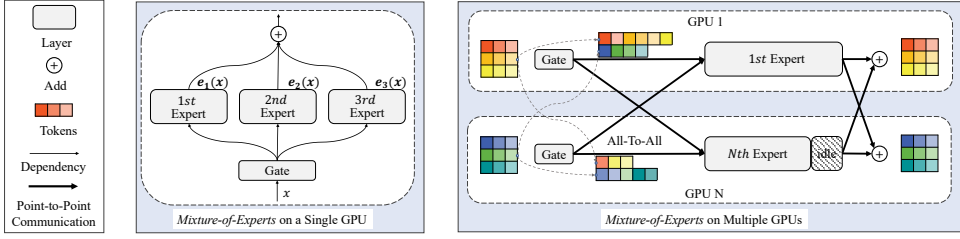
- * \mathcal{E} : A series of experts $\{e_1, \dots, e_N\}$.
- * \mathcal{G} : A set of GPUs, where $g \in \mathcal{G}$.
- * \mathcal{I} : The assignment of tokens, where $I_{e,g} \subset \mathcal{I}$.
 $I_{e,g}$ represents the number of tokens for expert e to GPU g .
- * \mathcal{P} : The expert-to-device mapping, where $(e, g) \in \mathcal{P}$.
 $(e, g) \in \mathcal{P}$ represents GPU g exists expert e .
- * TPS : Tokens-per-second for an expert.
- * $Bw_{g,g'}$: Bandwidth between GPU g and g' .
- * $BPS(\mathcal{G}')$: Bytes-per-second for AllReduce on a set of GPUs \mathcal{G}' .

2.1 Transformer with Mixture-of-Expert

The Transformer architecture [46] has demonstrated its superior performance in various tasks [6, 9, 10, 25, 27, 41, 49], which mainly consists of attention networks and feed-forward networks. Each attention network first linearly transforms the input tokens into corresponding queries (Q), keys (K) and values (V) and then performs the scaled dot-product on them as Equation 1, where d is the dimension of queries and keys. Meanwhile, this attention network benefits from capturing the dependencies between tokens within the sequence, and thus is effective in sequence transduction tasks.

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V \quad (1)$$

Each feed-forward network (FFN) is composed by two fully connected layers and an activation function, i.e., ReLU, formulated as Equation 2. Different from the attention network, FFNs model the relation of different feature dimensions within a single token by scaling it into larger dimension



(a) The workflow of an MoE layer (b) The distributed training of an MoE layer with Top-1 gate

Fig. 1. Illustrations of Mixture-of-Experts (MoE). Figure 1(a) represents the workflow of a typical MoE layer (detailed in Section 2.1). Figure 1(b) presents the distributed training of an MoE layer as expert parallelism, where experts are partitioned while non-MoE layers (e.g., self-attention, gate) are duplicated across devices. For example, according to the result of Gate, GPU 1 routes 6 input tokens to 1st Expert and the other 3 tokens to Nth Expert (detailed in Section 2.2).

space. For example, the output dimension of W_1 is always set as $4\times$ of its input dimension. Existing famous pre-trained models are usually stacked by a series of transformer layers to improve its model capacity as well as model quality.

$$\text{FFN}(x) = W_2 \cdot \text{ReLU}(W_1 \cdot x + b_1) + b_2 \quad (2)$$

Recently, scaling with more data and more parameters has driven significant model quality improvement [6, 10] while requires large amounts of computing resources for training. To scale models efficiently, researchers have adopted the sparse-gated Mixture-of-Experts (MoE) paradigm [12, 45] by replacing the FFN network with an MoE layer, where each input token activates only a subset of model parameters, and thus introducing model sparsity.

The key components of a MoE layer include a data-dependent sparse gate network $g(x)$ and a series of experts \mathcal{E} , as shown in Figure 1(a). For each input token x , the gate network first produces the probability of x with respect to all experts and then routes x to its corresponding experts. The Top-K gate [45] is formulated in Equation 3, which keeps only the top k values before the softmax function.

$$g(x) = \text{softmax}(\text{TopK}(x \cdot W_g)) \quad (3)$$

As soon as each expert e_i receives its input token x ($e_i \in \mathcal{E}$), it produces its corresponding output $e_i(x)$ and then the final output y is obtained by the linear combination of $e_i(x)$ weighted by the output of the gate network $g(x)_i$ as follows:

$$y = \sum_{i=1}^N g(x)_i \cdot e_i(x). \quad (4)$$

When the gate network's sparsity is static, e.g., Top-1 [12] or Top-2 [21], the computation and communication costs of given inputs nearly remain constant as the number of experts increases, allowing models to be scaled effectively with MoE.

To help readers better understand the routine of MoE layers, we also provide an illustration in Figure 1(b), where each GPU is associated with one expert and 9 tokens are fed to the gate network and routed to their target experts. For example, GPU 1 routes 6 and 3 tokens to the 1st and Nth experts, while GPU N routes 4 and 5 tokens to these two experts. This assignment would cause the workload imbalance problem under current expert placement (i.e., 10 tokens for GPU 1 and 8 tokens for GPU N). Since there is a reverse routing step after the expert computation, all GPUs have to wait for each other before executing the following layers. Such an imbalanced workload inevitably results in GPU under-utilization and low training efficiency.

2.2 Distributed Training

As models are sparsely scaled with MoE, multiple GPUs would be involved for model management and training acceleration [33, 35, 42]. We will analyze the existing popular parallelism strategies in the following:

Data Parallelism. Data parallelism (DP) is usually used to scale training when the model can fit in the device's available GPU memory, as the communication primitives (e.g., AllReduce) of DP achieve good scalability performance. In DP, training samples are partitioned while model parameters are duplicated for each device. Given a batch of training samples, each worker executes the forward and backward computation, synchronizes gradients globally (i.e., averaged), and updates local parameters based on the synchronized gradients. However, each device has to maintain a full replica of the model and thus DP can't be used to scale up to large models.

Model Parallelism. If the memory requirement of the given model exceeds the GPU memory, it should be partitioned across multiple devices. In tensor parallelism (TP), every single tensor can be partitioned over multiple devices to reduce the memory footprint of each GPU. For example, Megatron-LM [33] proposed to partition the attention network by exploiting its inherent parallelism in the multi-head attention operation where the queries (Q), keys (K) and values (V) matrices can be partitioned in a column-parallel fashion. Pipeline parallelism (PP) is an alternative partitioning method, where different groups of layers are placed on different devices. For example, GPipe [18] partitions the input batch into a number of micro batches and pipelines each device's computation across micro batches to improve the resource utilization of devices.

Expert Parallelism. GShard [21] first proposed expert parallelism (EP) as a specific model parallelism method for MoE models, where experts within an MoE layer are placed on different devices and non-MoE layers are replicated on devices as DP. The workflow of distributed training a Top-1 gate MoE layer is illustrated in Figure 1(b). After getting the target expert for each token, an All-to-All communication is involved to send tokens to target experts for processing. And another All-to-All communication is needed to send data back for the execution of data-parallel non-MoE layers. As MoE models often have numerous experts, e.g., 1024 Experts in Switch Transformer [12], EP can scale up with model size better than model parallelism.

2.3 Problem Formulation

To model the training cost, we consider a single MoE layer, which consists of a gate network and a set of experts \mathcal{E} ($e \in \mathcal{E}$). Given current expert-to-device mapping \mathcal{P} ($(e, g) \in \mathcal{P}$ represents GPU g has expert e), the assignment of tokens \mathcal{I} ($I_{e,g}$ represents the number of tokens for expert e to GPU g) at the current step, the training cost can be formulated as $T(\mathcal{I}, \mathcal{P})$. Our objective is to minimize this training cost, which can be expressed as follows:

$$\min T(\mathcal{I}, \mathcal{P}) = \min_{g \in \mathcal{G}} \sum_{(e,g) \in \mathcal{P}} \{T_C(I_{e,g}) + T_{A2A}(I_{e,g}) + T_{Sync}(\mathcal{P}, e)\} \quad (5)$$

The first two terms represent the cost of expert computation and All-to-All communication respectively, which are determined by the assignment of tokens, i.e., $I_{e,g}$. And the third term represents the cost of synchronization (T_{Sync}) because one expert may exist on multiple GPUs as data parallelism and needs to synchronize their gradients to maintain consistent. Each GPU sums up all of its local experts with regard to these three terms to obtain its execution time and the global training time $T(\mathcal{I}, \mathcal{P})$ is the maximum execution time among these GPUs. In order to minimize the training cost while not modify the definitions of models, we design a dynamic expert management mechanism to adaptively adjust the expert-to-device mapping \mathcal{P} and the assignment of tokens \mathcal{I} during training.

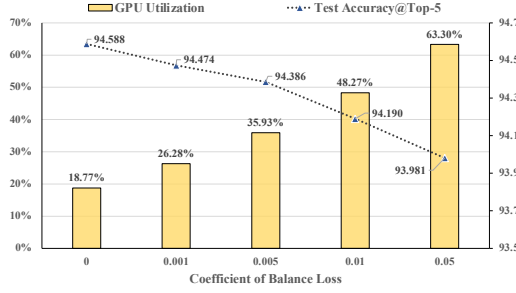


Fig. 2. The top-5 accuracy of Swin-MoE under different balance loss coefficients, where we do not restrict the capacity of each expert to ensure no token is absent.

Existing System-Friendly Optimizations. As discussed in Section 1, existing works focus on addressing the workload imbalance problem by modifying the definitions of models (e.g., balance loss and capacity), which doesn't need users to make modifications on the underlying DL systems and therefore is system-friendly.

- The *balance loss* depicts the imbalance level of \mathcal{I} and is widely used in MoE training [12, 21]. Once the gate network produces an imbalanced token assignment, it would be penalized by a large balance loss. In practice, the training process aims to minimize the weighted sum of balance loss and training loss. Thus, there exists a trade-off — emphasizing the balance loss would drive the assignment \mathcal{I} to be even but harms the model quality since the training loss would be higher, and vice versa.
- The *capacity* is a threshold that limits the number of input tokens for each expert (i.e., $I_{e,g}$). Tokens that exceeds the capacity will be skipped by the expert and directly forwarded to the next layer by the residual connection. In short, the capacity upper bounds the training cost of the heaviest expert to improve the overall efficiency, however, at the price of degraded model quality since a certain amount of tokens cannot be fully trained.

2.4 Observations and Opportunities

Observation 1: Limitations of the existing techniques. We first studied the effect of the existing techniques empirically and took Swin-MoE [25] as an example, where we varied the balance loss coefficient and did not restrict the capacity of each expert to ensure no token was absent. The results are summarized in Figure 2. By increasing the balance loss coefficient from 0 to 0.05, the GPU utilization is improved from 18.77% to 63.30% while the top-5 accuracy is decreased from 94.588% to 93.981%. It demonstrates that enforcing workload balance by adjusting the assignment of tokens \mathcal{I} would inevitably lead to the trade-off between system efficiency and model quality.

Observation 2: Dynamic imbalanced workloads. We also recorded the trace of training GPT-MoE models (64 expert per MoE layer) and studied the loads of experts during training. Results are summarized in Figure 3 and we have identified two key characteristics:

- **Skewness.** As shown in Figure 3(a), we visualize the computational loads of experts as the cumulative distribution function (CDF), where we observe that the Top-10 experts (10 out of 64) receive almost 75% tokens, leading to routing imbalance in MoE training. If the experts are evenly distributed among GPUs as in expert parallelism [21], such imbalanced workloads would result in severe resource under-utilization, as most experts need to wait for the slowest.
- **Smoothness and continuousness.** In Figure 3(b), we present evolution of expert loads throughout the entire training process, where different intervals represent different experts. We observe that the load of each expert is continuously changing during training, for example, from less to

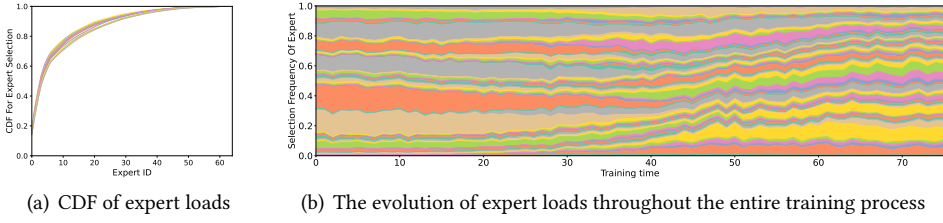


Fig. 3. Illustrations of expert loads. For a single step, we sort experts within each MoE layer based on their computational loads and visualize the corresponding cumulative distribution function (CDF) in Figure 3(a), where different colors indicate different MoE layers. For the entire training process, Figure 3(b) illustrates the dynamic changes in the load of each expert, where the different color areas represent the different experts.

more, from more to less, and from more to less and then more again, etc, which poses routing fluctuation when training MoE models. Fortunately, the load of each expert does not change dramatically in a short period of time, which means a smooth and continuous change.

Challenges. Motivated by these observations, our work explores how to develop a system that achieves workload balance by adjusting the expert-to-device mapping (i.e., \mathcal{P}) for better system efficiency while maintaining the model quality. Moreover, the system should dynamically adapt to the varying routing distribution during the training process. However, since both token routing and expert loads are decided by the data-dependent gating mechanism, we can not determine the optimal mapping ahead of its execution. The main challenge lies in designing and modifying the expert-to-device mapping efficiently under the fast GPU computation and rapid change in workloads. Another challenge is how to efficiently implement these irregular MoE operations.

Opportunities. Fortunately, as previously mentioned, the distribution changes smoothly and continuously. This means that the optimal expert-to-device mapping would not shift significantly in a short period of time. Therefore, it is feasible to refine the mapping based on the ad-hoc routing determination, without the need to predict the optimum for the next few training steps. Furthermore, the cost of computation and communication can be estimated before the actual execution. Based on the cost models, we can predict the benefits and overheads of different mapping candidates to find the best one.

3 FLEXMOE DESIGN

3.1 Overview

The workflow of our FLEXMoE is illustrated in Figure 4, where (1) input tokens are first processed by the gate network to determine their corresponding target experts and then (2) the router leverages a greedy algorithm to distribute tokens to different replicas of each expert. Finally, (3) each replica performs computations on its assigned tokens and delivers the output results back.

To handle the dynamic workload imbalance, we have designed a dynamic expert management mechanism that adaptively adjusts the expert-to-device mapping \mathcal{P} and the assignment of tokens \mathcal{I} during training. In addition to the regular workflow, we also introduce two new components: Scheduler and Policy Maker. (4) Scheduler monitors the real-time loads of experts and sends them to Policy Maker if the current imbalance metric (i.e., balance ratio in Equation 6) exceeds the predetermined threshold. Based on the received loads and the current placement of experts, Policy Maker produces modification instructions and sends them back to Scheduler. Scheduler then interacts with the executor of current DL frameworks and triggers modifications of expert placement at runtime.

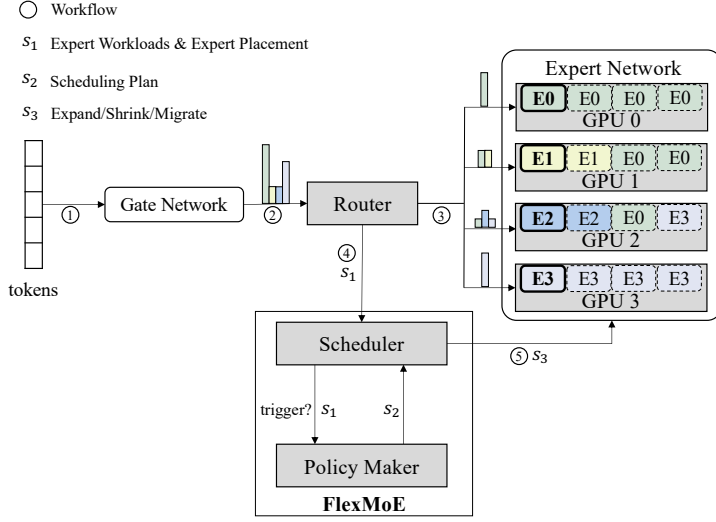


Fig. 4. FLEXMoE System Architecture.

To cope with the large decision space of flexible expert-to-device mapping and its dynamic modifications, FLEXMoE proposes the abstraction of vEXPERT to represent the minimum unit for scheduling, as detailed in Section 3.2. Moreover, FLEXMoE decouples the expert placement modifications using three placement modification primitives, including Expand, Shrink and Migrate, as explained in Section 3.3. Finally, we design a greedy heuristic algorithm to generate a sequential combination of these modification primitives to produce the adjustment plan of expert-to-device mapping, as described in Section 3.4.

3.2 Dynamic Expert Management and vEXPERT

To tackle the optimization problem related to expert-to-device mapping and its corresponding token routing, we introduce a novel abstraction called vEXPERT, which defines the minimum unit for scheduling GPU computations for an expert and enables dynamic expert management. The vEXPERT abstraction helps in determining how to duplicate experts, when to increase or decrease replicas, and how to partition tokens between replicas.

vEXPERT has the following characteristics:

- Each vExpert can only be assigned as the replica of exactly one expert and process part of tokens for the master expert.
- Each vExpert shares weights with other vExperts for the same expert on the same GPU, which means packing the same vExperts of the same GPU.
- Evenly workloads (i.e., tokens) partitioning among all vExperts for the same expert.

With this abstraction, we can significantly reduce the large search space of the optimization problem by making decisions at the vEXPERT level. Besides, we have designed three placement modification primitives to manage the dynamic expert management at the vEXPERT level, including Expand, Shrink and Migrate (see Section 3.3 for details). These primitives allow for arbitrary modifications of expert placement by composing them in different ways. Moreover, because the amount of data in each iteration is fixed ($\sum B_i = B$), we can calculate the ideal capacity of each vExpert ($B/(G * E)$) as a reference for decision-making, where B is the total number of input tokens, G is the number of GPUs and E is the number of vEXPERTS on each device.

Algorithm 1: Scheduler: Monitor Real-Time Workloads

Input: \mathcal{I} : the assignment of tokens;
 \mathcal{P} : the expert-to-device mapping;

```

1 while is training do
2    $balance\_ratio \leftarrow balance(\mathcal{I}, \mathcal{P})$  // Equation 6 ;
3   while  $balance\_ratio > threshold$  do
4      $plan \leftarrow MakeSchedulingPlan(\mathcal{I}, \mathcal{P})$  // Algorithm 2 ;
5     if  $plan$  is empty then
6       break ;
7      $\mathcal{P} \leftarrow Modify\ Expert\ Placement\ \mathcal{P}\ with\ plan$  ;
8      $balance\_ratio \leftarrow balance(\mathcal{I}, \mathcal{P})$  ;
9    $\mathcal{P} \leftarrow Modify\ Expert\ Placement\ \mathcal{P}\ with\ Migrate$  ;
  
```

3.3 Scheduler

Figure 4 illustrates the role of the Scheduler in connecting the Runtime Executor and the Policy Maker in FLEXMoE. Our design includes a metric, the balance ratio, which measures real-time workloads and determines when to trigger the Policy Maker for adjustment decisions. Additionally, we propose three atomic primitives for the Runtime Executor to use in scheduling experts. These primitives will be further explained in the following section.

Balance ratio. Owing to the synchronous execution mode in MoE layer, the slowest GPU will dominate the finish time of all-to-all communication as well as the global training. Thus, we design the balance ratio as Equation 6, which adds up the loads of all vEXPERTS on a single GPU and finds the max value among all GPUs to get the final results. While there may be other metrics that can be used to measure loads, we compare the balance ratio in Equation 6 with variance as a metric and provide a detailed analysis of our findings in Section 5.3.

$$balance(\mathcal{I}, \mathcal{P}) = \frac{1}{\sum_{(e,g) \in \mathcal{P}} I_{e,g} / |\mathcal{G}|} \max_{g \in \mathcal{G}} \sum_{(e,g) \in \mathcal{P}} I_{e,g} \quad (6)$$

Expand. When an expert receives increasing workloads, previously allocated vEXPERT resources may handle them slower than others and thus produce stragglers (over-utilized). In this scenario, FLEXMoE will call the Expand primitive to allocate an extra new vEXPERT resource for this expert at a time. The Expand primitive copies the expert parameter as well as the corresponding optimizer states from the source vEXPERT to the target vEXPERT via parameter sharing for intra-GPU communication and point-to-point communication of inter-GPU communication.

Shrink. In the opposite, when previously allocated vEXPERT resources are not fully utilized due to decreasing workloads (under-utilized), FLEXMoE will call the Shrink primitive to release an vEXPERT resource at a time. The Shrink primitive is executed by a tag without any communication, and thus introduces no overheads.

Migrate. If an expert retains multiple replicas on different GPUs, the training is slowed down due to the need for extra communications (i.e., gradients all-reduce) for synchronizations. The efficiency of this communication depends on the number of GPUs involved and their locality. To reduce the number of GPUs holding expert replicas and lower the synchronization cost, FLEXMoE will call the Migrate primitive exchanges the model states between two vEXPERTS.

We present the workflow of Scheduler in Algorithm 1, which monitors the real-time workloads, i.e., the assignment of tokens \mathcal{I} , and measures the balance ratio under current placements \mathcal{P} (Line 1-2). If current balance ratio exceeds the pre-defined *threshold*, Scheduler will repeatedly ask the Policy Maker for modification primitives until no beneficial modifications exist (Line 3-8). Then,

Algorithm 2: Policy Maker: vEXPERT-Based Scheduling

Input: \mathcal{I} : the assignment of tokens;
 \mathcal{P} : the expert-to-device mapping;

```

1 Function MakeSchedulingPlan( $\mathcal{I}, \mathcal{P}$ ):
2   Estimate time cost  $t_0$  with  $\mathcal{I}, \mathcal{P}$  by Equation 5;
3   for  $e \in \mathcal{E}$  do
4      $n_e \leftarrow$  number of vEXPERT allocated for  $e$  ;
5      $cap_e \leftarrow \mathcal{I}_e / n_e$  // capacity of vEXPERT for  $e$  ;
6    $e_0, e_1 \leftarrow \arg \max_{e \in \mathcal{E}} \{cap_e\}, \arg \min_{e \in \mathcal{E}} \{cap_e\}$  ;
7   Estimate time cost  $t_1$  after expanding  $e_0$  and shrinking  $e_1$ ;
8   if  $t_1 < t_0$  then
9     return {(Expand,  $e_0$ ), (Shrink,  $e_1$ )};
10  return {};

```

Scheduler turns to the Migrate operation to reduce the synchronization cost and continuously optimizes it at backend (Line 9). Intuitively, the vEXPERT-based Expand and Shrink primitives tackle the complicated expert-device mapping problem under the dynamic changing workloads, and the Migrate primitive continuously optimizes the placement of replicas for each expert.

3.4 Policy Maker

Using the modification primitives described above, FLEXMoE employs an efficient vEXPERT-based scheduling algorithm for generating sequential modification operations, as shown in Algorithm 2. Specifically, we leverage a cost-model driven search planning approach, which makes decisions based on feed-backs from simulating the training time.

To model the training cost of an MoE layer, we decompose it into three parts as Equation 5, including computation cost T_C , All-To-All communication cost T_{A2A} and expert synchronization cost T_{Sync} . For each part, we build cost models that take into account input variables (e.g., \mathcal{I} : the assignment of tokens and \mathcal{P} : the expert-to-device mapping) as well as environmental variables (e.g., TPS: tokens-per-second for an expert). By leveraging a profiling-based approach, we first profile the function's running time under different input sizes and then estimate the corresponding environmental variables. Besides, we also consider the cost of expert adjustments, which could be executed concurrently with model training.

Computation Cost. The computation cost refers to the time taken for an expert to perform the forward and backward computation of experts during training, which can be formulated as:

$$T_C(I_{e,g}) = \frac{I_{e,g}}{TPS} \quad (7)$$

where $I_{e,g}$ represents the number of input tokens received by expert e on GPU g , and TPS represents the throughput (tokens per second) of given GPUs to calculate an expert, which is obtained from profiling. The computation cost is proportional to the number of input tokens received by an expert, and inversely proportional to the throughput of the GPU. In other words, the more input tokens an expert receives, the longer it will take to compute the forward and backward passes, and the slower the GPU, the longer it will take to process each token.

All-To-All Cost. The All-To-All operation is involved to send tokens to target experts and send back their results after processing, which will be called for totally 4 times in each training step. We predict the All-To-All cost by a topology-aware model as Equation 8, where the expert e on GPU g has received $I_{e,g}.count(g')$ tokens from GPU g' and $Bw_{g,g'}$ is the profiled bandwidth between GPU g and GPU g' . Our cost model considers the intra-node bandwidth (e.g., PCIe, NvLink) and

Algorithm 3: Flexible Token Routing**Input:** \mathcal{I} : Expert Workloads; \mathcal{P} : Expert Placement; g : Current GPU;**Output:** r : Routing plans

```

1 /*iterate over all experts*/ ;
2 for  $e \in \mathcal{E}$  do
3    $n_e \leftarrow$  number of vEXPERT allocated for  $e$  ;
4    $cap_e \leftarrow \mathcal{I}_e / n_e$  // capacity of vEXPERT for  $e$  ;
5    $r_{e,g} \leftarrow \min(cap_e \times n_{e,g}, \mathcal{I}_{e,g})$  // locality first ;
6    $s_e \leftarrow \mathcal{I}_{e,g} - r_{e,g}$  // tokens for other GPUs ;
7    $c_{e,g} \leftarrow r_{e,g} - \mathcal{I}_{e,g}$  // local available capacity ;
8   for  $g' \in \mathcal{G} - \{g\}$  do
9      $c_{e,g'} \leftarrow \min(cap_e \times n_{e,g'}, \mathcal{I}_{e,g'}) - \mathcal{I}_{e,g'}$  ;
10     $r_{e,g'} \leftarrow s_e \times c_{e,g'} / \sum c_e$  // proportional to availability ;
11 return  $r \leftarrow \{r_e \mid e \in \mathcal{E}\}$  ;

```

inter-node bandwidth (e.g., IB, NIC) separately.

$$T_{A2A}(\mathcal{I}_{e,g}) = 4 * \sum_{g' \in \mathcal{G}} \frac{\mathcal{I}_{e,g}.count(g')}{Bw_{g,g'}} \quad (8)$$

Synchronization Cost. When expert e holds multiple replicas on different GPUs, their gradients must be synchronized via AllReduce communication, which is determined by the message size, the number of involved devices and their connected network. We enumerate different device groups and profile them before training to get their BPS (bytes-per-second). We predict the synchronization cost for the expert e by finding the device group that includes e (i.e., $\mathcal{P}.index(e)$) and then obtaining its corresponding BPS from profiling data. $size(e)$ represents the size of gradients for an expert.

$$T_{sync}(\mathcal{P}, e) = \frac{size(e.gradient)}{BPS(\mathcal{P}.index(e))} \quad (9)$$

Expert Adjustment Cost. FLEXMoE proposes three modification primitives, including Expand, Shrink and Migrate. The Expand and Migrate primitives involve transferring model states from the source GPU to the target GPU via NCCL Point-to-Point communication. On the other hand, the Shrink primitive is executed with no overheads by marking a tag. We predict the cost of transferring model states as $\frac{size(e.model_states)}{Bw_{g,g'}}$.

The scheduling policy of Policy Maker is summarized in Algorithm 2. Firstly, the algorithm gets the time cost t_0 of current placement as the baseline (Line 2), and then finds the expert id with the maximum workload and the expert id with the minimum workload (Line 3-6). After that, our Policy Maker estimates the time cost t_1 after applying the Expand and Shrink primitives (Line 7) and decides whether to return the modification by comparing t_0 and t_1 (Line 8-10).

4 IMPLEMENTATION

We have implemented the proposed mechanisms and algorithms on the top of PyTorch [38] by adding new customized operators and CUDA kernels. FLEXMoE is also a part of a novel distributed DL system Hetu [28, 29, 31]. To schedule dynamic dataflow more efficiently, FLEXMoE incorporates the following system-level optimizations:

Flexible Token Routing. Router should efficiently transfer input tokens from multiple devices (e.g., GPUs) to their target experts according to the complicated expert-to-device mapping \mathcal{P} . To accomplish this, FLEXMoE has designed a greedy policy as Algorithm 3, which prefers to route

Table 1. Models for Evaluation.

Model	Params.	#Layer	d_{Model}	d_{FFN}	#Expert
BERT-MoE-S	0.988B	12	768	3072	32
BERT-MoE-L	6.69B	24	1024	4096	64
GPT-MoE-S	0.988B	12	768	3072	32
GPT-MoE-L	39B	24	2048	8192	64
Swin-MoE-S	946M	24	-	-	32
Swin-MoE-L	1.83B	24	-	-	64

tokens to the local GPU and then scatters the remaining tokens to other GPUs in proportion to their available capacity. FLEXMoE also implements a efficient expert-wise layout transformation to arrange the inputs in a continuous space for each expert.

Paralleled Operation Modification. FLEXMoE employs a queue to sequentially insert modification primitives triggered by the scheduler, such as Expand, Shrink and Migrate. To reduce the time cost of adjustments and kernel launch, we merge several consecutive and parallelizable operations to run them concurrently. For example, if two operations share the same source and destination, they can be merged to increase the message size and improve the bandwidth utilization. Meanwhile, if they share neither source or destination, they can be executed concurrently to improve the utilization of clusters.

AllReduce Coordination. When the vEXPERTS of a single GPU are assigned to different experts, it should call the synchronizations separately for each expert, and may cause deadlock since the order of calls is inconsistent for different GPUs [4]. To avoid this deadlock problem, we assign a logical id to each expert and the logical id of each replica (vEXPERT) is same as its main expert. Then, each GPU invokes synchronizations in ascending order of experts' logical id.

Best-Effort Adjustment. Since the placement modifications may block the training process, it is not clear whether executing the current modification will be beneficial due to the dynamic workloads. To address this problem, FLEXMoE leverages a separate CUDA stream to conduct adjustments concurrently with the available network bandwidth and adopts the best-effort adjustment to avoid hindering the training process. the Scheduler interacts with the DL executor to determine whether to call the first operation in the candidate queue.

NCCL Group Management. FLEXMoE adopts NCCL [4] library to perform collective communication among GPUs and multiple NCCL groups are required to execute the complicated synchronization of experts. However, there is a maximum number of live NCCL groups that can remain, and it is inefficient to eliminate the groups once they have been utilized. FLEXMoE employs a *Least Recently Used (LRU)* cache to maintain nccl groups and therefore reduces the costs of group creations and eliminations.

5 EVALUATION

In this section, we present the detailed evaluation results to demonstrate the effectiveness and scalability of FLEXMoE.

5.1 Experimental Setup

Machine environment. We conduct experiments on Azure VMs [1], each equipped with 192-core AMD CPUs and 8 NVIDIA Ampere A100 GPUs. The GPUs are connected via NVLink 3.0 intra-node and the servers are connected via 8 InfiniBand NICs (8*200 Gbps totally). RDMA is used by default and the PyTorch version is 1.11.

Baselines. We compare FLEXMoE with other popular frameworks, including DeepSpeed [42] and FasterMoE [17]. Deepspeed used expert parallelism, which was first proposed by GShard [21]. FasterMoE proposed the shadowing strategy to replicate the popular expert among all GPUs.

Table 2. Comparison on model quality

	Metric	Masked LM		Language Modeling		Metric	Image Classification	
		BERT-MoE-S	BERT-MoE-L	GPT-MoE-S	GPT-MoE-L		Swin-MoE-S	Swin-MoE-L
DeepSpeed	PPL ↓	3.53	3.31	12.2	10.71	acc@1 ↑	77.316	77.022
						acc@5 ↑	93.838	93.642
FLEXMoE	PPL ↓	3.14	3.07	11.72	10.47	acc@1 ↑	77.754	77.109
						acc@5 ↑	94.042	93.663

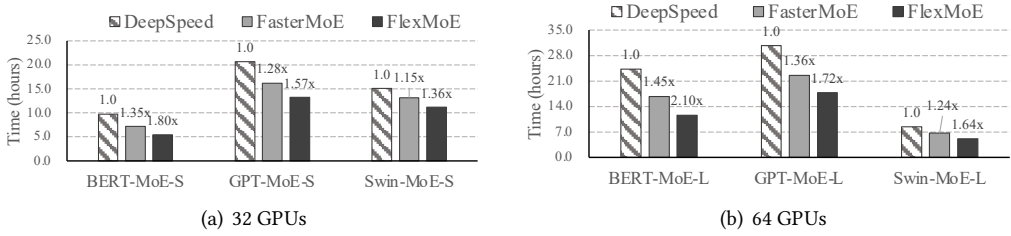


Fig. 5. Comparison on system efficiency

Benchmarks and datasets. Our evaluations are conducted by scaling representative transformer models in different application domains with the MoE architecture, including BERT [9] and GPT [40] in NLP and Swin [25] in CV, as shown in Table 1. We adopt Top-2 Gate for each model, which is adopted by widely used sparse MoE models (GShard [21], GLaM [11], V-MoE [43]), and set the capacity factor as 1.0 for each expert and balance loss as 0.001. We pretrain BERT-MoE with masked language modeling (MLM) and next-sentence-prediction (NSP) [9] tasks, GPT-MoE with language modeling task (LM) [40], Swin-MoE with image classification task [19, 25]. NLP models are trained on wikipedia [3] and vision models are trained on the ImageNet-1K [2]. All the hyper-parameters (e.g., learning rate) are fixed as for the same model.

5.2 Overall Performance

To evaluate the end-to-end efficiency of FLEXMoE, we compare it with DeepSpeed and FasterMoE. For DeepSpeed, we use the traditional expert parallelism approach and set the number of experts per GPU as 1 in every MoE layer. We evaluate two different width of models, X-MoE-S with 32 experts and X-MoE-L with 64 experts, on 32 GPUs and 64 GPUs, respectively.

Model Quality. To demonstrate the importance of not dropping tokens, we compare the model quality of FLEXMoE with DeepSpeed on various models and various tasks. We use the validation perplexity for language models (e.g., BERT-MoE and GPT-MoE), which is the lower the better, and top-1/top-5 accuracy of Imagenet-1K for vision models (e.g., Swin-MoE). As shown in Table 2, FLEXMoE outperforms DeepSpeed in almost all tasks, indicating that the limited capacity in existing MoE systems (e.g., DeepSpeed) can lead to model quality degradation. Considering the training cost, the Swin-MoE models (scaled based on the Swin-B model) are trained for 100 epochs (less than 300 epochs in the standard configuration) and thus the accuracy is slightly worse than the benchmark, where 94.04% of FLEXMoE v.s. 96.5% of the benchmark as for top-5 accuracy. And we believe it's enough to show the benefits of no dropping tokens. What's more, Swin-MoE-L performs slightly worse than Swin-MoE-S as the size of training data is small and thus it may lead to overfitting.

System Efficiency. We also evaluate FLEXMoE against other SOTA systems on efficiency. To measure efficiency, we record the required training time to achieve the target model quality, and the results are illustrated in Figure 5. Our experiments show that FLEXMoE achieves the best training

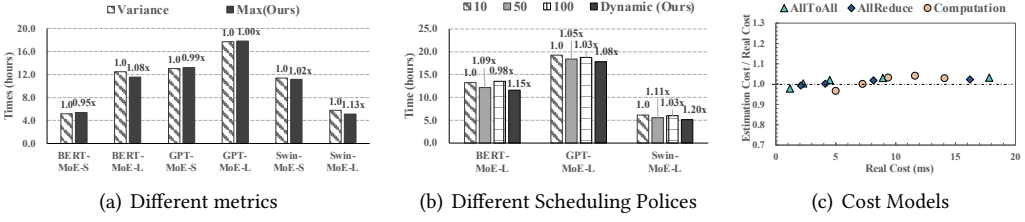


Fig. 6. Study of different metrics, different scheduling policies and cost models

efficiency, outperforming DeepSpeed by 1.70 \times on average and up to 2.10 \times , and FasterMoE by 1.30 \times on average and up to 1.45 \times . As mentioned above, although DeepSpeed obtains the smallest iteration time thanks to its limited capacity, it drops tokens to skip the expert network and thus requires more iterations to converge. FasterMoE proposes the dynamic shadowing strategy for loading balance, which replicates the popular expert on all GPUs. However, due to its coarse-grained expert management (i.e., on 1 GPU or on all GPUs), it falls back to a sub-optimal solution. As the number of GPUs increases, FasterMoE suffers from the global synchronization of expert replicas.

In addition to the expert networks, the time of models' training also consists of other parts, such as non-experts' computation, optimizers' update and communication. As FLEXMoE only optimizes the execution of the expert networks, we will mainly focus on analyzing the MoE layer in the following sections.

5.3 Ablation Study

To verify the effectiveness of FLEXMoE, we conduct several ablation studies, as demonstrated below.

Different Metrics. FLEXMoE utilizes the balance ratio as a key metric to trigger adjustments. To investigate the impact of alternative metrics, we also study the use of *Variance* as a ratio, formulated as $\sum_{g \in \mathcal{G}} (I_g - \bar{I})^2 / |\mathcal{G}|$, in addition to the *Max(ours)* ratio given by Equation 6. Results are presented in Figure 6(a) and show that *Max(ours)* outperforms *Variance* by 1.03 \times on average and up to 1.13 \times on Swin-MoE-L. Meanwhile, *Variance* also performs well on BERT-MoE-S and GPT-MoE-S. Because the training time is dominated by the slowest expert, *Max(ours)* is a simple but effective metric. *Variance* takes the global workload distribution into consideration, which triggers adjustment more frequently but often gets empty operations from Policy Maker as it is not always relevant to the actual training.

Different Scheduling Policies. FLEXMoE dynamically triggers adjustment according to the pre-defined balance ratio. We also conduct experiments on static scheduling policies, which triggers the adjustments in the fixed interval and executes them completely before training. As shown in Figure 6(b), our dynamic scheduling outperforms both large interval and small interval because of the dynamic changing workloads. Specifically, small interval will trigger adjustments frequently and thus introduce more adjustment costs, while large interval can not tackle the dynamic workloads well because it can not make adjustments in time. Our dynamic policy dynamically decides the adjustments based on current workloads, which is suitable for the dynamic workloads.

Evaluation on Cost Models Figure 6(c) demonstrates the effectiveness of our cost models, where we compare the estimation cost to real cost on different input sizes for computation/alltoall/allreduce respectively. It can be observed that our estimation results are very close to the real execution costs for all experimental models, where the average prediction error is less than 3%.

5.4 Token Efficiency and Expert Efficiency

In this section, we analyze both *token efficiency* and *expert efficiency* for different MoE training methods during the whole training process. *Token efficiency* refers to the fraction of input tokens

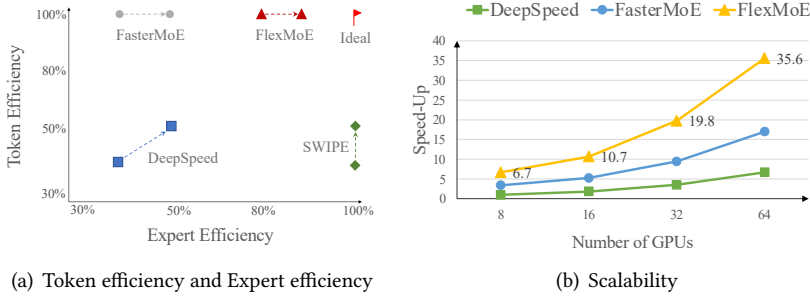


Fig. 7. Figure 7(a) shows the trends of token efficiency and expert efficiency for different methods during training. And Figure 7(b) shows the scalability for different methods.

that are processed by the expert network and *expert efficiency* refers to the meaningful computation of GPUs. 100% of both metrics is the ideal status of an MoE layer, shown as the red flag in Figure 7(a). The traditional expert-parallel method (e.g., DeepSpeed) obtains low token efficiency and expert efficiency as it drops tokens beyond capacity for loading balance. SWIPE, proposed by BaGuaLu [26], improves expert efficiency by modifying the gating algorithm to re-assign inputs to other experts for strict load balance. However, this approach changes the relations between tokens and experts, thus leads to low token efficiency. FasterMoE [17] replicates hot experts on each GPU and guarantees no tokens dropping in the implementation. However, it doesn't take load balance as its design goals and obtains low expert efficiency. Our system, FLEXMoE, guarantees 100% token efficiency and optimizes the allocation of computation resources to balance the workloads among GPUs, and is the closest to the ideal. With the training progressing, the imbalanced workloads are getting better due to the punishment of balance loss and all methods are moving towards to better efficiency.

5.5 Scalability

We also evaluate the scalability of FLEXMoE on 8, 16, 32 and 64 GPUs, which are conducted on a single MoE layer with 64 experts. The results are presented in Figure 7(b), which are normalized to the throughput of DeepSpeed-8GPUs. Results show FLEXMoE significantly outperforms DeepSpeed and FasterMoE. As experiments are conducted on a high-speed interconnected cluster (8*300Gbps intra-node and 8*200 Gbps inter-node), balanced computation among GPUs plays a key role and FLEXMoE is targeting as it.

6 CONCLUSION

In this paper, we presented FLEXMoE, a novel solution to address the dynamic imbalanced challenges encountered during the training of large-scale MoE models. By integrating a scheduling module on top of existing DNN frameworks, FLEXMoE monitors data traffic, creates scheduling plans, and dynamically adjusts the expert-to-device mapping during training. Our empirical results on six popular MoE models demonstrate that FLEXMoE outperforms DeepSpeed by an average of 1.70 \times and up to 2.10 \times , while also surpassing FasterMoE by an average of 1.30 \times and up to 1.45 \times .

7 ACKNOWLEDGMENTS

This work is supported by the National Key Research and Development Program of China (No. 2020AAA0105200), the National Natural Science Foundation of China (No. 61832001 and U22B2037) and PKU-Tencent joint research Lab. Bin Cui is the corresponding author.

REFERENCES

- [1] 2014. Azure A100 VMs. <https://docs.microsoft.com/en-us/azure/virtual-machines/nda100-v4-series>.
- [2] 2014. ImageNet-1K Dataset. <https://huggingface.co/datasets/imagenet-1k>.
- [3] 2014. Wikipedia Dataset. <https://huggingface.co/datasets/wikipedia>.
- [4] 2017. NCCL 2.0. <https://github.com/NVIDIA/nvcl>.
- [5] 2021. The AI2 Reasoning Challenge (ARC). <https://leaderboard.allenai.org/arc/submissions/public>.
- [6] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems* 33 (2020), 1877–1901.
- [7] Damai Dai, Li Dong, Shuming Ma, Bo Zheng, Zhifang Sui, Baobao Chang, and Furu Wei. 2022. StableMoE: Stable Routing Strategy for Mixture of Experts. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 7085–7095.
- [8] Ali Davoudian, Liu Chen, Hongwei Tu, and Mengchi Liu. 2021. A workload-adaptive streaming partitioner for distributed graph stores. *Data Science and Engineering* 6 (2021), 163–179.
- [9] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. 4171–4186.
- [10] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2021. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In *International Conference on Learning Representations*.
- [11] Nan Du, Yanping Huang, Andrew M Dai, Simon Tong, Dmitry Lepikhin, Yuanzhong Xu, Maxim Krikun, Yanqi Zhou, Adams Wei Yu, Orhan Firat, et al. 2022. Glam: Efficient scaling of language models with mixture-of-experts. In *International Conference on Machine Learning*. PMLR, 5547–5569.
- [12] William Fedus, Barret Zoph, and Noam Shazeer. 2022. Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity. *Journal of Machine Learning Research* 23, 120 (2022), 1–39.
- [13] Fangcheng Fu, Yuzheng Hu, Yihan He, Jiawei Jiang, Yingxia Shao, Ce Zhang, and Bin Cui. 2020. Don't waste your bits! squeeze activations and gradients for deep neural networks via tinscript. In *International Conference on Machine Learning*. PMLR, 3304–3314.
- [14] Fangcheng Fu, Jiawei Jiang, Yingxia Shao, and Bin Cui. 2019. An Experimental Evaluation of Large Scale GBDT Systems. *Proc. VLDB Endow.* 12, 11 (2019), 1357–1370.
- [15] Fangcheng Fu, Xupeng Miao, Jiawei Jiang, Huanran Xue, and Bin Cui. 2022. Towards Communication-efficient Vertical Federated Learning Training via Cache-enabled Local Update. *Proc. VLDB Endow.* 15, 10 (2022), 2111–2120.
- [16] Jia-Ke Ge, Yan-Feng Chai, and Yun-Peng Chai. 2021. WATuning: a workload-aware tuning system with attention-based deep reinforcement learning. *Journal of Computer Science and Technology* 36, 4 (2021), 741–761.
- [17] Jiaao He, Jidong Zhai, Tiago Antunes, Haojie Wang, Fuwen Luo, Shangfeng Shi, and Qin Li. 2022. FasterMoE: modeling and optimizing training of large-scale dynamic pre-trained models. In *Proceedings of the 27th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*. 120–134.
- [18] Yanping Huang, Youlong Cheng, Ankur Bapna, Orhan Firat, Dehao Chen, Mia Chen, HyoukJoong Lee, Jiquan Ngiam, Quoc V Le, Yonghui Wu, et al. 2019. Gpipe: Efficient training of giant neural networks using pipeline parallelism. *Advances in neural information processing systems* 32 (2019), 103–112.
- [19] Changho Hwang, Wei Cui, Yifan Xiong, Ziyue Yang, Ze Liu, Han Hu, Zilong Wang, Rafael Salas, Jithin Jose, Prabhat Ram, Joe Chau, Peng Cheng, Fan Yang, Mao Yang, and Yongqiang Xiong. 2022. Tutel: Adaptive Mixture-of-Experts at Scale. *CoRR abs/2206.03382* (2022).
- [20] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. *CoRR abs/2001.08361* (2020).
- [21] Dmitry Lepikhin, HyoukJoong Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping Huang, Maxim Krikun, Noam Shazeer, and Zhifeng Chen. 2021. GShard: Scaling Giant Models with Conditional Computation and Automatic Sharding. In *International Conference on Learning Representations*.
- [22] Mike Lewis, Shruti Bhosale, Tim Dettmers, Naman Goyal, and Luke Zettlemoyer. 2021. Base layers: Simplifying training of large, sparse models. In *International Conference on Machine Learning*. PMLR, 6265–6274.
- [23] Shen Li, Yanli Zhao, Rohan Varma, Omkar Salpekar, Pieter Noordhuis, Teng Li, Adam Paszke, Jeff Smith, Brian Vaughan, Pritam Damania, et al. 2020. PyTorch Distributed: Experiences on Accelerating Data Parallel Training. *Proceedings of the VLDB Endowment* 13, 12 (2020), 3005–3018.
- [24] Junyang Lin, Rui Men, An Yang, Chang Zhou, Ming Ding, Yichang Zhang, Peng Wang, Ang Wang, Le Jiang, Xianyan Jia, et al. 2021. M6: A chinese multimodal pretrainer. *CoRR abs/2103.00823* (2021).

- [25] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. 2021. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 10012–10022.
- [26] Zixuan Ma, Jiaao He, Jiezhong Qiu, Huanqi Cao, Yuanwei Wang, Zhenbo Sun, Liyan Zheng, Haojie Wang, Shizhi Tang, Tianyu Zheng, et al. 2022. BaGuaLu: targeting brain scale pretrained models with over 37 million cores. In *Proceedings of the 27th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*. 192–204.
- [27] Atsuki Maruta and Makoto P Kato. 2022. Intent-Aware Data Visualization Recommendation. *Data Science and Engineering* 7, 4 (2022), 301–315.
- [28] Xupeng Miao, Xiaonan Nie, Hailin Zhang, Tong Zhao, and Bin Cui. 2023. Hetu: A highly efficient automatic parallel distributed deep learning system. *Science China Information Sciences* 66, 1 (2023), 1–2.
- [29] Xupeng Miao, Yining Shi, Hailin Zhang, Xin Zhang, Xiaonan Nie, Zhi Yang, and Bin Cui. 2022. HET-GMP: A Graph-based System Approach to Scaling Large Embedding Model Training. In *Proceedings of the 2022 International Conference on Management of Data*. 470–480.
- [30] Xupeng Miao, Yujie Wang, Youhe Jiang, Chunan Shi, Xiaonan Nie, Hailin Zhang, and Bin Cui. 2022. Galvatron: Efficient Transformer Training over Multiple GPUs Using Automatic Parallelism. *Proceedings of the VLDB Endowment* 16, 3 (2022), 470–479.
- [31] Xupeng Miao, Hailin Zhang, Yining Shi, Xiaonan Nie, Zhi Yang, Yangyu Tao, and Bin Cui. 2021. HET: scaling out huge embedding model training via cache-enabled distributed framework. *Proceedings of the VLDB Endowment* 15, 2 (2021), 312–320.
- [32] Deepak Narayanan, Aaron Harlap, Amar Phanishayee, Vivek Seshadri, Nikhil R Devanur, Gregory R Ganger, Phillip B Gibbons, and Matei Zaharia. 2019. PipeDream: generalized pipeline parallelism for DNN training. In *Proceedings of the 27th ACM Symposium on Operating Systems Principles*. 1–15.
- [33] Deepak Narayanan, Mohammad Shoeybi, Jared Casper, Patrick LeGresley, Mostofa Patwary, Vijay Korthikanti, Dmitri Vainbrand, Prethvi Kashinkunti, Julie Bernauer, Bryan Catanzaro, et al. 2021. Efficient large-scale language model training on gpu clusters using megatron-lm. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*. 1–15.
- [34] Xiaonan Nie, Shijie Cao, Xupeng Miao, Lingxiao Ma, Jilong Xue, Youshan Miao, Zichao Yang, Zhi Yang, and Bin Cui. 2021. Dense-to-sparse gate for mixture-of-experts. *CoRR* abs/2112.14397 (2021).
- [35] Xiaonan Nie, Yi Liu, Fangcheng Fu, Jinbao Xue, Dian Jiao, Xupeng Miao, Yangyu Tao, and Bin Cui. 2023. Angel-PTM: A Scalable and Economical Large-scale Pre-training System in Tencent. *arXiv preprint arXiv:2303.02868* (2023).
- [36] Xiaonan Nie, Xupeng Miao, Zhi Yang, and Bin Cui. 2022. TSPPLIT: Fine-grained GPU Memory Management for Efficient DNN Training via Tensor Splitting. In *International Conference on Data Engineering*. IEEE, 2615–2628.
- [37] Xiaonan Nie, Pinxue Zhao, Xupeng Miao, Tong Zhao, and Bin Cui. 2022. HetuMoE: An Efficient Trillion-scale Mixture-of-Expert Distributed Training System. *CoRR* abs/2203.14685 (2022).
- [38] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems* 32 (2019), 8026–8037.
- [39] Yun Peng, Byron Choi, and Jianliang Xu. 2021. Graph learning for combinatorial optimization: a survey of state-of-the-art. *Data Science and Engineering* 6, 2 (2021), 119–141.
- [40] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language Models are Unsupervised Multitask Learners. *OpenAI blog* 1, 8 (2019), 9.
- [41] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *Journal of Machine Learning Research* 21 (2020), 1–67.
- [42] Samyam Rajbhandari, Conglong Li, Zhewei Yao, Minjia Zhang, Reza Yazdani Aminabadi, Ammar Ahmad Awan, Jeff Rasley, and Yuxiong He. 2022. DeepSpeed-MoE: Advancing Mixture-of-Experts Inference and Training to Power Next-Generation AI Scale. In *International Conference on Machine Learning*. PMLR, 18332–18346.
- [43] Carlos Riquelme, Joan Puigcerver, Basil Mustafa, Maxim Neumann, Rodolphe Jenatton, André Susano Pinto, Daniel Keysers, and Neil Houlsby. 2021. Scaling vision with sparse mixture of experts. *Advances in Neural Information Processing Systems* 34 (2021), 8583–8595.
- [44] Stephen Roller, Sainbayar Sukhbaatar, Jason Weston, et al. 2021. Hash layers for large sparse models. *Advances in Neural Information Processing Systems* 34 (2021), 17555–17566.
- [45] Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc V. Le, Geoffrey E. Hinton, and Jeff Dean. 2017. Outrageously Large Neural Networks: The Sparsely-Gated Mixture-of-Experts Layer. In *International Conference on Learning Representations*.
- [46] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017), 6000–6010.

- [47] Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2019. Superglue: A stickier benchmark for general-purpose language understanding systems. *Advances in neural information processing systems* 32 (2019), 3266–3280.
- [48] Zhuping Wang, Jingcheng Liu, Chao Huang, Hao Zhang, and Huaicheng Yan. 2023. Global output feedback adaptive stabilization for systems with long uncertain input delay. *Science China Information Sciences* 66, 1 (2023), 119201.
- [49] Hua-Peng Wei, Ying-Ying Deng, Fan Tang, Xing-Jia Pan, and Wei-Ming Dong. 2022. A Comparative Study of CNN-and Transformer-Based Visual Style Transfer. *Journal of Computer Science and Technology* 37, 3 (2022), 601–614.
- [50] Barret Zoph, Irwan Bello, Sameer Kumar, Nan Du, Yanping Huang, Jeff Dean, Noam Shazeer, and William Fedus. 2022. ST-MoE: Designing Stable and Transferable Sparse Expert Models. *CoRR* abs/2202.08906 (2022).

Received July 2022; revised October 2022; accepted November 2022