

LLM-Barber: Block-Aware Rebuilder for Sparsity Mask in One-Shot for Large Language Models

Yupeng Su^{*1}, Ziyi Guan^{*2}, Xiaoqun Liu¹, Tianlai Jin¹, Dongkuan Wu¹,
Graziano Chesi², Ngai Wong², Hao Yu¹

¹School of Microelectronics, Southern University of Science and Technology, Shen Zhen, China
²Department of Electrical and Electronic Engineering, University of Hong Kong, Hong Kong, China

Abstract

Large language models (LLMs) have grown significantly in scale, leading to a critical need for efficient model pruning techniques. Existing post-training pruning techniques primarily focus on measuring weight importance on converged dense models to determine salient weights to retain. However, they often overlook the changes in weight importance during the pruning process, which can lead to performance degradation in the pruned models. To address this issue, we present **LLM-Barber** (Block-Aware Rebuilder for Sparsity Mask in One-Shot), a novel one-shot pruning framework that rebuilds the sparsity mask of pruned models without any retraining or weight reconstruction. LLM-Barber incorporates block-aware error optimization across Self-Attention and MLP blocks, ensuring global performance optimization. Inspired by the recent discovery of prominent outliers in LLMs, LLM-Barber introduces an innovative pruning metric that identifies weight importance using weights multiplied by gradients. Our experiments show that LLM-Barber can efficiently prune models like LLaMA and OPT families with 7B to 13B parameters on a single A100 GPU in just 30 minutes, achieving state-of-the-art results in both perplexity and zero-shot performance across various language benchmarks. Code is available at <https://github.com/YupengSu/LLM-Barber>.

Introduction

Large language models (LLMs) have become a cornerstone in natural language processing (NLP) due to their impressive performance on various tasks. However, as these models increase in size and complexity, their deployment poses significant challenges due to extensive computational and storage demands. For instance, models such as GPT-175B (Brown et al. 2020), with 175 billion parameters, require vast resources, making it impractical for many applications. Therefore, efficient model compression strategies are crucial for deploying these powerful models in practical applications.

Model compression techniques commonly employ quantization and pruning to enhance model efficiency. Quantization reduces the precision of model parameters, while pruning removes less critical parameters. Traditional pruning methods, such as magnitude-based pruning (Han et al.

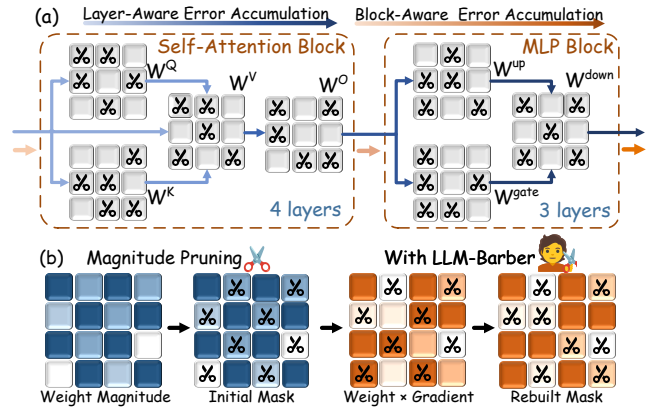


Figure 1: The benefits of integrating LLM-Barber into the pruning process: (a) Transition from the layer-aware to block-aware error accumulation to achieve an optimized global solution. (b) Rebuilding sparsity mask using a novel pruning metric based on weights multiplied by gradients.

2015), directly trim weights based on their absolute values. While effective for smaller models, these methods often struggle with large-scale LLMs, resulting in suboptimal sparsity due to their inability to capture the complex interactions and importance of weights in such massive scalability. To address these limitations, recent post-training pruning techniques such as SparseGPT and Wanda have emerged.

Current pruning methods face two major challenges. First, as depicted in Figure 1(a), traditional layer-aware pruning methods focus on individual layers and neglect inter-layer dependencies within, leading to higher error accumulation (represented by blue arrows). In contrast, block-aware pruning, by considering groups of layers, captures inter-layer interactions to reduce error accumulation (represented by orange arrows). Second, as shown in Figure 1(b), conventional methods typically build the pruning mask once, ignoring the changes of weight significance in post-pruning stage. This oversight can lead to improper identification of salient weights, resulting in performance degradation.

To address these limitations, we propose LLM-Barber, a novel and straightforward approach designed to rebuild sparsity mask of pruned networks without requiring for retraining or weight reconstruction. Firstly, unlike layer-aware methods that are confined to local optimization and thus

^{*}These authors contributed equally.

prone to significant error accumulation, as depicted by the blue arrows in Figure 1(a), LLM-Barber integrates pruning across both Self-Attention and MLP block. This approach mitigates error accumulation, as evidenced by the lighter orange arrows, facilitating global optimization and improved model performance. Secondly, LLM-Barber identifies weights that, although initially non-salient without a sparsity mask, gain significance in post-pruning. As shown in Figure 1(b), varying color shades represent the relative importance scores of different weights. LLM-Barber accurately identifies and rebuilds masks for high-score (deeply shaded) weights while pruning newly identified low-score (lightly shaded) weights. Thirdly, LLM-Barber utilizes a first-order Taylor series as pruning metric, which leverages weight-gradient multiplication for precise rebuilding of the sparsity mask, thereby enabling accurate pruning decisions and significantly reducing computational complexity compared to methods using second-order information. Finally, our one-shot approach surpasses iterative fine-tuning in efficiency while maintaining comparable accuracy. By rebuilding sparsity masks of pruned model at once, LLM-Barber offers a faster and more precise pruning solution. We validate the effectiveness of LLM-Barber through extensive experiments on the LLaMA and OPT families. Our method consistently outperforms existing post-training pruning techniques in both perplexity and zero-shot performance, establishing LLM-Barber as a state-of-the-art approach in this domain. To sum up, the key contributions are fourfold:

- **Block-Aware Global Optimization:** We are among the first to introduce a block-aware reconstruction problem that integrates sparsity across the Self-Attention and MLP blocks, achieving global optimization in pruning.
- **Rebuilding Sensitive Regions:** We identify non-salient weights that restore significance in post-pruning stage and rebuild sparsity mask to retain these weights, while simultaneously adjusting the mask to prune newly non-salient weights. This targeted rebuilding enhances overall model performance by optimally reallocating sparsity.
- **Innovative Pruning Metric:** We propose an innovative pruning metric based on the product of weights and gradients, leveraging first-order Taylor series for importance evaluation to reduce computational complexity compared to second-order Hessian-based approaches.
- **Versatility and Efficiency:** LLM-Barber demonstrates its efficiency and effectiveness across various pruning techniques, consistently achieving state-of-the-art performance in perplexity and zero-shot tasks, thereby establishing new yardsticks in LLM post-training pruning.

Related Work

LLMs Pruning and Sparsity

Network pruning reduces deep neural networks by removing unnecessary weights. For LLMs, pruning methods are divided into parameter efficient fine-tuning and post-training approaches. Parameter Efficient Fine-tuning (PEFT) begins with a initialized sparse network and refines it through iterative processes (Liu et al. 2019). LoRA (Hu et al.

2021) adapts pre-trained models to specific tasks or domains by injecting trainable rank decomposition matrices. Dynamic Sparse No Training (Zhang et al. 2024) minimizes the reconstruction error by iteratively pruning and growing weights. However, fine-tuning requires ample data and often leads to performance decline. Recent studies advance toward one-shot post-training pruning, showing substantial improvements. Post-training pruning removes weights from a pre-trained model. SparseGPT (Frantar and Alistarh 2023) uses Hessian-based metrics and subsequent residual weight updates, while Wanda (Sun et al. 2023) introduces a first-order pruning metric using weight-activation products. LLM-Barber employs a block-aware reconstruction approach and rebuilds masks with a novel pruning metric.

Model Compression Strategy

Compression is key to reducing the memory and computational demands of model. This work highlights the block-aware compression strategy over traditional layer-aware approaches. Layer-aware compression strategy began with Optimal Brain Damage (LeCun, Denker, and Solla 1989) and Optimal Brain Surgeon (Hassibi, Stork, and Wolff 1993). Building on these foundations, recent works like GPTQ (Frantar et al. 2023) further enhance layer-aware quantization using second-order information. Block-aware compression strategy generally offer better accuracy recovery in pruned models compared to layer-aware methods. For instance, APTQ (Guan et al. 2024) applies global quantization to attention mechanisms, enhancing model robustness, while BESA (Xu et al. 2024) uses block-wise sparsity allocation. Our method leverages block-aware pruning to optimize global performance across Self-Attention and MLP blocks, effectively balancing efficiency and accuracy.

Block-Aware Rebuilder for Sparsity Mask

Preliminaries

LLM pruning removes weights from dense networks to minimize output discrepancies, which is computationally intensive cross large-scale models, leading to address a layer-aware reconstruction problem (Hassibi, Stork, and Wolff 1993). This section reviews and reanalyses layer-aware reconstruction error and Taylor expansion at dense networks.

Layer-aware Reconstruction Error. For linear projection layer weight \mathbf{W} of shape (C_{out}, C_{in}) , where C_{out}, C_{in} indicates the output and input channels. With N calibration samples and sequence length L , the input activation is denoted as \mathbf{X} with the shape of $(C_{in}, N \times L)$. Layer-aware reconstruction error \mathbf{E} is defined as the ℓ_2 norm difference between output of dense and sparse layers:

$$\mathbf{E}(\widehat{\mathbf{W}}, \mathbf{X}) = \|\mathbf{W}\mathbf{X} - \widehat{\mathbf{W}}\mathbf{X}\|_2^2, \quad (1)$$

where $\widehat{\mathbf{W}}$ is the element-wise product of \mathbf{W} and a binary sparsity mask $\mathbf{M}(i, j) \in \{0, 1\}$ of shape (C_{out}, C_{in}) , in context with mask selection and without weight reconstruction:

$$\mathbf{E}(\mathbf{M}, \mathbf{X}) = \|\mathbf{W}\mathbf{X} - (\mathbf{W} \odot \mathbf{M})\mathbf{X}\|_2^2. \quad (2)$$

The objective is to search for an optimal sparsity mask \mathbf{M} where the pruning process reduces model complexity while preserving its predictive accuracy.

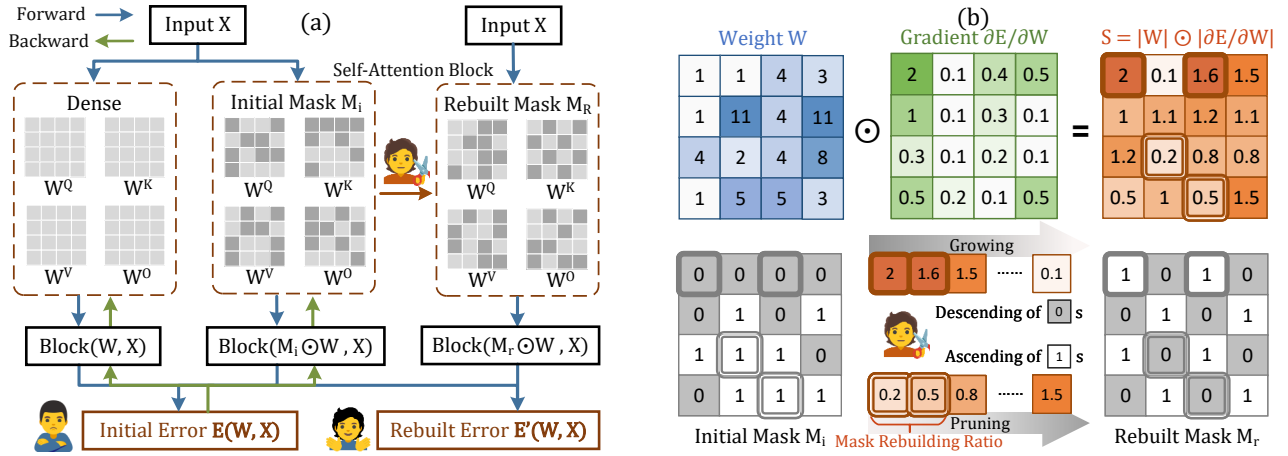


Figure 2: The workflow of LLM-Barber. (a) illustrates the process of block-aware reconstruction error calculation and gradient back-propagation for each linear weight. (b) shows pruning metric computation and procedure for rebuilding sparsity mask.

Taylor Expansion at Dense Networks. For a dense network $\widehat{W}_{\text{dense}}$ at a local minimum, the reconstruction error can be expanded into its Taylor series with respect to \widehat{W} , ignoring terms beyond second order:

$$E(\widehat{W}, X) = E(\widehat{W}_{\text{dense}}, X) + \frac{\partial E}{\partial W} \Delta W + \frac{1}{2} \Delta W^T H \Delta W. \quad (3)$$

Without a sparsity mask in dense model, we can simply assign all-one matrix to the mask M , thereby yielding the zeroth-order terms $E(\widehat{W}_{\text{dense}}, X) = 0$. The first-order derivative $\partial E / \partial W$ vanishes when training converged (Frantarh and Alistarh 2023; Hassibi, Stork, and Wolff 1993), leaving only the computationally expensive second-order terms involving a large Hessian matrix which are challenging for layer-wise reconstruction and channel-wise independence assumption.

Block-Aware Rebuilder for Sparsity Mask

In this work, we depart from existing post-training pruning methods in three key aspects: Firstly, to address the layer-aware reconstruction problem that leads to exponentially accumulating biases, we adopt a block-aware reconstruction error and apply a divide-and-conquer strategy to mitigate errors and computational costs. Secondly, inspired by the pruning-and-growing¹ operation (Mocanu et al. 2018; Zhang et al. 2024), we address the limitations of mask selection of pruning in dense networks due to the changeable significance of weights, by re-evaluating weight importance score in sparse networks and rebuilding the sparsity mask through targeted growth of salient weights and pruning of non-salient weights. Thirdly, our analysis reveals that with the advancement of LLMs, mask selection becomes increasingly critical in weight reconstruction. Thus, we prioritize the rebuilding of the sparsity mask and strip the reconstruction of weights. Based on these insights, we propose **LLM-Barber**, a **Block-Aware Rebuilder for Sparsity Mask in One-Shot** without any need for fine-tuning or retraining.

¹“Pruning” and “Growing” will both be used in the following text to refer to rebuild for sparsity masks. Pruning changes mask values from 1 to 0, whereas growing changes them from 0 to 1.

Block-Aware Reconstruction Error. Building on the definitions in Eq. (1) and Eq. (2), we define the block-aware reconstruction error for a Self-Attention or MLP block:

$$E(M, X) = \|\text{Block}(W, X) - \text{Block}(W \odot M, X)\|_2^2. \quad (4)$$

Evaluating reconstruction error across blocks, denoted as $\text{Block}(\cdot)$, allows us to achieve a globally optimal solution in Self-Attention and MLP blocks rather than layer-wise. This new block-aware reconstruction formulation offers a more effective and cohesive strategy for selecting sparsity masks.

Taylor Expansion at Sparse Networks. Migrating Eq. (3) at sparse networks $\widehat{W}_{\text{sparse}}$ with an initialization sparsity mask M_i , we can obtain the Taylor series expansion as:

$$E(\widehat{W}, X) = E(\widehat{W}_{\text{sparse}}, X) + \frac{\partial E}{\partial W} \Delta W + \frac{1}{2} \Delta W^T H \Delta W. \quad (5)$$

The zeroth-order term of the Taylor expansion represents the reconstruction error after mask initialization:

$$E(\widehat{W}_{\text{sparse}}, X) = \|\text{Block}(W, X) - \text{Block}(W \odot M_i, X)\|_2^2. \quad (6)$$

Assuming non-negligible zeroth-order terms, the first-order gradient in sparse networks remains significant even after convergence and can be efficiently accessed via PyTorch’s Autograd. First-order information provides computational efficiency and operates independently of any reconstruction error. Therefore, second-order terms can be omitted when significant gradients are present, leading to the following change in reconstruction error due to mask rebuilding:

$$\Delta E = (\partial E / \partial W) \cdot \Delta W, \quad (7)$$

which delineates the importance score of weights during sparsity mask rebuilding.

Pruning Metric. For a Self-Attention or MLP Block with weights W , first-order information suffices for block-aware reconstruction error in sparse networks. The change in weight magnitude during sparsity mask adjustment matches the weight’s original magnitude ($|\Delta W_{ij}| = |W_{ij}|$). We thus assess the impact of mask rebuilding on reconstruction error by computing the product of the weight’s magnitude and its gradient. The **importance score** for W_{ij} is:

$$S_{ij} = |W_{ij}| \cdot |(\partial E / \partial W_{ij})|, \quad (8)$$

where $|\cdot|$ represents the absolute value operator, and \mathbf{E} is denoted as the block-aware reconstruction error:

$$\mathbf{E}(\mathbf{M}, \mathbf{X}) = \|\text{Block}(\mathbf{W}, \mathbf{X}) - \text{Block}(\mathbf{W} \odot \mathbf{M}, \mathbf{X})\|_2^2. \quad (9)$$

This method prioritizes weights with both substantial magnitudes and significant gradients, allowing for the preservation of critical weights while pruning less important ones.

Pruning Granularity. Choosing the right pruning granularity is crucial (Sun et al. 2023). Traditional methods operate on a layer-wise (Han et al. 2015), input-wise (Frantar and Alistarh 2023), or output-wise (Sun et al. 2023) basis, which mainly address the layer-aware reconstruction problem, known for its output channel independence. Wanda’s output-wise ranking yields superior results compared to other methods. However, in LLM-Barber’s block-aware framework, output channel independence is no longer applicable. Thus, LLM-Barber extends consideration to block-wise granularity, prioritizing all linear layers within a block. Our analysis of the four distinct granularity levels shows that optimal granularity depends on the specific sparse mask initialization, with detailed results discussed in the Ablation Study.

Mask Rebuilding. With the block-aware reconstruction error \mathbf{E} , gradient information, sparsity metric, and granularity established, we proceed to rebuild the sparsity mask for each layer. Consider a cluster of weights \mathbf{W}^c under a specific sparsity granularity and its corresponding sparsity mask \mathbf{M}^c . We define the **growing criterion** and the **pruning criterion**:

$$gi, gj = \operatorname{argmax} |\mathbf{W}^c| \cdot |(\partial \mathbf{E} / \partial \mathbf{W}^c)|, \text{ if } \mathbf{M}^c_{gi, gj} = 0, \quad (10)$$

$$pi, pj = \operatorname{argmin} |\mathbf{W}^c| \cdot |(\partial \mathbf{E} / \partial \mathbf{W}^c)|, \text{ if } \mathbf{M}^c_{pi, pj} = 1. \quad (11)$$

The growing and pruning weights form a **mask rebuilding pair**, representing the interchange within the sparsity mask. The value of each pair is defined as the difference between the importance scores of the growing and pruning weights.

Our experiments show that LLM-Barber identifies varying proportions of salient weights depending on the sparse mask’s initialization method. To control the extent of mask rebuilding, we introduce a hyperparameter α , called the mask rebuilding ratio. The number of mask rebuilding pairs N is calculated as:

$$N = \{i \mid \mathbf{S}_i^{\text{grow}} - \mathbf{S}_i^{\text{prune}} > 0\} \cdot \alpha, \quad (12)$$

where $\mathbf{S}^{\text{grow}} - \mathbf{S}^{\text{prune}}$ represents the value of the mask rebuilding pairs, where \mathbf{S}^{grow} is arranged in descending order, and $\mathbf{S}^{\text{prune}}$ in ascending order. The subscript i denotes the number of values that exceeds zero, indicating that the growing weight is more important than the pruning weight.

Procedure

Like most post-training pruning methods (Han et al. 2015; Sun et al. 2023), LLM-Barber is executed within a single global LLM forward pass, with local backward passes for gradient computation in each block. Figure 2 illustrates the LLM-Barber workflow, which consists of four stages:

- **Sparsity Mask Initialization.** We apply a post-training pruning technique to initialize a preliminary sparsity mask from the dense network.
- **Block-aware Reconstruction Error Computation.** We use a block-aware reconstruction error to evaluate the discrepancy between the dense and sparse model outputs.

Algorithm 1: Pseudocode of LLM-Barber.

Input: Calibration samples \mathbf{X} , a block’s weights $\{\mathbf{W}^l\}_{l=1}^L$ and initial masks $\{\mathbf{M}_i^l\}_{l=1}^L$, mask rebuilding ratio α .

Output: Rebuilt sparsity masks $\{\mathbf{M}_r^l\}_{l=1}^L$.

```

1:  $\mathbf{E}_i \leftarrow \text{Block}(\mathbf{W}, \mathbf{X}) - \text{Block}(\mathbf{W} \odot \mathbf{M}_i, \mathbf{X})$ 
2:  $\{\mathbf{G}^l\}_{l=1}^L \leftarrow$  backpropagation for gradients via  $\mathbf{E}_i$ .
3: for  $l$  in  $\{1, 2, \dots, L\}$  do
4:    $\mathbf{M}_r^l \leftarrow \mathbf{M}_i^l$   $\triangleright$  Initialize rebuilt sparsity mask.
5:    $\mathbf{S}^l \leftarrow |\mathbf{W}^l| \cdot |\mathbf{G}^l|$   $\triangleright$  Obtain importance score.
6:    $N \leftarrow \mathbf{S}^l, \alpha$  via Eq.(12)  $\triangleright$  Obtain rebuilding number.
7:   for  $n$  in  $\{1, 2, \dots, N\}$  do
8:     Obtain growing index  $gi, gj$  via Eq.(10)
9:     Obtain pruning index  $pi, pj$  via Eq.(11)
10:     $\mathbf{M}_r^{pi, pj} = 1$   $\triangleright$  Weight Growing.
11:     $\mathbf{M}_r^{gi, gj} = 0$   $\triangleright$  Weight Pruning.
12:   end for
13: end for
14:  $\mathbf{E}_r \leftarrow \text{Block}(\mathbf{W}, \mathbf{X}) - \text{Block}(\mathbf{W} \odot \mathbf{M}_r, \mathbf{X})$ 
15: Identify improvement by comparing error  $\mathbf{E}_i$  and  $\mathbf{E}_r$ .
16: return rebuilt sparsity masks  $\{\mathbf{M}_r^l\}_{l=1}^L$ .

```

- **Back-propagation for Gradients.** Gradients are automatically derived via back-propagation, and the product of weights and gradients serves as the pruning metric.

- **Sparsity Mask Rebuilding.** Masks are sorted based on the pruning metric, unpruned weights in ascending order while pruned weights in descending order. We rebuild the weight masks by growing newly significant weights and pruning those became non-salient with a certain percentage named mask rebuilding ratio.

Ultimately, we recalculated the block-aware reconstruction error to assess the enhancement brought by our LLM-Barber. Experimental results reveal a substantial decrease in error, markedly decelerating the rate of error accumulation.

Structured N:M Sparsity

While LLM-Barber primarily targets unstructured sparsity, it can be adapted for structured N:M sparsity. Here, groups of M weights are pruned to retain only N non-zero weights. During mask rebuilding, LLM-Barber divides each M -group into N pairs (each pair contains one pruned and one non-pruned weight), then sorts these pairs by output channel to determine which are mask rebuilding pairs. This method creates an optimal sparse mask, leveraging N:M sparsity while maintaining model performance.

Experiment

Experiment Settings

Setup. LLM-Barber is implemented in Pytorch and utilized public model checkpoints from the HuggingFace library² on a single 80GB NVIDIA A100 GPU. After mask initialization, LLM-Barber uniformly rebuilds sparsity masks in sequence, performing in one-shot without any fine-tuning.

²huggingface.co/meta-llama, huggingface.co/huggyllama

	LLaMA1		LLaMA2		LLaMA3	OPT	
Method	7B	13B	7B	13B	8B	6.7B	13B
Dense	5.677	5.091	5.472	4.884	6.136	10.86	10.13
Magnitude	17.26	20.14	16.03	6.827	205.5	9.7e2	1.2e4
w/ LLM-Barber	7.332	<u>6.089</u>	7.170	5.955	10.98	13.12	15.52
SparseGPT	7.201	6.194	7.005	6.036	9.399	<u>11.59</u>	<u>11.15</u>
SparseGPT w/o WR	7.545	6.311	7.413	6.134	9.994	13.13	15.76
w/ LLM-Barber	7.159	6.125	7.004	5.929	9.348	11.95	11.93
Wanda	7.254	6.152	6.920	5.972	9.821	11.98	11.93
w/ LLM-Barber	<u>7.118</u>	6.091	<u>6.868</u>	<u>5.918</u>	9.451	11.95	11.71

Table 1: WikiText-2 perplexity comparison for pruning LLM models at 50% sparsity rate. **Bold** results show improvements of integrating LLM-Barber. Underscored results indicate best performance in each LLM. WR represents weight reconstruction.

Models & Datasets. LLM-Barber is evaluated on the LLaMA family, including LLaMA-7B/13B (Touvron et al. 2023a), LLaMA2-7B/13B (Touvron et al. 2023b), and LLaMA3-8B (AI@Meta 2024), as well as the OPT model series: OPT-6.7B/13B (Zhang et al. 2022). Notably, LLM-Barber is broadly applicable to any Transformer-based LLMs with Self-Attention and MLP blocks. Following previous works, we use 128 segments of 2048 tokens from the C4 dataset (Raffel et al. 2020) for mask rebuilding.

Evaluation. To comprehensively assess LLM-Barber, we conduct rigorous evaluations on *perplexity* and *zero-shot accuracy*. Perplexity is measured on the validation sets of benchmarks such as WikiText-2 (Merity et al. 2017), PTB (Marcus et al. 1994), and C4 (Raffel et al. 2020). Zero-shot accuracy is assessed using the EleutherAI LM Harness (Gao et al. 2023) across six benchmarks: BoolQ (Clark et al. 2019), RTE (Wang et al. 2018), HellaSwag (Zellers et al. 2019), ARC Easy and Challenge (Clark et al. 2018), and OpenbookQA (Mihaylov et al. 2018).

Baselines. The results of LLM-Barber are compared with the following established post-training pruning methods:

- Magnitude Pruning (Han et al. 2015) eliminates weights based only on their magnitudes;
- SparseGPT (Frantar and Alistarh 2023) identifies weights importance by using second-order information;
- Wanda (Sun et al. 2023) determines weights to be pruned by the weight magnitude multiplied by input activation.

Main Results

Unstructured Sparsity. LLM-Barber effectively prunes the LLaMA and OPT models, achieving 50% unstructured sparsity without requiring supplementary weight reconstruction, as detailed in Table 1. LLM-Barber demonstrates the capability to rebuild the sparsity masks initialized by other pruning methods in a single forward pass, significantly outperforming conventional pruning baselines. In the LLaMA3-8B model, LLM-Barber creates a new sparsity mask that reduces perplexity to 9.451, a substantial improvement over the Wanda baseline of 9.821. Notably, LLM-Barber achieves robust improvements even with poorly performing initial

sparsity masks, such as magnitude pruning, where it reduces perplexity from 205.5 to 10.98, which is an impressive and substantial enhancement.

Sparsity	60%	70%	80%	90%
Magnitude	3.39e4	1.62e6	8.55e7	2.26e7
w/ LLM-Barber	28.14	2.08e2	1.09e3	7.55e4
Wanda	23.57	1.28e2	8.54e2	1.28e4
w/ LLM-Barber	22.04	1.07e2	5.70e2	7.07e3

Table 2: WikiText-2 perplexity performance for pruning LLaMA3-8B at varying sparsity rate.

Varying Sparsity Levels. We conduct experiments on varying sparsity levels for unstructured pruning in LLaMA3-8B as shown in Table 2. LLM-Barber consistently increases perplexity across all initialization methods, with magnitude pruning showing the most significant improvement.

Method	Sparsity	V1-7B	V1-13B	V2-7B
Magnitude	4:8	16.83	13.72	15.91
w/ LLM-Barber	4:8	8.852	7.137	9.345
SparseGPT	4:8	8.608	7.437	8.495
w/ LLM-Barber	4:8	8.191	7.085	8.003
Magnitude	2:4	42.56	18.32	37.76
w/ LLM-Barber	2:4	11.04	9.006	13.47
SparseGPT	2:4	11.55	9.116	10.94
w/ LLM-Barber	2:4	10.14	8.517	9.806

Table 3: WikiText-2 perplexity comparison for pruning LLaMA family with structured N:M pattern.

Structured N:M Sparsity. In contrast to unstructured sparsity, employing N:M fine-grained sparsity can provide more tangible acceleration benefits when leveraging NVIDIA Ampere’s sparse tensor cores (Choquette et al.

Model	Method	BoolQ	RTE	HellaSwag	ARC-e	ARC-c	OBQA	Mean
LLaMA-7B	Dense	75.08	66.79	56.96	75.29	41.89	34.40	58.40
	Magnitude	54.61	54.51	45.47	58.75	33.45	22.60	44.89
	w/ LLM-Barber	71.47	59.93	50.63	69.11	35.41	28.20	52.46
	Wanda	70.98	55.23	<u>51.90</u>	69.44	36.86	28.60	52.17
	w/LLM-Barber	73.12	56.68	51.80	70.32	36.95	28.80	52.95
LLaMA2-13B	Dense	80.58	65.34	60.05	79.38	48.46	35.20	61.50
	Magnitude	70.53	55.96	54.42	57.68	38.40	27.80	50.79
	w/ LLM-Barber	81.10	60.29	55.67	75.34	40.28	31.20	57.31
	Wanda	80.95	60.28	<u>56.98</u>	76.30	<u>42.26</u>	31.20	57.99
	w/ LLM-Barber	80.98	62.82	55.99	76.39	42.13	31.40	58.29
LLaMA3-8B	Dense	81.35	69.68	60.19	80.09	50.60	34.80	62.79
	Magnitude	42.87	53.07	29.85	46.59	25.09	22.00	36.57
	w/ LLM-Barber	72.72	54.87	51.00	72.05	37.46	27.40	52.58
	Wanda	78.41	60.29	<u>51.20</u>	71.38	<u>40.10</u>	<u>29.40</u>	55.13
	w/LLM-Barber	78.59	61.37	51.01	71.46	39.85	29.00	55.21

Table 4: Zero-shot performance comparison for pruning LLaMA series on six evaluated tasks at 50% sparsity rate. **Bold** results show improvements of integrating LLM-Barber. Underscored results indicate the best performance in each zero-shot tasks.

2021). Therefore, we also evaluate the effectiveness of our LLM-Barber in partial LLaMA models on the N:M fine-grained sparsity pattern as shown in Table 3.

Zero-shot Performance. Following previous works (Frantar and Alistarh 2023; Sun et al. 2023), we evaluated the LLaMA models on six diverse zero-shot tasks. The results are summarized in Table 4, where models are pruned to unstructured 50% sparsity. Averaging the accuracy across the six evaluated tasks, it becomes apparent that LLM-Barber possesses the capability to identify a more effective network than those obtained via the initialization methods. For tasks such as BoolQ, RTE, and ARC-e, LLM-Barber consistently outperforms the baseline pruning techniques across the entire LLaMA model suite. However, it is worth noting that there is no single universally superior performer for the remaining tasks in the evaluation set, with the initial pruning methods sometimes matching or even marginally exceeding the results obtained through LLM-Barber.

Mask Rebuilding Ratio Selection

A critical aspect of the LLM-Barber method lies in determining the optimal mask rebuilding ratio to achieve peak accuracy. One effective strategy involves analyzing the distribution of value magnitudes within the mask rebuilding pairs, corresponding to differences between the growing and pruning importance score. Fortunately, this distribution often exhibits a distinct pattern of outliers, facilitating rapid identification of an appropriate mask rebuilding ratio.

We have plotted the score distribution of mask rebuilding pairs for various pruning granularity and initialization methods on LLaMA-7B, as well as perplexity performance corresponding to mask rebuilding ratios in Figure 3. The results reveal a strong correlation between the distribution of out-

liers in mask rebuilding pairs and the optimal mask rebuilding ratio. For instance, outliers are significantly distributed within the top 10% with Magnitude pruning, thus selecting a 10% mask rebuilding ratio yields the overall optimal solution. Similarly, with the initialization method Wanda, a significant outlier distribution around the 1% mark corresponds to optimal results near a 1% mask rebuilding ratio. Thus, LLM-Barber can preemptively narrow the search range for the optimal mask rebuilding ratio by analyzing outlier distributions. This approach allows LLM-Barber to adapt flexibly to various reconstruction masks without extensive searching.

It is worth noting that LLM-Barber identifies a larger proportion of outliers in less effective initial masks, which corresponds to a more substantial rebuilding of the mask. This is why LLM-Barber provides a more significant improvement for masks with poorer initial sparsity mask with a more aggressive mask rebuilding strategy.

Ablation Study

Given the outstanding potential of LLM-Barber, we have comprehensively analyzed the impact of three critical factors to assess the robustness and effectiveness of our LLM-Barber method: pruning granularity, pruning metric, and calibration data size. These factors were chosen to gain a deeper understanding of how different configurations influence the performance of pruned models and to demonstrate LLM-Barber’s versatility across various settings.

Pruning Granularity. LLM-Barber dynamically selects different pruning granularities to adapt to varying initialization methods. In this paper, we evaluated the impact of four levels of granularities: block-wise, layer-wise, input-wise, and output-wise pruning, as shown in Table 5. For magnitude pruning, the block-wise granularity yields the best

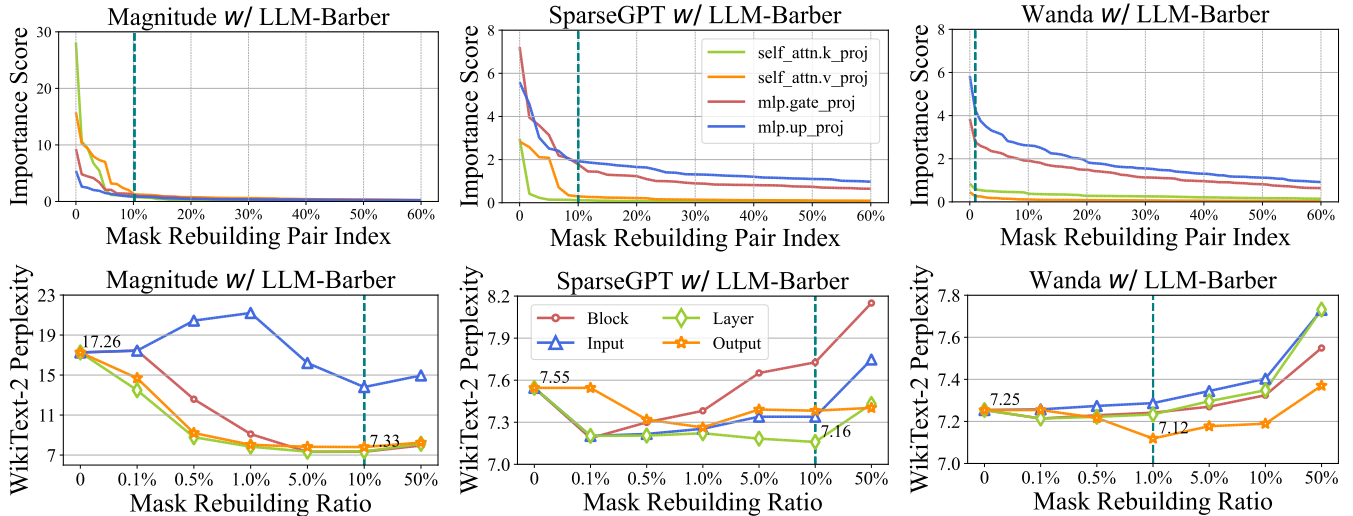


Figure 3: Relationship between importance score distribution of mask rebuilding pair and WikiText-2 perplexity results at varying pruning granularities in LLaMA-7B. The green dashed line indicates the best choice of mask rebuilding ratio.

performance, while the output-wise granularity delivers the lowest perplexity for Wanda pruning. By tailoring the pruning granularity to the particular pruning approach, LLM-Barber can consistently achieve optimal model performance.

Method	Pruning Granularity				
	w/ LLM-Barber	Block	Layer	Input	Output
Magnitude (205.5)		10.98	11.06	72.48	11.81
SparseGPT (9.994)		9.380	9.348	9.418	9.567
Wanda (9.821)		9.626	9.633	9.849	9.451

Table 5: Ablation of pruning granularity in LLaMA3-8B. Perplexity in brackets show the WikiText-2 score without LLM-Barber. **Bold** results show best granularity of this row.

Pruning Metric. We analyze the effect of different pruning metrics, focusing on weight magnitude, gradient magnitude, and the product of weight and gradient. As seen in Table 6, the product of weight and gradient consistently outperforms other metrics, achieving the lowest perplexity of 10.98 for Magnitude, 9.348 for SparseGPT, and 9.451 for Wanda. This confirms the effectiveness of our product-based pruning metric across various methods.

Method	Pruning Metric			
	w/ LLM-Barber	$ \mathbf{W} $	$ \partial\mathbf{E}/\partial\mathbf{W} $	$ \mathbf{W} \cdot \partial\mathbf{E}/\partial\mathbf{W} $
Magnitude (205.5)		186.5	14.43	10.98
SparseGPT (9.994)		9.544	9.457	9.348
Wanda (9.821)		9.880	9.699	9.451

Table 6: Ablation of pruning metric in LLaMA3-8B. **Bold** results support the pruning metric $|\mathbf{W}| \cdot |\partial\mathbf{E}/\partial\mathbf{W}|$.

Calibration Data Size. We explore how varying the size of calibration data influences the performance of LLM-Barber. Figure 4 demonstrates that as the calibration sample

size increases, LLM-Barber maintains robust performance. Notably, with more than four samples, our method achieves better perplexity than SparseGPT and Wanda. LLM-Barber outperforms SparseGPT even with just a single sample, underscoring its robustness across different sample sizes.

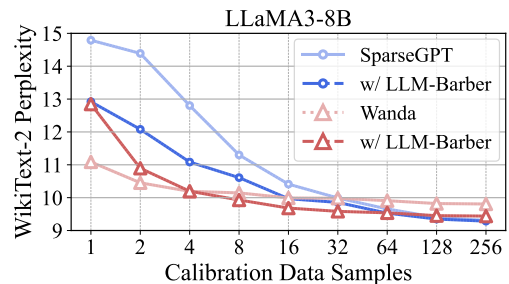


Figure 4: Ablation of calibration data size in LLaMA3-8B. LLM-Barber is robust across varying calibration data size.

Conclusion

We propose LLM-Barber (Block-Aware Rebuilder for Sparsity Mask in One-Shot), a novel framework that rebuilds sparsity mask to optimize LLM post-training pruning. By integrating a block-aware approach across Self-Attention and MLP blocks, LLM-Barber effectively reallocates sparsity to improve accuracy without the need for extensive fine-tuning. Specifically, LLM-Barber identifies novel importance score after mask initialization and rebuilds the sparsity mask with mask rebuilding pairs, simultaneously applying new sparsity masks to weights that have become less critical, thereby optimizing overall model performance. By utilizing the novel pruning metric as the product of weights and gradients, our approach enables precise and efficient reallocation of sparsity mask. Extensive experiments on pruning LLaMA series models demonstrate that LLM-Barber achieves state-of-the-art results in both perplexity and zero-shot performance within the domain of post-training pruning.

References

- AI@Meta. 2024. Llama 3 Model Card.
- Brown, T.; Mann, B.; Ryder, N.; Subbiah, M.; Kaplan, J. D.; Dhariwal, P.; Neelakantan, A.; Shyam, P.; Sastry, G.; Askell, A.; et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33: 1877–1901.
- Choquette, J.; Gandhi, W.; Giroux, O.; Stam, N.; and Krashinsky, R. 2021. Nvidia a100 tensor core gpu: Performance and innovation. *IEEE Micro*, 41(2): 29–35.
- Clark, C.; Lee, K.; Chang, M.-W.; Kwiatkowski, T.; Collins, M.; and Toutanova, K. 2019. BoolQ: Exploring the Surprising Difficulty of Natural Yes/No Questions. In *Proceedings of the 2019 Conference of the North*.
- Clark, P.; Cowhey, I.; Etzioni, O.; Khot, T.; Sabharwal, A.; Schoenick, C.; and Taffjord, O. 2018. Think you have Solved Question Answering? Try ARC, the AI2 Reasoning Challenge. *arXiv: Artificial Intelligence, arXiv: Artificial Intelligence*.
- Frantar, E.; and Alistarh, D. 2023. SparseGPT: Massive language models can be accurately pruned in one-shot. In *International Conference on Machine Learning*, 10323–10337. PMLR.
- Frantar, E.; Ashkboos, S.; Hoefler, T.; and Alistarh, D. 2023. GPTQ: Accurate Post-training Compression for Generative Pretrained Transformers. *ICLR*.
- Gao, L.; Tow, J.; Abbasi, B.; Biderman, S.; Black, S.; DiPofi, A.; Foster, C.; Golding, L.; Hsu, J.; Le Noac’h, A.; Li, H.; McDonell, K.; Muennighoff, N.; Ociepa, C.; Phang, J.; Reynolds, L.; Schoelkopf, H.; Skowron, A.; Sutawika, L.; Tang, E.; Thite, A.; Wang, B.; Wang, K.; and Zou, A. 2023. A framework for few-shot language model evaluation.
- Guan, Z.; Huang, H.; Su, Y.; Huang, H.; Wong, N.; and Yu, H. 2024. APTQ: Attention-aware Post-Training Mixed-Precision Quantization for Large Language Models. In *2024 61th ACM/IEEE Design Automation Conference (DAC)*. IEEE.
- Han, S.; Pool, J.; Tran, J.; and Dally, W. 2015. Learning both weights and connections for efficient neural network. *Advances in neural information processing systems*, 28.
- Hassibi, B.; Stork, D.; and Wolff, G. 1993. Optimal brain surgeon: Extensions and performance comparisons. *Advances in neural information processing systems*, 6.
- Hu, E. J.; Wallis, P.; Allen-Zhu, Z.; Li, Y.; Wang, S.; Wang, L.; Chen, W.; et al. 2021. LoRA: Low-Rank Adaptation of Large Language Models. In *International Conference on Learning Representations*.
- LeCun, Y.; Denker, J.; and Solla, S. 1989. Optimal brain damage. *Advances in neural information processing systems*, 2.
- Liu, Z.; Mu, H.; Zhang, X.; Guo, Z.; Yang, X.; Cheng, K.-T.; and Sun, J. 2019. Metapruning: Meta learning for automatic neural network channel pruning. In *Proceedings of the IEEE/CVF international conference on computer vision*, 3296–3305.
- Marcus, M.; Kim, G.; Marcinkiewicz, M. A.; MacIntyre, R.; Bies, A.; Ferguson, M.; Katz, K.; and Schasberger, B. 1994. The Penn Treebank. In *Proceedings of the workshop on Human Language Technology - HLT '94*.
- Merity, S.; Xiong, C.; Bradbury, J.; and Socher, R. 2017. Pointer Sentinel Mixture Models. *International Conference on Learning Representations*.
- Mihaylov, T.; Clark, P.; Khot, T.; and Sabharwal, A. 2018. Can a Suit of Armor Conduct Electricity? A New Dataset for Open Book Question Answering. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*.
- Mocanu, D. C.; Mocanu, E.; Stone, P.; Nguyen, P. H.; Gibescu, M.; and Liotta, A. 2018. Scalable training of artificial neural networks with adaptive sparse connectivity inspired by network science. *Nature communications*, 9(1): 2383.
- Raffel, C.; Shazeer, N.; Roberts, A.; Lee, K.; Narang, S.; Matena, M.; Zhou, Y.; Li, W.; and Liu, P. J. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140): 1–67.
- Sun, M.; Liu, Z.; Bair, A.; and Kolter, J. Z. 2023. A Simple and Effective Pruning Approach for Large Language Models. In *The Twelfth International Conference on Learning Representations*.
- Touvron, H.; Lavril, T.; Izacard, G.; Martinet, X.; Lachaux, M.-A.; Lacroix, T.; Rozière, B.; Goyal, N.; Hambro, E.; Azhar, F.; et al. 2023a. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Touvron, H.; Martin, L.; Stone, K.; Albert, P.; Almahairi, A.; Babaei, Y.; Bashlykov, N.; Batra, S.; Bhargava, P.; Bhosale, S.; et al. 2023b. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Wang, A.; Singh, A.; Michael, J.; Hill, F.; Levy, O.; and Bowman, S. 2018. GLUE: A Multi-Task Benchmark and Analysis Platform for Natural Language Understanding. In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*.
- Xu, P.; Shao, W.; Chen, M.; Tang, S.; Zhang, K.; Gao, P.; An, F.; Qiao, Y.; and Luo, P. 2024. BESA: Pruning Large Language Models with Blockwise Parameter-Efficient Sparsity Allocation. In *The Twelfth International Conference on Learning Representations*.
- Zellers, R.; Holtzman, A.; Bisk, Y.; Farhadi, A.; and Choi, Y. 2019. HellaSwag: Can a Machine Really Finish Your Sentence? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*.
- Zhang, S.; Roller, S.; Goyal, N.; Artetxe, M.; Chen, M.; Chen, S.; Dewan, C.; Diab, M.; Li, X.; Lin, X. V.; et al. 2022. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*.
- Zhang, Y.; Zhao, L.; Lin, M.; Yunyun, S.; Yao, Y.; Han, X.; Tanner, J.; Liu, S.; and Ji, R. 2024. Dynamic Sparse No Training: Training-Free Fine-tuning for Sparse LLMs. In *The Twelfth International Conference on Learning Representations*.