

Block-wise Optimization for Lottery-Ticket Adaptation in Large Language Models

Anonymous ACL submission

Abstract

We propose **Block-wise Optimization for Lottery-Ticket Adaptation (BOLA)**, a novel and simple sparse fine-tuning framework designed to enhance parameter efficiency in adapting large language models to new domains. Unlike conventional parameter-efficient fine-tuning (PEFT) methods such as LoRA and DoRA, which update all parameters, BOLA introduces a block-wise sparse selection mechanism. This mechanism searches for only domain-relevant subsets of the parameters and updates them. By integrating lottery ticket-style search with block-level granularity, BOLA mitigates catastrophic forgetting and enables interpretable, efficient adaptation while remaining compatible with existing PEFT techniques. Experiments on the math and commonsense reasoning benchmark demonstrate that BOLA achieves competitive performance with LoRA and DoRA. Our experiment code is available at <https://anonymous.4open.science/r/peft-B728>.

1 Introduction

With the advent of large language models (LLMs), it is no longer necessary to train separate models for each individual NLP task. Instead, a single general-purpose model can perform a wide range of tasks simply by providing suitable instructions. It is well-established that the performance of these models follows empirical scaling laws: as the number of parameters increases, model quality improves according to a power-law relationship (Kaplan et al., 2020). However, training such LLMs requires massive computational resources, imposing significant demands in terms of hardware and energy consumption. For example, if the number of trainable parameters is ϕ , the model state is 16 bits, and the optimizer state is 32 bits, then $16 \times \phi$ bytes of computing resources are required (Suhoi, 2024). Consequently, training an LLM with 8B parameters typically requires approximately 128 GB or

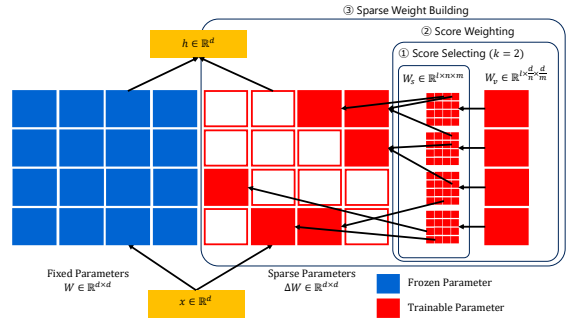


Figure 1: Overview of the proposed BOLA, which constructs trainable sparse weights $\Delta W \in \mathbb{R}^{d \times d}$ for fine-tuning, where W_s are block-wise score weights, and W_v are block-wise value weights. Blue indicates frozen parameters, red indicates trainable parameters, and white indicates zero parameters.

more of GPU memory. This substantial memory footprint poses a challenge for researchers and engineers attempting to fine-tune such models for specific domains or tasks.

To address this issue, parameter-efficient fine-tuning (PEFT) methods (Houlsby et al., 2019) have been proposed, that fine-tune pretrained models using only a minimal number of trainable parameters. Among these, reparameterized PEFT approaches such as LoRA (Hu et al., 2022) and DoRA (Liu et al., 2024a) employ low-rank decomposition to reduce the number of trainable parameters. Moreover, since these methods do not alter the model architecture, they have gained significant popularity due to their comparable performance to full fine-tuning but with reduced GPU memory requirements. Notably, QLoRA (Dettmers et al., 2023) enables the fine-tuning of 8B-scale LLMs using as little as 24 GB of GPU memory, depending on the precision and optimization techniques employed. These advancements significantly improve the accessibility and efficiency of adapting large-scale

models to a wide range of domains and tasks.

Since PEFT methods are dense methods that update all parameters of the target module, they are prone to catastrophic forgetting (Ramasesh et al., 2022; Dong et al., 2024; Luo et al., 2025), a phenomenon in which the pretrained knowledge of the original model is overwritten. However, it has been found that sparse methods (Panda et al., 2024; Xu and Zhang, 2024; Rios et al., 2025) also work well. Sparse methods offer several features over dense methods:

- They update only subsets of parameters, minimizing contamination of the original weights and helping to avoid catastrophic forgetting.
- They provide interpretability by identifying which parameters encode knowledge relevant to the target task or domain.

Sparse methods typically adapt models through a two-stage pipeline consisting of first searching for the optimal parameter subset, and then fine-tuning those parameters. However, this search phase introduces additional computational overhead, increasing the complexity of implementation.

To tackle this problem, we propose **Block-wise Optimization for Lottery-Ticket Adaptation (BOLA)** as shown in Figure 1, which is designed to simultaneously search for and update block-wise target parameters, inspired by Lottery Ticket Hypothesis (LTH) (Frankle and Carbin, 2019; Ramanujan et al., 2020). We evaluate BOLA across several language models and domains, demonstrating the relationship between training domains and training parameters through score maps. The experimental results show that BOLA is highly compatible with LoRA and DoRA, and can serve as a substitute for them despite being a sparse method. The sparsity of BOLA minimizes contamination of the original weights and helps to avoid catastrophic forgetting. The summary of our contributions are as follows:

- We introduce BOLA as a novel sparse-PEFT method capable of mitigating catastrophic forgetting during model merging, which is difficult to achieve with dense-PEFT methods such as LoRA and DoRA.
- We show that BOLA achieves comparable performance to LoRA and DoRA while reducing catastrophic forgetting.

- We demonstrate that BOLA can provide model interpretability by making it easier to identify where domain-specific knowledge is acquired by visualizing domain-specific parameter block.

2 Related Works

Lottery Ticket Hypothesis (LTH) (Frankle and Carbin, 2019) posits that within a randomly initialized neural network, there exist subnetworks that, when trained in isolation, can achieve test accuracy comparable to that of the original network. This hypothesis enables pruning pretrained models to arbitrary sizes. These subnetworks, called winning tickets, are identified by examining the magnitude of the trained dense weights to prune unimportant connections (Frankle and Carbin, 2019; Ramanujan et al., 2020). Several more-efficient methods have been proposed to identify winning tickets earlier and more effectively during the dense training phase (Chen et al., 2022; Yuan et al., 2025). Methods that perform training and pruning simultaneously have also been proposed (You et al., 2025).

Parameter-Efficient Fine-Tuning (PEFT) methods (Houlsby et al., 2019) are designed to reduce the high expense of fine-tuning large language models. Since the number of trainable parameters directly affects GPU memory consumption during training (Suhoi, 2024), PEFT methods attempt to reduce the number of trainable parameters using various techniques. These methods can be divided into several categories (Han et al., 2024), among which reparameterized PEFT is notable for its high performance and lack of extra inference burden. Reparameterized PEFT represents the change in fine-tuned weight W' change with $W' = W_o + \Delta W$. Moreover, the storage cost is reduced because it is sufficient to store only the weights ΔW . We further classify reparameterized PEFT into dense-PEFT and sparse-PEFT, based on whether trainable weights ΔW are constructed as a dense or sparse matrix.

Dense-PEFT employs a dense matrix for trainable weights $\Delta W \in \mathbb{R}^{d \times d}$. LoRA (Hu et al., 2022) applies low-rank matrix decomposition to approximate $\Delta W = BA$, where $B \in \mathbb{R}^{d \times r}$ and $A \in \mathbb{R}^{r \times d}$. The rank r is much smaller than the dimension d , which leads to a reduction in the number of trainable parameters from $2rd$ to d^2 . DoRA (Liu et al., 2024a) decomposes ΔW into its magnitude and direction as $\Delta W = m\Delta W'/\|\Delta W'\|$,

thereby enhancing performance by bridging the gap between the learning patterns of full fine-tuning (Full-FT) and LoRA while maintaining a low computational cost. Variants of LoRA (Edalati et al., 2022; Zhang et al., 2023; Kopiczko et al., 2024; Hayou et al., 2024) have also been proposed. Additionally, methods to construct the trainable weight ΔW using Fourier transforms (Liu et al., 2024b; Gao et al., 2024) or specially designed operations have been proposed (Jiang et al., 2024; Tan et al., 2024). These dense-PEFT methods update all parameters of the target module and leads to catastrophic forgetting (Ramasesh et al., 2022; Dong et al., 2024; Luo et al., 2025). Therefore, it is not suited for model merging that combines models from multiple domains.

Sparse-PEFT employs a sparse matrix for trainable weights $\Delta W \in \mathbb{R}^{d \times d}$. A sparse matrix ΔW is a matrix in which most elements are zero. The sparse function constructs ΔW from a set of indices I and corresponding values V , representing only non-zero entries, thereby inducing sparsity. LoTA (Panda et al., 2024) identifies the index through lottery ticket scores obtained in a prior stage before fine-tuning. Although SpARTA (Rios et al., 2025) randomly selects indices, they similarly require a prior stage before fine-tuning. These sparse-PEFT methods require computational or design cost to determine which parameters to update before fine-tuning.

Our method simultaneously performs the search and update of trainable parameters inspired by LTH. Therefore it does not require a prior stage to identify indices and can be easily applied as a replacement for popular dense-PEFT methods such as LoRA and DoRA. To the best of our knowledge, this is the first PEFT framework that simultaneously searches for and updates target parameters. We validate the performance and efficacy of our method through comprehensive experiments.

3 BOLA: Our Sparse PEFT

Figure 1 introduces our newly proposed sparse-PEFT, **Block-wise Lottery Ticket Adaptation (BOLA)**. To mitigate the computational burden of individually identifying important parameters, BOLA partitions the weights into several blocks and conducts exploration at the block level, as illustrated in Figure 1. BOLA searches for and updates block-wise target parameters, also known as lottery tickets (Frankle and Carbin, 2019), based on

the block-wise scores of weights. BOLA employs the reparametrization PEFT framework, where the fine-tuned weights $W' \in \mathbb{R}^{d \times d}$ can be represented as:

$$W' = W_o + \Delta W, \quad (1)$$

where $W_o \in \mathbb{R}^{d \times d}$ is the frozen pretrained weights, $\Delta W \in \mathbb{R}^{d \times d}$ is the trainable block sparse weights, d is the dimension of input and output. As depicted in Figure 1, the block sparse weight is one in which nonzero elements are clustered into dense blocks, while the remainder consists of zero blocks. To construct block sparse weights, both the values V and their corresponding insertion indices I are required. In this approach, since the sparse weights are constructed at the block level, the trainable weight ΔW can be represented as:

$$\Delta W = \text{BlockSparse}(I, V), \quad (2)$$

where $\text{BlockSparse}(I, V)$ is a function that generates a block sparse weight given indices $I \in \mathbb{R}^{l \times k}$ and values $V \in \mathbb{R}^{l \times k \times (d/n) \times (d/m)}$. In this context, l denotes the number of trainable blocks, k denotes the number of indices to insert for each trainable block, and n and m specify the division numbers along the input and output dimensions, respectively. The values with the same position are accumulated. Both I and V are derived from trainable score weights $W_s \in \mathbb{R}^{k \times n \times m}$ and trainable value weights $W_v \in \mathbb{R}^{k \times (d/n) \times (d/m)}$ as:

$$I, S = \text{TopK}(W_s)_{n \times m}, \quad (3)$$

$$P = \text{softmax}(S), \quad (4)$$

$$\hat{S} = P + \text{StopGradient}(S - P), \quad (5)$$

$$V = \alpha \cdot \hat{S} \cdot W_v, \quad (6)$$

where $\text{TopK}(\cdot)$ returns the k largest values and their indices, $\text{StopGradient}(\cdot)$ stops the gradient of input, and α is scalar factor. Equation 3 identifies the block positions with the k largest scores across the $n \times m$ grid for each candidate in $W_s \in \mathbb{R}^{l \times n \times m}$, where each score reflects the importance of a block. Assuming that blocks with larger scores correspond to more important block parameters, we identify the k largest indices $I \in \mathbb{R}^{l \times k}$ by computing the TopK over $W_s \in \mathbb{R}^{l \times n \times m}$ along with $n \times m$ dimension, which encodes the scores of each block. The corresponding values to be weighted by these scores and inserted at these positions are derived from $W_v \in \mathbb{R}^{l \times (d/n) \times (d/m)}$. That is, the indices

are obtained from W_s , while the values themselves are sourced from W_v . The same block can be selected multiple times, in which case its associated values are cumulatively aggregated. The weights W_s are initialized using a uniform Kaiming distribution (He et al., 2015) with mean of one, while W_v is initialized to zero. This initialization ensures that the trainable weights ΔW start from zero, similar to LoRA variants (Hu et al., 2022; Edalati et al., 2022; Zhang et al., 2023; Kopiczko et al., 2024; Hayou et al., 2024; Liu et al., 2024a), which is key for stable training (Liao et al., 2023).

Since this involves a sparse parameter update, gradients do not propagate through the score parameters W_s outside the top- k selection in Equation (3). Inspired by the motivation behind the Edge-popup algorithm (Ramanujan et al., 2020), we adopt straight-through estimation (Bengio et al., 2013), as in Equation (5), to allow gradients to flow through all score parameters, including those not selected in the forward pass. This enables dynamic block selection during training. The scores computed by Equation (5) are used to weight each block as described in Equation (6). Since normalizing the scores excessively reduces the block values and leads to unstable training, we employ the unnormalized scores for the weighting. As shown in Algorithm 2, this mechanism not only enables gradient flow through all parameters W_s but also allows the model to dynamically re-evaluate and update block importance throughout training, leading to more flexible and adaptive sparse representations.

The number of trainable parameters $\phi = lnm + ldd/nm = l(nm + dd/nm)$ is determined by the weights $W_s \in \mathbb{R}^{l \times n \times m}$ and $W_v \in \mathbb{R}^{l \times (d/n) \times (d/m)}$ where n and m are block size and k is the number of trainable blocks. Setting $nm = d$ yields $\phi = ld + ld = 2ld$ such that the number of trainable parameters is similar to $\phi = 2rd$ as in LoRA (Hu et al., 2022). As a result, this formulation enables straightforward replacement of LoRA, thereby positioning BOLA as a compatible and efficient alternative for parameter-efficient fine-tuning.

4 Experiments

We conduct a series of experiments to demonstrate the effectiveness of our method across various language models and tasks. Specifically, we use RoBERTa as a small-scale model and LLaMA3 as a large-scale model. For evaluation, we consider

```
def build_delta_weights(
    w_s: torch.Tensor,
    w_v: torch.Tensor,
    alpha: float = 1.0,
) -> torch.Tensor:
    l, n, m = w_s.shape
    l, d_o, d_i = w_v.shape
    s = w_s.view(1, -1)
    # STE-trick
    s_prb = torch.softmax(s, dim=-1)
    s_hat = s_prb + (s - s_prb).detach()
    # selecting
    topk = torch.topk(s_hat, k, dim=-1)
    # weighting
    scores = topk.values.view(1, -1, 1, 1)
    values = w_v.view(1, 1, d_o, d_i)
    values = alpha * scores * values
    # building
    indices = torch.stack(
        torch.unravel_index(
            topk.indices.view(-1),
            shape=(n, m)
        ),
        dim=0,
    ) # (2, l*k)
    values = values.view(-1, d_o, d_i) # (l*k, d_o, d_i)
    w_delta = block_sparse(indices, values)
    return w_delta
```

Figure 2: Pytorch-based Pseudocode

datasets from the GLUE benchmark, along with benchmarks targeting mathematical and common-sense reasoning. We adopt LoRA (Hu et al., 2022), DoRA (Liu et al., 2024a), and SpaRTA (Rios et al., 2025) as strong baselines. To ensure a fair comparison, we configured the hyperparameters such that the number of trainable parameters matched that of LoRA and SpaRTA. However, in DoRA, the number of trainable parameters was increased by the dimension d of the magnitude. Furthermore, we maintain hyperparameter settings as consistent as possible across different experimental conditions to facilitate reproducibility (see Appendix A for details). For model selection, we choose the model that achieved the highest validation performance among five different learning rates. All experiments are conducted on NVIDIA H100 80GB GPUs. We used four GPUs for RoBERTa-125M and eight GPUs for Llama3-8B. RoBERTa-125M is trained using float32 precision owing to its smaller size, while LLaMA3-8B is trained using bfloat16 precision to reduce memory consumption and accelerate training on H100 GPUs. All experiments are implemented using PyTorch, Hugging Face Transformers (Wolf et al., 2020), and the PEFT (Mangrulkar et al., 2022) library.

4.1 GLUE Benchmark

We evaluate BOLA in comparison with LoRA, DoRA, and SpaRTA using the relatively small language model RoBERTa-125M (Liu et al., 2019) for

Method	#Params.	MNL	SST2	MRPC	CoLA	QNLI	QQP	RTE	SST-B	Avg.
Full-FT*	125M	87.6	94.8	90.2	63.6	92.8	91.9	78.7	91.2	86.4
Full-FT	125M	86.3 \pm 0.0	93.5 \pm 0.3	90.0 \pm 0.4	62.9 \pm 1.4	92.2 \pm 0.1	91.4 \pm 0.0	78.7 \pm 1.1	91.0 \pm 0.0	85.2 \pm 0.5
LoRA	0.3M	87.5 \pm 0.0	95.0 \pm 0.1	89.2 \pm 0.2	63.1 \pm 0.4	92.7 \pm 0.0	90.9 \pm 0.1	77.1 \pm 1.3	90.9 \pm 0.0	85.8 \pm 0.2
DoRA	0.3M	84.3 \pm 0.2	94.2 \pm 0.1	89.3 \pm 0.6	64.6 \pm 0.5	91.6 \pm 0.3	87.7 \pm 0.0	79.1 \pm 1.4	91.0 \pm 0.1	85.2 \pm 0.1
SpaRTA	0.3M	86.7 \pm 0.1	94.9 \pm 0.3	89.5 \pm 0.1	60.9 \pm 0.2	92.5 \pm 0.2	90.4 \pm 0.0	76.7 \pm 1.1	90.5 \pm 0.2	85.3 \pm 0.1
BOLA	0.3M	86.9 \pm 0.1	94.8 \pm 0.3	89.5 \pm 0.1	62.2 \pm 1.3	92.8 \pm 0.2	90.8 \pm 0.0	76.8 \pm 1.0	90.6 \pm 0.0	85.2 \pm 0.2

Table 1: RoBERTa-125M with different adaptation methods on the GLUE benchmark. #Params. indicates that the number of trainable parameters excluded with classification head module. We report the overall (matched and mismatched) accuracy for MNLI, Matthew’s correlation for CoLA, Pearson correlation for STS-B, and accuracy for other tasks. For all metrics, higher values indicate better performance. The highest values are indicated in bold. * indicates numbers published in prior works.

Method	#Params.	MultiArith	GSM8K	AddSub	AQuA	SingleEq	SVAMP	Avg.
LoRA	56.6M	89.6 \pm 2.3	55.9 \pm 1.8	79.2 \pm 0.7	23.6 \pm 1.7	85.6 \pm 0.2	63.7 \pm 2.3	66.3 \pm 1.4
DoRA	57.4M	89.4 \pm 2.0	56.6 \pm 1.9	78.9 \pm 0.9	24.1 \pm 1.3	84.8 \pm 0.4	63.6 \pm 1.9	66.2 \pm 1.3
SpaRTA	56.6M	71.3 \pm 3.0	50.6 \pm 0.6	80.4 \pm 2.8	30.7 \pm 2.5	79.6 \pm 1.6	58.8 \pm 1.5	61.9 \pm 0.8
BOLA	56.6M	85.4 \pm 7.7	56.8 \pm 1.8	83.7 \pm 0.7	34.1 \pm 0.8	82.7 \pm 1.8	60.3 \pm 1.2	67.2 \pm 2.0

Table 2: Llama3-8B with different adaptation methods on the math benchmark. #Params. indicates that the number of trainable parameters. We report the accuracy on each dataset and its higher value is better for all metrics. The highest values are indicated in bold.

Method	#Params.	BoolQ	PIQA	SIQA	HS	WG	ARC-e	ARC-c	OBQA	Avg.
LoRA	56.6M	74.4 \pm 0.6	88.5 \pm 0.2	80.2 \pm 0.3	95.9 \pm 0.1	85.0 \pm 0.3	90.7 \pm 0.2	79.5 \pm 0.3	85.8 \pm 0.3	85.0 \pm 0.1
DoRA	57.4M	74.5 \pm 0.3	88.6 \pm 0.4	80.3 \pm 0.1	95.9 \pm 0.2	84.6 \pm 0.6	90.6 \pm 0.3	79.5 \pm 0.2	86.0 \pm 0.6	85.0 \pm 0.1
SpaRTA	56.6M	74.5 \pm 0.3	87.4 \pm 0.4	79.6 \pm 0.4	95.1 \pm 0.2	85.4 \pm 0.3	89.9 \pm 0.3	77.7 \pm 0.2	84.7 \pm 0.9	84.3 \pm 0.2
BOLA	56.6M	74.2 \pm 0.3	87.1 \pm 0.3	79.0 \pm 0.4	95.1 \pm 0.3	85.6 \pm 0.5	89.1 \pm 0.2	78.3 \pm 0.4	86.0 \pm 0.6	84.3 \pm 0.1

Table 3: Llama3-8B with different adaptation methods on the commonsense benchmark. #Params. indicates that the number of trainable parameters. We report the accuracy on each dataset and its higher value is better for all metrics. The highest values are indicated in bold. Note that HS and WG are abbreviations for HellaSwag and WinoGrande, respectively.

English GLUE benchmark (Wang et al., 2018). Although GLUE consists of eight downstream tasks, consistency with previous work (Hu et al., 2022), We report the mean over 3 random seeds; the results on eight downstream tasks and their average. We also report the number of trainable parameters, excluding the classification head to ensure a fair comparison of adaptation modules.

We present the results in Table 1. The results for each run are taken from the best epoch. On average, BOLA achieved performance comparable to Full-FT, but still fell slightly short of LoRA. This limitation arises from the fact that BOLA can update only kd parameters, which hinders the effective incorporation of domain knowledge into RoBERTa-125M.

4.2 Math and Commonsense Benchmark

We evaluate BOLA in comparison with LoRA, DoRA, and SpaRTA using the LLM LLaMA3-8B (Llama-Team, 2024) and benchmarks for math and commonsense reasoning. The math and commonsense benchmarks consist of eight and six tasks, respectively, each with predefined training and testing sets. Following the experimental setup of (Hu et al., 2023), we merge the training data from all eight commonsense tasks to construct a unified training dataset and evaluate the individual testing dataset for each task. We use the math 10K and commonsense 170K datasets¹ to fine-tune the models.

The results of the math benchmark are shown in Table 2. BOLA demonstrated improved performance over the baseline on the average. BOLA

¹<https://github.com/AGI-Edgerunners/LLM-Adapters>

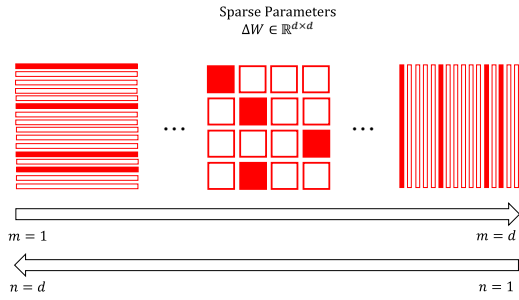


Figure 3: Transition diagram of block shapes determined by the hyperparameters n and m . Red and white blocks represent selected and unselected parameters, respectively. We define configurations as “structured” when either $n = d$ or $m = d$

outperforms SpARTA, another sparse method. The results of the commonsense benchmark are shown in Table 3. BOLA demonstrated improved performance over the baseline on the WG and OBQA tasks. On average, both of sparse methods BOLA and SpARTA achieves similar performance.

5 Ablation Study

In this section, we introduce the impact of hyperparameter in BOLA.

5.1 Structured vs. Unstructured

The hyperparameters n and m determine the number of block divisions and the input and output dimensions, thereby defining the shape of the trainable block weights. Figure 4 presents the performance trends as n is varied under the constraint $nm = d$. In this study, we define configurations as “structured” when either $n = d$ or $m = d$, meaning that block weights span the full input or output dimension. All other configurations are considered “unstructured”.

We evaluate the impact of block weight shapes on RoBERTa-125M performance using the GLUE benchmark, as shown in Figure 4. Overall, it is observed that increasing the block size n generally leads to improved performance. This indicates that a higher output dimensionality plays a crucial role, and therefore, it is preferable to set the value of n to its maximum. Since the structured configuration with the maximum n performs well across many downstream tasks, it is more effective to focus on tuning the number of base blocks k and the scaling

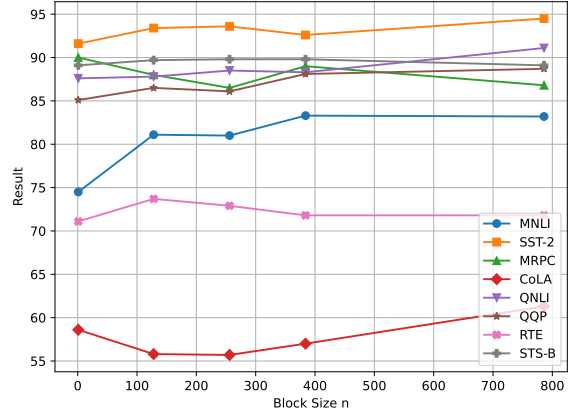


Figure 4: Evaluation results for various block sizes n for RoBERTa-125M on GLUE tasks. The horizontal axis corresponds to size m , which is varied from 1 to $d = 768$ under the constraint that $nm = d = 768$ remain constant. Each marker represents the evaluation result for a downstream task included in the GLUE benchmark. The number of selected blocks k is set 1.

factor α . This is analogous to tuning the rank r and scaling factor α in methods such as LoRA and DoRA.

6 Analysis of Sparse Update

In this section, we present two complementary benefits of sparse-PEFT compared with dense-PEFT, namely, enhanced interpretability and reduced performance dropout during model merging.

6.1 Interpretability of Model Structure

The sparsity of BOLA enhances model interpretability by making it easier to identify where domain-specific knowledge is acquired. Figure 5 shows the Llama3-8B (Llama-Team, 2024) weights ΔW of LoRA and BOLA, trained on math and commonsense datasets. LoRA updates all parameters, whereas BOLA searches for and updates only a small, important subset. Therefore, in BOLA, it becomes evident that most of the trainable weights ΔW are zero.

BOLA divides the weights into n and m blocks in the input and output dimensions, respectively, and only the k blocks with the highest scores are selected and updated. Interestingly, BOLA tends to utilize all of k blocks for the math domain. In contrast, BOLA tends to utilize significantly fewer blocks for the commonsense domain. This suggests that the model tends to acquire a greater amount of knowledge in the math domain, where its initial performance is weaker, while acquiring only

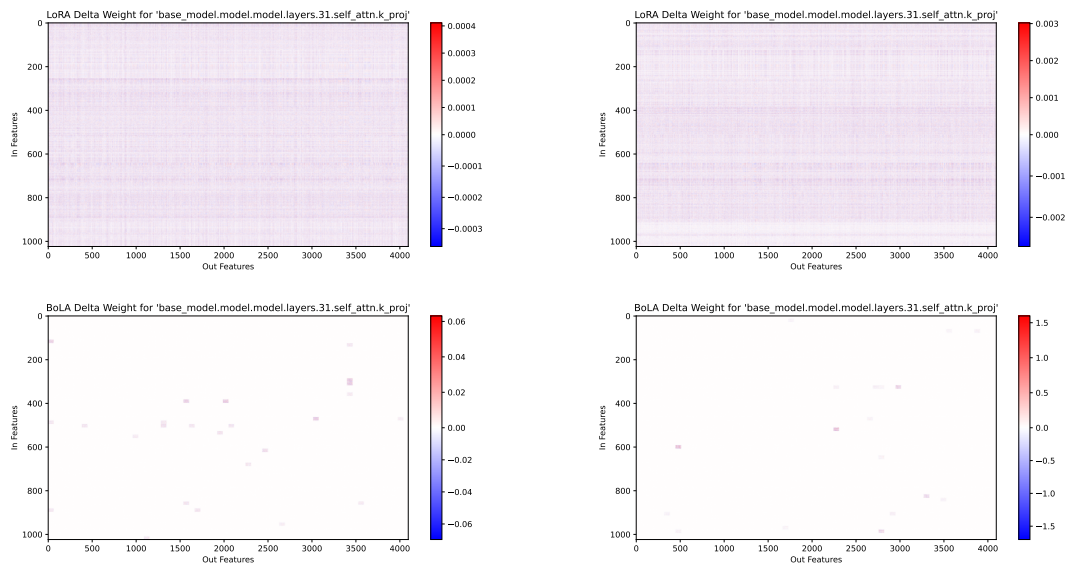


Figure 5: Visualization map of the trained weights ΔW in the last layer of Llama3-8B. The top shows LoRA weights and the bottom shows BOLA weights. The maps on the left and right were trained on mathematical and commonsense reasoning data, respectively. Color intensity indicates the magnitudes of the weights. In BOLA, although the number of selectable blocks is $k = 32$, it may appear to be fewer due to index duplication.

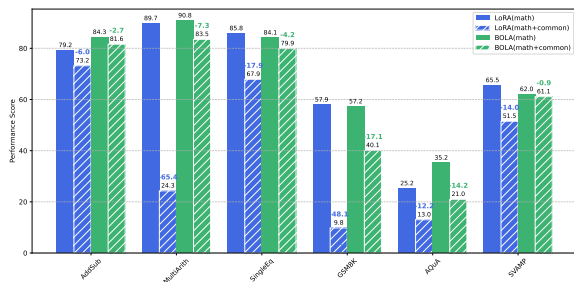


Figure 6: Performance comparison of model merging on the math benchmark. Bar plots show the performance of math-domain LoRA and BOLA models (filled), and their counterparts after merging with commonsense-domain models (hatched).

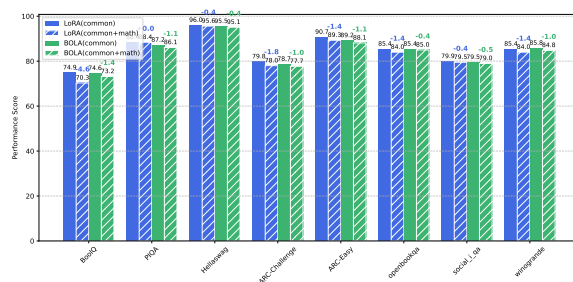


Figure 7: Performance comparison of model merging on the commonsense benchmark. Bar plots show the performance of commonsense-domain LoRA and BOLA models (filled), and their counterparts after merging with math-domain models (hatched).

435 limited knowledge in the commonsense domain, where it initially performs well.

436
437 As shown in Figure 5, while LoRA has full
438 access to all parameters, BOLA achieves com-
439 parable performance by accessing only $k/nm =$
440 $32/4096 = 0.78\%$ of the total parameters. There-
441 fore, only a small subset of parameters is crucial
442 for domain-specific adaptation in LLMs. This ob-
443 servation is consistent with prior work suggesting
444 that certain parameters play specialized roles (Yu
445 et al., 2025).

446 6.2 Sparse Model Merge

447 One notable property of sparse-PEFT is the mini-
448 mal overlap between weights trained on different

domains, which contributes to reduced interference
during model merging. LoTA (Panda et al., 2024)
leverages this property to demonstrate that cata-
strophic forgetting can be mitigated during model
merging. Catastrophic forgetting (Ramasesh et al.,
2022; Dong et al., 2024; Luo et al., 2025) refers
to the degradation in performance on previously
learned tasks that occurs when LLMs are fine-tuned
on multiple domains or tasks in a sequential manner.
Figures 6 and 7 show that sparse-PEFT BOLA sim-
ilarly reduces catastrophic forgetting during model
merging. We perform a linear model merge (Worts-
man et al., 2022), which is the most simple ap-
proach, the models used in this table were trained
as described in section 4. In the math benchmark,

449
450
451
452
453
454
455
456
457
458
459
460
461
462
463

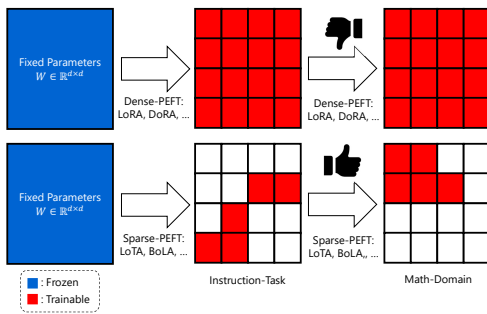


Figure 8: Workflow diagram illustrating model development using dense-PEFT and sparse-PEFT. During domain-specific training, dense-PEFT updates all model parameters, which may lead to overwriting the original instruct model. In contrast, sparse-PEFT selectively updates only a small subset of parameters, thereby preserving the instruct model and reducing the risk of unintended changes.

LoRA exhibits significantly degraded performance across all datasets when merging weights trained on math with those trained on commonsense, resulting in an average drop of 27.2 points. Notably, Multi-Arith and GSM8K experiences substantial performance degradation. In contrast, BOLA exhibits a performance improvement of 7.7 points on average after merging, demonstrating its robustness to catastrophic forgetting. No notable performance degradation occurs on MultiArith and GSM8K. In the commonsense benchmark, LoRA and BOLA experience significant performance degradation across all datasets when merging weights trained on math with those trained on commonsense. LoRA and BOLA exhibits only a minor performance declines of 1.4 and 0.9 points on average after merging. Since commonsense knowledge also encompasses information applicable to math tasks, it is expected to be less susceptible to forgetting due to overwriting. BOLA models trained independently on the math and commonsense domains can be merged without significant performance degradation.

Unlike dense-PEFT approaches, which often suffer from catastrophic forgetting and require costly fine-tuning, sparse-PEFT enables seamless integration of instruction-following and domain-specific capabilities. As shown in Figure 8, a key challenge in domain adaptation using dense-PEFT is that all parameters are overwritten during training, potentially erasing the model’s original instruction-following behavior. In contrast, sparse-PEFT up-

dates only a small subset of parameters, allowing domain adaptation while preserving the core characteristics of the instruction model.

Furthermore, if the instruction model itself is trained using sparse-PEFT, subsequent domain-specific adaptation can be achieved with even fewer parameter updates. This layered sparsity enables efficient and modular fine-tuning, making sparse-PEFT a highly practical approach for developing locally specialized models.

In summary, sparse-PEFT offers a robust and efficient approach to developing domain-specialized models without compromising the integrity or performance of the base model.

7 Conclusion

This work introduces **Block-wise Lottery Ticket Adaptation (BOLA)**, as a novel and simple sparse-PEFT for large language models designed to simultaneously search for and update parameters, a capability not found in previous methods. Furthermore, BOLA is highly compatible with popular dense-PEFT methods LoRA and DoRA, across various downstream tasks and model architectures. Therefore, BOLA serves as a costless alternative to LoRA and DoRA, as its ability to merge fine-tuned weights back into the pretrained model guarantees no extra inference overhead. Furthermore, unlike dense-PEFT, BOLA provides interpretability in knowledge acquisition and supports model merging without catastrophic forgetting. For future work, we wish to explore the generalizability of our method in other domains beyond audio and vision.

Limitations

Our PEFT approach demonstrates the potential to alleviate catastrophe forgetting during model merging. However, when the characteristics of the target domains differ significantly, there is no guarantee that the trained parameters will remain mutually exclusive. Another challenge is that our method has more hyperparameters to tune than existing PEFT approaches, making the adjustment process more difficult. As a future direction, we plan to extend the method to allow users to explore and update critical parameters within user-specified regions, enabling more effective control over parameter redundancy.

541
542
543
544
545

546
547
548
549
550

551
552
553
554

555
556
557
558
559
560

561
562
563
564

565
566
567

568
569
570
571

572
573
574
575

576
577
578
579

580
581
582
583

584
585
586
587
588
589

590
591
592
593
594

References

Yoshua Bengio, Nicholas Léonard, and Aaron Courville. 2013. [Estimating or propagating gradients through stochastic neurons for conditional computation](#). *Preprint*, arXiv:1308.3432.

Xiaohan Chen, Jason Zhang, and Zhangyang Wang. 2022. [Peek-a-boo: What \(more\) is disguised in a randomly weighted neural network, and how to find it efficiently](#). In *International Conference on Learning Representations*.

Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. [QLoRA: Efficient finetuning of quantized LLMs](#). In *Thirty-seventh Conference on Neural Information Processing Systems*.

Guanting Dong, Hongyi Yuan, Keming Lu, Chengpeng Li, Mingfeng Xue, Dayiheng Liu, Wei Wang, Zheng Yuan, Chang Zhou, and Jingren Zhou. 2024. [How abilities in large language models are affected by supervised fine-tuning data composition](#). *Preprint*, arXiv:2310.05492.

Ali Edalati, Marzieh Tahaei, Ivan Kobyzev, Vahid Par-tovi Nia, James J. Clark, and Mehdi Rezagholizadeh. 2022. [Krona: Parameter efficient tuning with kro-necker adapter](#). *Preprint*, arXiv:2212.10650.

Jonathan Frankle and Michael Carbin. 2019. [The lottery ticket hypothesis: Finding sparse, trainable neural networks](#). *Preprint*, arXiv:1803.03635.

Ziqi Gao, Qichao Wang, Aochuan Chen, Zijing Liu, Bingzhe Wu, Liang Chen, and Jia Li. 2024. [Parameter-efficient fine-tuning with discrete fourier transform](#). *Preprint*, arXiv:2405.03003.

Zeyu Han, Chao Gao, Jinyang Liu, Jeff Zhang, and Sai Qian Zhang. 2024. [Parameter-efficient fine-tuning for large models: A comprehensive survey](#). *Preprint*, arXiv:2403.14608.

Soufiane Hayou, Nikhil Ghosh, and Bin Yu. 2024. [LoRA+: Efficient low rank adaptation of large mod-els](#). In *Forty-first International Conference on Ma-chine Learning*.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2015. [Delving deep into rectifiers: Surpassing human-level performance on imagenet classification](#). *Preprint*, arXiv:1502.01852.

Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. [Parameter-efficient transfer learning for NLP](#). In *Proceedings of the 36th International Conference on Machine Learning*.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. [LoRA: Low-rank adaptation of large language models](#). In *International Conference on Learning Representations*.

Zhiqiang Hu, Lei Wang, Yihuai Lan, Wanyu Xu, Ee-Peng Lim, Lidong Bing, Xing Xu, Soujanya Poria, and Roy Lee. 2023. [LLM-adapters: An adapter fam-ily for parameter-efficient fine-tuning of large lan-guage models](#). In *Proceedings of the 2023 Confer-ence on Empirical Methods in Natural Language Pro-cessing*, pages 5254–5276, Singapore. Association for Computational Linguistics.

Ting Jiang, Shaohan Huang, Shengyue Luo, Zi-han Zhang, Haizhen Huang, Furu Wei, Weiwei Deng, Feng Sun, Qi Zhang, Deqing Wang, and Fuzhen Zhuang. 2024. [Mora: High-rank updat-ing for parameter-efficient fine-tuning](#). *Preprint*, arXiv:2405.12130.

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](#). *Preprint*, arXiv:2001.08361.

Dawid Jan Kopiczko, Tijmen Blankevoort, and Yuki M Asano. 2024. [VeRA: Vector-based random matrix adaptation](#). In *The Twelfth International Conference on Learning Representations*.

Baohao Liao, Shaomu Tan, and Christof Monz. 2023. [Make your pre-trained model reversible: From pa-rameter to memory efficient fine-tuning](#). *Preprint*, arXiv:2306.00477.

Shih-Yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, Kwang-Ting Cheng, and Min-Hung Chen. 2024a. [DoRA: Weight-decomposed low-rank adaptation](#). In *Forty-first In-ternational Conference on Machine Learning*.

Weiyang Liu, Zeju Qiu, Yao Feng, Yuliang Xiu, Yuxuan Xue, Longhui Yu, Haiwen Feng, Zhen Liu, Juyeon Heo, Songyou Peng, Yandong Wen, Michael J. Black, Adrian Weller, and Bernhard Schölkopf. 2024b. [Parameter-efficient orthogonal finetuning via butter-fly factorization](#). In *The Twelfth International Con-ference on Learning Representations*.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-dar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [Roberta: A robustly optimized bert pretraining ap-proach](#). *Preprint*, arXiv:1907.11692.

Llama-Team. 2024. [The llama 3 herd of models](#). *Preprint*, arXiv:2407.21783.

Yun Luo, Zhen Yang, Fandong Meng, Yafu Li, Jie Zhou, and Yue Zhang. 2025. [An empirical study of catas-trophic forgetting in large language models during continual fine-tuning](#). *Preprint*, arXiv:2308.08747.

Sourab Mangrulkar, Sylvain Gugger, Lysandre De-but, Younes Belkada, Sayak Paul, and Benjamin Bossan. 2022. [Peft: State-of-the-art parameter-efficient fine-tuning methods](#). <https://github.com/huggingface/peft>.

A Hyperparameters

727

For all experiments, we select the learning rate that achieves the highest validation performance from the range $\{1e-3, 5e-4, 1e-4, 5e-5, 1e-5\}$, and report the corresponding results.

728

729

Hyperparameter	LoRA	DoRA	SpaRTA	BOLA
Epochs			30	
Dropout			0.1	
Sequence Length			512	
Optimizer			AdamW	
LR	1e-3	1e-3	1e-3	1e-3
LR Scheduler			Linear	
Batch Size			64	
Warmup Ratio			0.06	
Weight Decay			0.1	
Target Modules			Q,K	
r	8	8		
(n, m)				(24, 32)
k			12,288	4
α	16	16	4	2
Dropout			0.1	

Table 4: Hyperparameter configurations of LoRA, DoRA, SpaRTA, and BOLA for RoBERTa-125M on the GLUE benchmark.

Hyperparameter	LoRA	DoRA	SpaRTA	BOLA
Epochs			3	
Dropout			0.1	
Sequence Length			512	
Optimizer			AdamW	
LR	1e-4	1e-4	5e-4	5e-4
LR Scheduler			Linear	
Batch Size			128	
Warmup Ratio			0.06	
Weight Decay			0.0	
Target Modules			Q,K,V,U,D	
r	32	32		
(n, m)				(64, 64)
k			353,894	4
α	64	64	4	2
Dropout			0.1	

Table 5: Hyperparameter configurations of LoRA, DoRA, SpaRTA, and BOLA for Llama3-8B on the math benchmark.

Hyperparameter	LoRA	DoRA	SpaRTA	BOLA
Epochs			3	
Dropout			0.1	
Sequence Length			512	
Optimizer			AdamW	
LR	1e-4	1e-4	1e-4	5e-4
LR Scheduler			Linear	
Batch Size			128	
Warmup Ratio			0.06	
Weight Decay			0.0	
Target Modules			Q,K,V,U,D	
r	32	32		
(n, m)				(64, 64)
k			353,894	32
α	64	64	4	2
Dropout			0.1	

Table 6: Hyperparameter configurations of LoRA, DoRA, SpaRTA, and BOLA for Llama3-8B on the commonsense benchmark.