

# WaLeM: Layer Pruning via Weight-Aware Learnable Merging for Large Language Model Compression

Anonymous ACL submission

## Abstract

Layer pruning has emerged as a promising technique for compressing Large Language Models (LLMs) by reducing their depth. However, existing methods predominantly rely on direct layer removal or coarse-grained layer merging, where uniform merging coefficients are applied to entire layers. These approaches neglect the distinct functional roles of internal components, often leading to severe feature blurring and performance degradation. To address this, we propose Weight-aware Learnable Merging (WaLeM), a novel framework that transitions from coarse-grained, heuristic layer pruning to fine-grained, optimization-driven layer merging. In WaLeM, we first employ Centered Kernel Alignment (CKA) combined with dynamic programming to globally identify redundant layers for merging, ensuring structural consistency. Subsequently, we introduce a learnable merging mechanism that assigns adaptive, component-specific coefficients to Transformer’s weight matrices, optimized via knowledge distillation. Extensive experiments on various models demonstrate that WaLeM significantly outperforms state-of-the-art baselines. Notably, WaLeM preserves complex reasoning capabilities on benchmarks, like GSM8K, even at high compression rates, offering a superior trade-off between efficiency and performance.

## 1 Introduction

The rapid evolution of Large Language Models (LLMs), represented by the GPT (Brown et al., 2020; Achiam et al., 2023) and Llama (Touvron et al., 2023; Grattafiori et al., 2024) series, has fundamentally transformed natural language processing. When scaled to billions of parameters, these models exhibit remarkable performance across a wide range of tasks, including complex reasoning (Wei et al., 2023; Yao et al., 2023), code generation (Chen et al., 2021; Jimenez et al., 2024), and mathematical problem solving (Cobbe et al., 2021;

Lightman et al., 2023). However, the prohibitive computational costs and memory footprints associated with such massive scales pose significant impediments to their deployment, particularly in resource constrained environments such as edge computing platforms.

To address these challenges, substantial research efforts have been directed towards model compression (Cheng et al., 2020; Zhu et al., 2024; Wang et al., 2024). Generally, compression methodologies encompass quantization (Liu et al., 2023; Dettmers et al., 2023), knowledge distillation (Hinton et al., 2015; Gu et al., 2024), and pruning (Ma et al., 2023; Zhu et al., 2024). Pruning can be further delineated into unstructured pruning, which sparsify individual weights (Frantar and Alistarh, 2023), and structured pruning that eliminate coherent components such as attention heads or entire layers (Ma et al., 2023). Layer pruning, a specific form of structured pruning, has gained prominence for its ability to reduce model depth directly. Unlike other pruning approaches, it delivers tangible reduction in inference latency while preserving architectural compatibility without requiring specialized hardware support (Ling et al., 2024). Capitalizing on these advantages, this study focus on the layer-wise pruning methodology.

Current research on layer pruning primarily addresses two fundamental challenges: selecting specific redundancy layers and determining the pruning execution strategy (Pei et al., 2025). For redundancy layer selection, various studies have introduced distinct metrics, including block influence, perplexity and angular distance (Men et al., 2024; Song et al., 2024; Gromov et al., 2025). Regarding determining the pruning execution strategy, existing methods generally fall into two categories: direct removal (Men et al., 2024; Song et al., 2024) and layer merging (Choudhry et al., 2025; Liu et al., 2025). Compared with layer merging, direct removal usually risks losing valuable

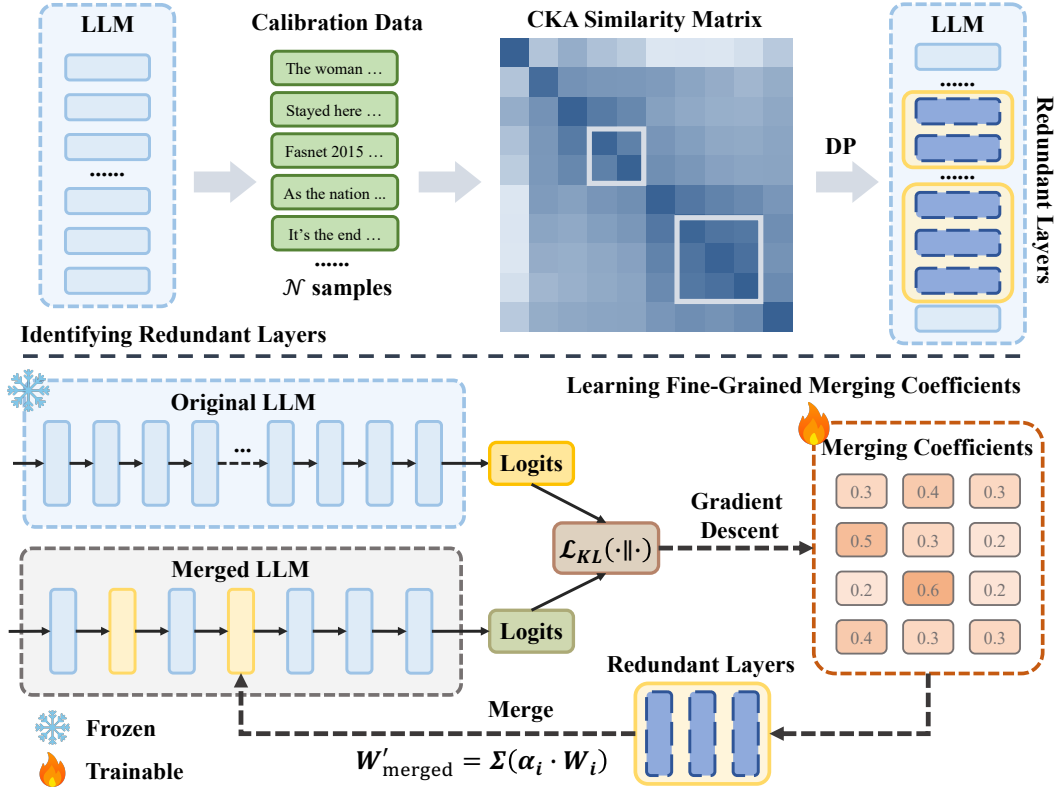


Figure 1: Overview of the Weight-aware Learnable Merging (WaLeM) framework. The framework involves two phases. (1) Identifying redundant layers, which employs Centered Kernel Alignment (CKA) to quantify representational similarity across all layers and applies dynamic programming (DP) to locate blocks with redundant layers; (2) Learning fine-grained merging coefficients, which treats the merging coefficients as learnable parameters and train them through a lightweight knowledge-distillation process.

083 model structures and interrupts the flow of semantic  
084 information (Liu et al., 2025). Consequently,  
085 we adopt layer merging strategy to condense multiple  
086 layers and preserve information. However, current  
087 merging-based approaches still exhibit a fundamental  
088 drawback: the merging coefficients are determined in  
089 an empirical or heuristic manner, which restricts  
090 these approaches to adopting a single, uniform  
091 coefficient at each layer-merging step (Yang et al.,  
092 2024; Liu et al., 2025). Such coarse-grained  
093 merging methods fail to account for the heterogeneous  
094 functional responsibilities of sub-modules inside  
095 each Transformer layer (Clark et al., 2019; Geva  
096 et al., 2021; Meng et al., 2023), inevitably  
097 causing feature blurring and functional confusion.  
098

099 Therefore, we present **Weight-aware Learnable**  
100 **Merging (WaLeM)**, a framework that transition  
101 model compression from coarse-grained, heuristic  
102 layer pruning to fine-grained, optimization-driven  
103 and information-preserving layer merging. In detail,  
104 WaLeM begins by using Centered Kernel Alignment  
105 (Kornblith et al., 2019) to quantify representational  
106 similarity across all layers. Then, a

dynamic programming algorithm is employed to  
107 analyze this similarity and identify optimal merging  
108 blocks, which can bypass the sub-optimality of  
109 previous greedy methods (Liu et al., 2025; Pei et al.,  
110 2025). Finally, based on the identified optimal  
111 merge blocks, we implement fine-grained learnable  
112 merging by assigning independent coefficients to  
113 different Transformer’s weight matrices. These  
114 coefficients are treated as learnable parameters and  
115 are determined via a lightweight knowledge  
116 distillation phase that minimizes Kullback–Leibler  
117 divergence against the original LLM to achieve  
118 alignment. Such a way allows the merging process  
119 to adaptively inherit information based on  
120 component sensitivity rather than relying on a  
121 fixed coefficient. Extensive experiments on five  
122 LLMs across diverse tasks demonstrate that WaLeM  
123 outperforms established baselines, maintaining  
124 robust capabilities even at high compression rates.  
125

In summary, the main contributions of our work  
126 are as follows:  
127

- We propose **Weight-aware Learnable Merging**  
128 **(WaLeM)**, a novel framework designed  
129 to compress LLMs through information-  
130

preserving layer merging. WaLeM employs a fine-grained, learnable mechanism that assigns adaptive coefficients to heterogeneous weight components, which allows for precise information inheritance.

- In WaLeM, we develop a globally optimal strategy for layer selection based on Centered Kernel Alignment (CKA) and dynamic programming. By rigorously quantifying the representational similarity across the network, our approach can identify redundant layers from a global perspective.
- We validate the efficacy of WaLeM through comprehensive experiments on multiple benchmarks and various large language models. The results demonstrate that our method consistently outperforms existing baselines by effectively compressing models while preserving their original performance.

## 2 Related Work

**Metrics for Layer Redundancy.** Identifying redundant layers constitutes a critical step in layer pruning. To address this, numerous studies have established metrics to quantify the importance of specific layers. For instance, ShortGPT (Men et al., 2024) introduces Block Influence (BI), which assesses the similarity between a layer’s input and output. Similarly, MKA (Liu et al., 2025) employs manifold learning and the Normalized Pairwise Information Bottleneck to identify layers suitable for merging. Other strategies, such as SLEB (Song et al., 2024), leverage perplexity to evaluate layer importance. Furthermore, angular distance, Magnitude and first-order Taylor saliency has been used to detect less informative layers (Gromov et al., 2025; Kim et al., 2024). Collectively, these metrics guide the pruning process to ensure the preservation of functional integrity.

**Pruning via Layer Merging.** To mitigate the information loss inherent in direct layer removal, recent studies have pivoted towards layer merging, a paradigm that consolidates adjacent layer parameters rather than discarding them. Prominent approaches, such as MKA (Liu et al., 2025), LaCo (Yang et al., 2024), and Layer-Merge (Choudhry et al., 2025), have advanced this direction by employing manifold learning, iterative parameter aggregation, and PCA-based activation analysis, respectively. Despite these innovations, existing techniques predominantly rely on coarse-

grained merging strategies, which neglect the functional heterogeneity of components within Transformer layers. Concurrent with and independent of our work, FuseGPT (Pei et al., 2025) utilizing learnable low-rank matrices to graft pruned layers onto adjacent ones. Our approach differs in two key aspects: (1) Granularity: Instead of low-rank grafting, WaLeM performs weight-aware merging, assigning adaptive coefficients to specific functional components; (2) Efficiency: While FuseGPT relies on an iterative greedy search to handle rank shift, WaLeM utilizes Dynamic Programming to efficiently derive a globally optimal merging strategy in a single pass.

## 3 Methodology

The proposed WaLeM framework mainly proceeds in two phases: identifying redundant layers and learning the merging coefficients, which is illustrated in Figure 1.

### 3.1 Identifying Redundant Layers

We consider a LLM comprising  $L$  Transformer layers. Given a dataset  $\mathcal{D} = \{x_i\}_{i=1}^N$ , we extract the output of the  $l$ -th Transformer layer as the feature representation. Flattening that feature representation across the batch and sequence dimensions yields  $\mathbf{X}^{(l)}$ :

$$\mathbf{X}^{(l)} \in \mathbb{R}^{M \times d} \quad (1)$$

where  $M$  represents the token count and  $d$  is the hidden dimension.

To quantify geometric similarity between the  $l$ -th layer and the  $k$ -th layer, we employ Linear Centered Kernel Alignment (CKA). In detail, we first compute the Gram matrices:

$$\mathbf{K}^{(l)} = \mathbf{X}^{(l)}(\mathbf{X}^{(l)})^\top, \quad \mathbf{K}^{(k)} = \mathbf{X}^{(k)}(\mathbf{X}^{(k)})^\top \quad (2)$$

Subsequently, we calculate the Hilbert-Schmidt Independence Criterion (HSIC) between the two Gram matrices:

$$\text{HSIC}(\mathbf{K}^{(l)}, \mathbf{K}^{(k)}) = \frac{1}{(M-1)^2} \text{tr}(\mathbf{K}^{(l)} \mathbf{H} \mathbf{K}^{(k)} \mathbf{H}) \quad (3)$$

where  $\mathbf{H} = \mathbf{I}_M - \frac{1}{M} \mathbf{1}\mathbf{1}^\top$  is the centering matrix used to remove the mean shift. The CKA similarity score  $s_{lk}$  between the  $l$ -th layer and  $k$ -th layer is defined as the normalized HSIC value:

$$\begin{aligned} s_{lk} &= \text{CKA}(\mathbf{X}^{(l)}, \mathbf{X}^{(k)}) \\ &= \frac{\text{HSIC}(\mathbf{K}^{(l)}, \mathbf{K}^{(k)})}{\sqrt{\text{HSIC}(\mathbf{K}^{(l)}, \mathbf{K}^{(l)}) \cdot \text{HSIC}(\mathbf{K}^{(k)}, \mathbf{K}^{(k)})}} \end{aligned} \quad (4)$$

Computing this for all layer pairs yields a symmetric CKA similarity matrix  $\mathbf{S} \in \mathbb{R}^{L \times L}$ , where  $\mathbf{S}[l, k] = s_{lk}$ , capturing the model’s hierarchical evolution.

Building on the CKA similarity matrix  $\mathbf{S}$ , we further formulate identifying redundant layers as a constrained optimization problem on  $\mathbf{S}$  to remove exactly  $K$  layers by merging non-overlapping diagonal blocks.

In detail, for a candidate block of size  $B$  ending at the  $k$ -th layer (spanning  $[k - B + 1, k]$ ), we define the redundancy score  $r(k, B)$  as:

$$r(k, B) = B^\alpha \cdot (\bar{s}_{k,B} - \gamma) \cdot \mathcal{W}_{pos}(\bar{l}) \quad (5)$$

where  $\bar{s}_{k,B}$  is the average similarity within the upper triangle of the block’s sub-matrix,  $\alpha$  is a hyperparameter incentivizing larger block sizes, and  $\gamma$  represents the predefined similarity threshold. To prioritize deeper layers (Gromov et al., 2025), we apply a position-dependent weight based on the block’s center  $\bar{l} = k - \frac{B-1}{2}$ :

$$\mathcal{W}_{pos}(\bar{l}) = 1 + \beta \cdot \frac{\bar{l}}{L} \quad (6)$$

where  $L$  is the total number of layers. This weighting mechanism ensures that blocks located in the latter stages of the network contribute more significantly to the overall score.

Finally, a dynamic programming algorithm is employed to derive the optimal strategy. Let  $v_{k,j}$  denote the maximum redundancy score achievable within the first  $k$  layers, subject to a constraint that exactly  $j$  layers are removed. We define the state transition equation as:

$$v_{k,j} = \max \left( v_{k-1,j}, \max_{B \in [B_{\min}, B_{\max}]} \{v_{k-B,j-(B-1)} + r(k, B)\} \right) \quad (7)$$

The first term corresponds to preserving the  $k$ -th layer (i.e., not merging it into a preceding block), where the state simply inherits the score from  $v_{k-1,j}$ . The second term considers the case where the  $k$ -th layer serves as the end of a merging block of size  $B$ . Since merging  $B$  layers into a single layer results in a reduction of  $B - 1$  layers, we transition from the state  $v_{k-B,j-(B-1)}$  and add the corresponding redundancy gain  $r(k, B)$ .

By leveraging a classical dynamic programming algorithm, the problem can be solved exactly and efficiently with an  $\mathcal{O}(L^2)$  time complexity. Moreover, such a way enables the identification of redundant layers from a global, rather than local or heuristic, perspective.

### 3.2 Learning Fine-Grained Merging Coefficients

Once the redundant layers are identified, we determine merging coefficients via a learning phase. Consider a target redundant block of size  $B$  involving Transformer layers with indices  $l \in [k - B + 1, k]$ . A set of weight types  $\mathcal{T}$  (e.g.,  $\mathbf{W}_q, \mathbf{W}_k, \mathbf{W}_v, \mathbf{W}_o$  in attention) are defined to be merged.

For each block, we initialize a learnable coefficient matrix  $\mathbf{F} \in \mathbb{R}^{|\mathcal{T}| \times B}$ . To ensure normalized contributions, we apply a row-wise Softmax operation over the layer dimension. The merging coefficient  $\alpha_{\tau,l}$  for a specific weight type  $\tau \in \mathcal{T}$  and the  $l$ -th layer in the block is computed as:

$$\alpha_{\tau,l} = \frac{\exp(\mathbf{F}_{\tau,l})}{\sum_{j=1}^B \exp(\mathbf{F}_{\tau,j})} \quad (8)$$

The merged Transformer weight matrix  $\widetilde{\mathbf{W}}_\tau$  is then obtained by the weighted aggregation of the original parameters  $\mathbf{W}_\tau^{(l)}$ :

$$\widetilde{\mathbf{W}}_\tau = \sum_{j=1}^B \alpha_{\tau,j} \cdot \mathbf{W}_\tau^{(k-B+j)} \quad (9)$$

To preserve the representational capacity of the original LLM  $\mathcal{M}$ , we align the merged LLM  $\mathcal{M}'$  using a calibration dataset  $\mathcal{D}_{cal}$  (Kovalev and Tikhomirov, 2025). We minimize the Kullback-Leibler (KL) divergence between the output probability distributions of the two models. The optimization objective is formulated as:

$$\mathcal{L} = \sum_{x \in \mathcal{D}_{cal}} D_{KL}(P_{\mathcal{M}}(y|x) \| P_{\mathcal{M}'}(y|x)) \quad (10)$$

where  $P_{\mathcal{M}}$  and  $P_{\mathcal{M}'}$  denote the softmax probabilities generated by the original model and the merged model, respectively. We update  $\mathbf{F}$  via backpropagation while keeping other model parameters frozen.

Finally, consistent with prior studies (Ma et al., 2023), we implement a fast post-fusion recovery phase to mitigate performance degradation. We employ Low-Rank Adaptation (LoRA) (Hu et al., 2021) to efficiently heal the merged model. This lightweight fine-tuning step allows the compressed model to recover its capabilities with minimal computational overhead.

Note that, the comprehensive pseudocode of the WaLeM framework is shown in the Algorithm 1 at the Appendix.

Model	Method	Pruning Rate	PPL (↓)	HellaSwag	PIQA	Winogrande	MMLU	ARC-e	ARC-c	RACE	GSM8K	Avg. Score	
Llama3-8B	Dense	0.00%	7.44	79.11	80.96	73.40	62.12	77.65	54.10	40.38	49.89	64.70	
	MKA	18.75%	1469.73	57.39	69.97	59.75	58.45	61.91	44.03	31.96	2.50	48.25	
	ShortGPT	18.75%	39.72	68.31	72.80	<u>70.17</u>	60.55	59.47	43.94	34.74	5.61	51.95	
	SLEB	18.75%	14.99	63.66	74.86	56.27	25.85	62.25	35.07	32.73	2.58	44.16	
	LaCo	18.75%	3167.02	30.45	57.89	59.19	35.09	36.36	28.58	24.11	0.30	34.00	
	LLMPruner	18.75%	29.50	69.71	76.17	67.09	<b>62.04</b>	70.16	<b>48.55</b>	38.76	23.96	57.06	
	Ours	18.75%	<b>13.63</b>	<b>74.52</b>	<b>76.55</b>	<b>72.23</b>	<u>61.28</u>	<b>70.96</b>	<u>47.95</u>	<b>41.82</b>	<b>36.32</b>	<b>60.20</b>	
	MKA	25.00%	3557.07	52.94	67.74	60.38	<b>62.20</b>	59.00	41.72	29.95	0.45	46.80	
	ShortGPT	25.00%	-	31.41	61.15	55.01	34.59	39.65	33.11	24.69	0.23	34.98	
	SLEB	25.00%	23.07	57.96	72.14	52.41	26.63	57.87	32.85	31.29	1.90	41.63	
	LaCo	25.00%	89.61	60.01	67.90	65.04	60.58	52.27	38.48	32.25	2.12	47.33	
	LLMPruner	25.00%	42.06	67.34	73.50	65.82	<u>61.04</u>	<b>65.99</b>	<b>45.56</b>	37.99	2.20	52.43	
	Ours	25.00%	<b>17.19</b>	<b>70.63</b>	<b>73.67</b>	<b>70.72</b>	59.71	<u>64.81</u>	<u>44.20</u>	<b>40.67</b>	<b>22.82</b>	<b>55.90</b>	
	MKA	31.25%	9999.44	48.19	63.87	60.85	<b>59.03</b>	52.90	39.76	29.86	0.00	44.31	
	ShortGPT	31.25%	-	33.83	60.45	58.41	36.68	38.24	32.08	26.79	0.99	35.93	
	SLEB	31.25%	<u>34.24</u>	52.88	69.26	53.51	26.90	52.36	31.31	29.67	2.35	39.78	
	LaCo	31.25%	166.27	48.36	63.71	63.61	25.38	41.46	31.57	29.38	1.29	39.03	
	LLMPruner	31.25%	59.94	63.31	71.38	64.40	54.27	<b>64.14</b>	<b>44.71</b>	36.36	<b>14.94</b>	51.69	
	Ours	31.25%	<b>21.84</b>	<b>66.15</b>	<b>72.52</b>	<b>70.09</b>	<u>58.78</u>	<u>61.49</u>	<u>42.75</u>	<b>38.85</b>	<u>5.53</u>	<b>52.02</b>	
	Llama2-13B	Dense	0.00%	7.74	79.65	80.30	72.53	52.39	76.47	48.81	40.48	23.65	59.29
		MKA	20.00%	1411.24	55.44	69.64	62.83	50.69	58.33	40.53	32.15	0.00	46.20
		ShortGPT	20.00%	15.35	<u>72.62</u>	75.63	69.93	52.22	67.00	43.26	<b>38.09</b>	2.35	52.64
		SLEB	20.00%	<u>11.77</u>	70.69	76.33	63.69	30.02	70.62	43.00	37.70	2.81	49.36
		LaCo	20.00%	66.70	61.98	69.53	65.35	37.96	57.74	37.80	33.68	1.67	45.71
LLMPruner		20.00%	38.53	72.09	76.82	66.30	<u>52.59</u>	71.09	<u>46.42</u>	37.42	<u>11.45</u>	54.27	
Ours		20.00%	<b>11.63</b>	<b>76.33</b>	<b>78.13</b>	<b>71.35</b>	<b>53.57</b>	<b>71.84</b>	<b>48.55</b>	<b>40.57</b>	<b>13.50</b>	<b>56.73</b>	
MKA		25.00%	4433.19	50.91	66.59	61.17	50.06	54.29	37.88	31.67	0.00	44.07	
ShortGPT		25.00%	40.40	68.08	72.74	69.14	50.10	60.10	41.89	<u>37.32</u>	0.76	50.02	
SLEB		25.00%	<u>13.90</u>	67.23	<b>76.66</b>	60.38	25.31	65.45	39.33	35.31	1.44	46.39	
LaCo		25.00%	130.90	55.80	66.59	64.25	34.45	54.17	35.75	30.53	0.30	42.73	
LLMPruner		25.00%	67.54	67.89	75.14	64.88	<b>52.27</b>	<u>68.56</u>	<u>44.03</u>	36.75	<b>6.82</b>	<u>52.04</u>	
Ours		25.00%	<b>13.44</b>	<b>74.34</b>	<u>76.39</u>	<b>71.82</b>	<u>51.69</u>	<b>69.61</b>	<b>45.56</b>	<b>40.10</b>	<u>6.67</u>	<b>54.52</b>	
MKA		30.00%	8453.01	47.57	64.58	59.75	<u>51.44</u>	52.23	36.09	31.29	0.00	42.87	
ShortGPT		30.00%	65.10	58.31	69.86	68.11	48.44	53.45	36.77	31.96	0.99	45.99	
SLEB		30.00%	<u>16.54</u>	64.49	<u>74.59</u>	<u>61.09</u>	23.37	63.09	36.60	34.45	1.67	44.92	
LaCo		30.00%	205.28	52.07	63.98	63.14	36.87	51.30	35.15	29.57	0.53	41.58	
LLMPruner		30.00%	85.14	65.53	72.85	64.01	51.05	65.74	<b>43.34</b>	<u>36.27</u>	3.79	<u>50.32</u>	
Ours		30.00%	<b>15.11</b>	<b>71.60</b>	<b>75.24</b>	<b>70.01</b>	<b>52.76</b>	<b>66.50</b>	<u>43.26</u>	<b>38.85</b>	<b>4.40</b>	<b>52.83</b>	

Table 1: Zero-shot performance of Llama models with distinct pruning strategies. “Dense” denotes the unpruned models and “PPL” denotes the perplexity evaluated on Wikitext2. **Bold** fonts indicate the best results, while underlined values represent the sub-optimal performance. Results for Llama2-7B are detailed in Appendix A.3.

## 4 Experiments

### 4.1 Experimental Setups

**LLMs.** We perform experiments on five different-sized LLMs: Llama2-7B (Touvron et al., 2023), Llama2-13B (Touvron et al., 2023), Llama3-8B (Grattafiori et al., 2024), Qwen2.5-7B (Qwen et al., 2025) and Qwen2.5-14B (Qwen et al., 2025).

**Datasets and Metrics.** We utilize the C4 (Rafael et al., 2023) for identifying redundant layers (500 samples) and learning merging coefficients (10k samples), and utilize Alpaca (Taori et al., 2023) for performance recovery. Using LM Evaluation Harness (Gao et al., 2024), we evaluate on **Perplexity** (WikiText2 (Merity et al., 2016)), **Commonsense Reasoning** (HellaSwag (Zellers et al., 2019), PIQA (Bisk et al., 2019), WinoGrande (Sakaguchi et al., 2019), MMLU (Hendrycks et al., 2021), ARC-easy (Clark et al., 2018), ARC-challenge (Clark et al., 2018),

RACE (Lai et al., 2017)), and **Mathematical Reasoning** (GSM8K (Cobbe et al., 2021)).

**Baselines.** We benchmark WaLeM against five state-of-the-art methods. For layer removal, we compare it with LLM-Pruner (Ma et al., 2023), ShortGPT (Men et al., 2024), and SLEB (Song et al., 2024), which utilize depth-based, block influence, and greedy strategies, respectively. For layer merging, we compare it with LaCo (Yang et al., 2024) and MKA (Liu et al., 2025), which employ representation thresholds and manifold learning to guide merging. Specific implementation details are available in Appendix A.1.

### 4.2 Main Results

Following prior studies (Zhang et al., 2024; Zhong et al., 2025; Lu et al., 2025), we evaluate our proposed WaLeM framework across pruning ratios ranging from approximately 20% to 30%. As presented in Table 1 and Table 2, WaLeM demon-

Model	Method	Pruning Rate	PPL ( $\downarrow$ )	HellaSwag	PIQA	Winogrande	MMLU	ARC-e	ARC-c	RACE	GSM8K	Avg. Score
Qwen2.5-7B	Dense	0.00%	9.41	80.58	80.47	71.19	71.78	81.40	55.03	46.22	81.80	71.06
	MKA	21.43%	182964.77	45.12	61.81	<u>58.09</u>	<b>40.31</b>	40.40	35.49	28.61	1.14	38.87
	ShortGPT	21.43%	<u>16.64</u>	<b>65.26</b>	<b>76.93</b>	56.04	28.12	<u>68.27</u>	<u>40.44</u>	<b>37.42</b>	<u>4.25</u>	<u>47.09</u>
	LaCo	21.43%	797.66	53.05	61.97	53.59	26.48	38.22	30.97	25.26	1.21	36.34
	Ours	21.43%	<b>13.22</b>	<b>66.84</b>	<u>76.39</u>	<b>60.22</b>	<u>39.78</u>	<b>69.49</b>	<b>45.05</b>	<u>37.32</u>	<b>15.69</b>	<b>51.35</b>
	MKA	25.00%	607147.37	42.11	60.72	<b>57.38</b>	<b>37.10</b>	37.92	33.53	26.99	0.99	37.09
	ShortGPT	25.00%	<u>17.96</u>	<u>61.02</u>	<u>74.59</u>	53.99	27.86	<u>63.72</u>	<u>36.43</u>	<u>35.12</u>	<u>2.96</u>	<u>44.46</u>
	LaCo	25.00%	1357.73	48.24	60.77	50.04	23.07	35.14	30.29	22.58	1.21	33.92
	Ours	25.00%	<b>13.71</b>	<b>65.06</b>	<b>76.12</b>	<u>57.06</u>	<u>30.63</u>	<b>68.27</b>	<b>42.66</b>	<b>37.22</b>	<b>8.95</b>	<b>48.25</b>
	MKA	28.57%	1876698.27	38.60	59.14	<u>56.35</u>	23.00	34.01	31.91	25.36	1.59	33.75
	ShortGPT	28.57%	<u>20.44</u>	56.95	<u>73.45</u>	54.93	<b>26.58</b>	61.49	<u>36.09</u>	<u>34.64</u>	<u>2.35</u>	<u>43.31</u>
	LaCo	28.57%	1694.25	44.19	58.92	51.22	23.44	34.85	30.29	23.25	0.61	33.35
Ours	28.57%	<b>17.95</b>	<b>62.12</b>	<b>74.97</b>	<b>56.43</b>	<u>26.14</u>	<b>65.87</b>	<b>41.38</b>	<b>35.22</b>	<b>4.17</b>	<b>45.79</b>	
Qwen2.5-14B	Dense	0.00%	6.66	84.37	81.94	75.77	78.83	81.65	62.29	45.36	44.43	69.33
	MKA	20.83%	400283.16	45.90	65.02	58.96	<b>77.80</b>	58.33	41.38	31.39	0.00	47.35
	ShortGPT	20.83%	-	71.47	<u>77.91</u>	60.30	42.60	66.46	43.34	35.41	3.03	50.07
	LaCo	20.83%	1909.78	47.30	61.81	<u>62.04</u>	<u>53.79</u>	42.38	33.79	29.19	0.23	41.32
	Ours	20.83%	<b>11.95</b>	<b>72.31</b>	<b>78.18</b>	<b>63.22</b>	<u>44.94</u>	<b>72.90</b>	<b>49.23</b>	<b>38.28</b>	<b>26.16</b>	<b>55.65</b>
	MKA	25.00%	1289431.81	42.51	63.17	<u>59.59</u>	<b>53.32</b>	52.15	38.91	30.53	0.00	42.52
	ShortGPT	25.00%	-	67.00	75.14	56.20	32.97	62.79	37.20	34.55	3.11	46.12
	LaCo	25.00%	12723.14	44.42	62.35	<b>59.91</b>	42.08	39.65	35.58	30.53	0.15	39.33
	Ours	25.00%	<b>13.64</b>	<b>69.60</b>	<b>77.37</b>	<u>58.56</u>	39.96	<b>69.87</b>	<b>45.05</b>	<b>36.08</b>	<b>9.33</b>	<b>50.73</b>
	MKA	31.25%	7517022.80	36.18	57.94	57.85	<b>32.08</b>	39.52	34.30	26.32	0.00	35.52
	ShortGPT	31.25%	-	46.86	<u>68.82</u>	51.38	23.64	<u>52.78</u>	30.55	24.11	0.53	37.33
	LaCo	31.25%	<u>775.68</u>	40.32	<u>62.84</u>	56.27	22.99	35.52	30.72	25.26	0.68	34.33
Ours	31.25%	<b>18.76</b>	<b>63.76</b>	<b>74.27</b>	<b>59.35</b>	<u>31.84</u>	<b>64.31</b>	<b>40.78</b>	<b>37.03</b>	<b>4.78</b>	<b>47.02</b>	

Table 2: Zero-shot performance of Qwen models with distinct pruning strategies. “Dense” denotes the unpruned models and “PPL” denotes the perplexity evaluated on Wikitext2. **Bold** fonts indicate the best results, while underlined values represent the sub-optimal performance.

(a) Results on Llama3-8B					
Strategy	Pruning Rate				
	18.75%	21.88%	25.00%	28.13%	31.25%
<b>BI</b>	<b>21.16</b>	76.44	100.54	190.62	363.86
<b>DP</b>	44.24	<b>59.33</b>	<b>89.27</b>	<b>142.42</b>	<b>263.25</b>

(b) Results on Llama2-13B					
Strategy	Pruning Rate				
	15.00%	20.00%	25.00%	30.00%	35.00%
<b>BI</b>	11.54	13.86	28.21	<b>41.22</b>	60.39
<b>DP</b>	<b>10.90</b>	<b>13.53</b>	<b>18.05</b>	49.38	<b>59.22</b>

Table 3: Perplexity comparison between the proposed Global Optimal Fusion strategy and ShortGPT Block Influence across various pruning ratios.

strates superior performance consistently over baseline methods across diverse zero-shot downstream tasks. Specifically, WaLeM achieves the highest average scores while maintaining exceptionally low perplexity on the Wikitext2 dataset. Notably, this advantage becomes increasingly pronounced at higher pruning ratios. While baseline methods often exhibit severe performance degradation or exploding perplexity, WaLeM preserves the model’s linguistic coherence and knowledge retention.

Beyond standard multiple-choice benchmarks, we explicitly evaluate mathematical reasoning capabilities using the GSM8K dataset, which is often

overlooked in previous pruning studies. WaLeM exhibits remarkable resilience in this challenging tasks. For instance, even with approximately 20% of parameter pruned, our method still retains significant reasoning abilities, particularly on Llama3-8B, achieving results that other baselines fail to match.

Furthermore, we visualize the CKA similarity matrices for Llama3-8B, Llama2-13B, and Qwen2.5-14B in Figure 2. These visualizations explicitly reveal the inherent redundancy present within the original LLMs. By annotating the specific merging strategies with red blocks on the similarity matrices, we demonstrate that our proposed WaLeM method can efficiently capture this redundancy. Subsequently, we present the CKA matrices of the pruned models using these strategies in Figure 3. Comparing this with Figure 2, we observe that the original redundancy in the deeper layers of Llama2-13B, as well as in the middle and deeper layers of Qwen2.5-14B, has been effectively mitigated. This suggests a more efficient information flow within the compressed model.

## 5 Analyses

### 5.1 Effectiveness of Identifying Redundant Layers

We investigate the effectiveness of the strategy of identifying redundant layers armed by our dy-

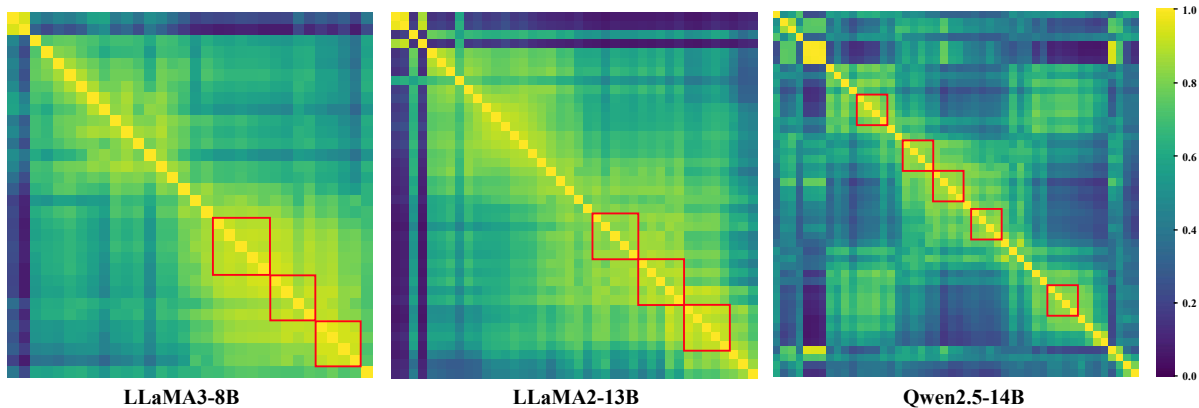


Figure 2: CKA similarity matrices for Llama3-8B, Llama2-13B, and Qwen2.5-14B. The axes represent layer indices arranged from shallow to deep. Red highlighted regions denote the layer fusion strategies adopted at an approximate pruning ratio of 30%.

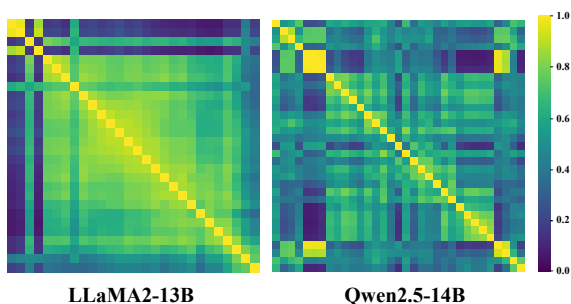


Figure 3: Post-pruning CKA similarity matrices for Llama2-13B and Qwen2.5-14B with an approximate 30% pruning ratio.

392 namic programming approach. We benchmark our  
 393 method against the Block Influence metric pro-  
 394 posed in ShortGPT (Men et al., 2024). Both strate-  
 395 gies serve as the initialization for the subsequent  
 396 fine-grained learnable merging training. We evalu-  
 397 ate the pruned models based on perplexity using the  
 398 Wikitext2 dataset. Results are presented in Table 3.  
 399 As shown, our proposed strategy consistently out-  
 400 performs Block Influence on Llama3-8B, achiev-  
 401 ing significant margins in most of the pruning ratios.  
 402 Similar performance advantages are also observed  
 403 on the Llama2-13B model.

## 404 5.2 Significance of Fine-Grained Merging

405 To validate the efficacy of our proposed strategy, we  
 406 conduct a controlled experiment where merging is  
 407 restricted to the final layers to ensure a fair compar-  
 408 ison. In this setting, we benchmark our learnable  
 409 merging coefficients against two baselines, MKA  
 410 and LaCo. Specifically, MKA directly derives the  
 411 merging coefficients via similarity, whereas LaCo  
 412 integrates layer differences. As shown in Table 4,  
 413 our method achieves significantly lower perplexity

Model	Method	Ratio	PPL ( $\downarrow$ )	Avg. Score
Llama3-8B	LaCo	18.75%	494.73	37.00
	MKA		1469.73	48.19
	<b>Ours</b>		<b>75.07</b>	<b>48.63</b>
	LaCo	25.00%	393.57	38.30
	MKA		3557.07	<b>46.94</b>
	<b>Ours</b>		<b>126.64</b>	<b>46.69</b>
Llama2-7B	LaCo	18.75%	100.70	43.75
	MKA		1012.85	42.30
	<b>Ours</b>		<b>72.35</b>	<b>44.58</b>
	LaCo	25.00%	145.84	40.88
	MKA		1837.25	39.67
	<b>Ours</b>		<b>102.55</b>	<b>41.30</b>
Qwen2.5-7B	LaCo	21.43%	611424.10	34.99
	MKA		190022.30	38.93
	<b>Ours</b>		<b>3458.57</b>	<b>40.11</b>
	LaCo	25.00%	478749.02	34.70
	MKA		607147.37	<b>37.09</b>
	<b>Ours</b>		<b>9401.84</b>	<b>36.07</b>

Table 4: Comparison of our learnable fusion strategy against similarity-based (MKA) and difference-based (LaCo) baselines under a fixed layer-merging scope.

414 compared to the baselines, often reducing it by an  
 415 order of magnitude. This substantial improvement  
 416 underscores the superiority of the learned coeffi-  
 417 cients, which is further evidenced by consistently  
 418 higher average scores across downstream tasks.

## 419 5.3 Sensitivity to Merging Granularity

420 We further investigate the impact of block size on  
 421 the layer merging strategy. In this experiment, with  
 422 the pruning rate fixed at approximately 20%, we  
 423 constrain the block size using minimum and max-  
 424 imum hyperparameters and evaluate the pruned  
 425 model based on perplexity and downstream task  
 426 performance. As illustrated in Figure 4, we observe

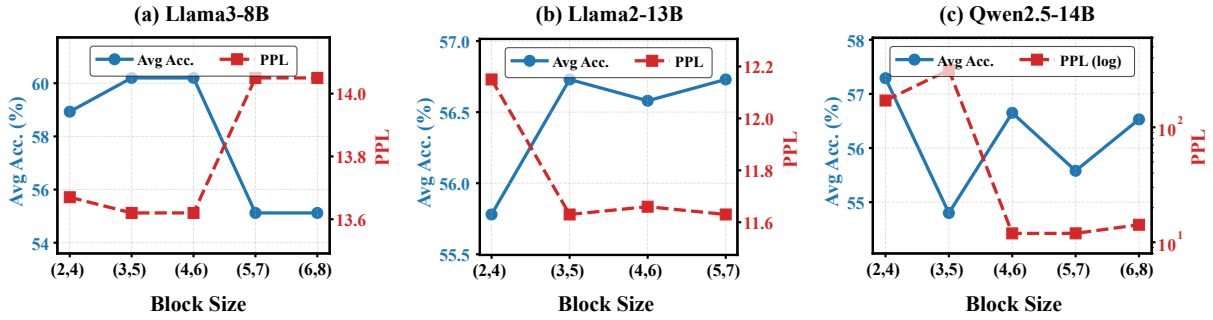


Figure 4: The impact of merging block size on model performance. Results are shown for (a) LLaMA3-8B, (b) LLaMA2-13B, and (c) Qwen2.5-14B. The tuples (min, max) on the x-axis denote the minimum and maximum block size constraints used for layer merging.

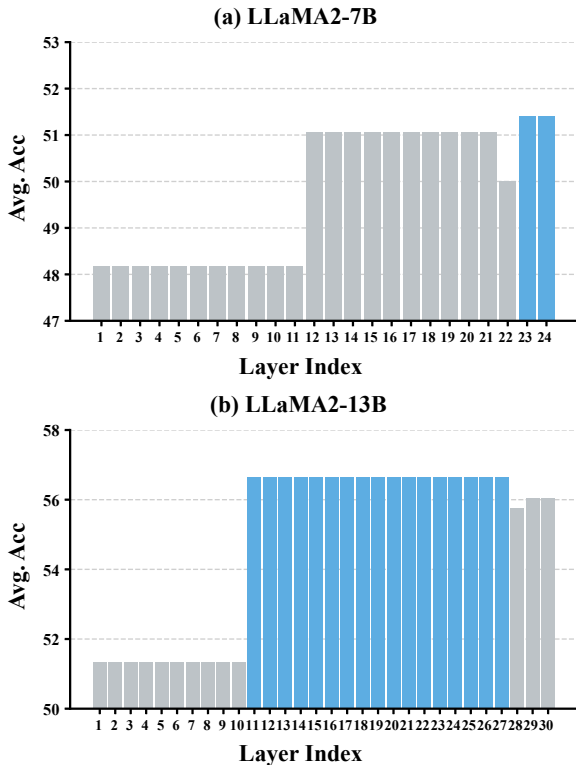


Figure 5: Impact of Merging Starting Depth. (a) and (b) show the results for LLaMA2-7B and LLaMA2-13B, respectively. The x-axis represents the starting layer index. Blue bars indicate the best performance achieved.

that a block size range of 3 to 6 maintains relatively higher performance and lower perplexity. We attribute this to a trade-off between merging efficiency and feature preservation. Excessively small block sizes result in a larger number of blocks dispersed across the model; this leads to more merged layers, causing error accumulation and amplification during forward propagation, which destabilizes performance. Conversely, overly large block sizes compress too many layers into a single layer, resulting in severe information loss and feature

blurring. Consequently, we conclude that the block size range of (3, 6) is optimal, as it typically yields two merged blocks, each containing three to five layers, thereby effectively preserving essential features. Further details regarding Figure 4 can be found in Appendix A.4.

#### 5.4 Impact of Merging Starting Depth

We also explore how the starting depth of merging affects model performance. Under a 20% pruning rate, we adjust the hyperparameter controlling the start point to influence the merging strategy discovery phase, effectively constraining the merging process to occur only after a designated layer. This design aligns with established observations (Cai et al., 2025; Xiao et al., 2024; Ma et al., 2025) that earlier layers in large language models are more crucial than later ones. The results in Figure 5 confirm this hypothesis, demonstrating that deferring merging to later layers results in superior efficacy on downstream tasks compared to compressing the earlier layers.

## 6 Conclusion

In this paper, we presented WaLeM, a framework designed to improve layer pruning by addressing the limitations of coarse-grained merging. By distinguishing the functional roles of weight matrices and integrating learnable merging coefficients, WaLeM enables more precise layer merging. Empirical results demonstrate that WaLeM consistently outperforms state-of-the-art baselines, particularly on reasoning-intensive benchmarks like GSM8K. Furthermore, our analysis underscores the critical impact of merging granularity and position, validating fine-grained layer merging as a robust solution for efficient LLM deployment.

## 473 Limitations

474 Despite WaLeM’s promising performance in LLM  
475 compression, we acknowledge several limitations  
476 remain. First, unlike training-free approaches,  
477 WaLeM entails a lightweight optimization phase  
478 for merging coefficients. While significantly more  
479 efficient than full fine-tuning, this process intro-  
480 duces computational overhead and necessitates ac-  
481 cess to calibration data. Furthermore, the layer se-  
482 lection mechanism relies on Centered Kernel Align-  
483 ment (CKA). Since CKA involves Gram matrix  
484 computation, its complexity scales quadratically  
485 with sequence length. This computational bottle-  
486 neck may prohibit scaling to extremely long con-  
487 text windows without approximation techniques.  
488 Finally, constrained by computational resources,  
489 our empirical validation is currently restricted to  
490 models up to the 14B parameter scale. Verify the  
491 efficacy of WaLeM on larger-scale foundation mod-  
492 els is a critical direction for future work.

## 493 Declaration of LLM Usage

494 Large language models were used in a lim-  
495 ited, assistive role during the preparation of this  
496 manuscript, primarily to improve language and  
497 clarity. All scientific content, interpretations, and  
498 conclusions are the sole work of the authors, who  
499 reviewed and approved the final version. The litera-  
500 ture review and reference curation were conducted  
501 independently by the authors using published and  
502 verifiable sources.

## 503 References

504 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama  
505 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,  
506 Diogo Almeida, Janko Altenschmidt, Sam Altman,  
507 Shyamal Anadkat, et al. 2023. Gpt-4 technical report.  
508 *arXiv preprint arXiv:2303.08774*.

509 Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng  
510 Gao, and Yejin Choi. 2019. [Piqa: Reasoning about  
511 physical commonsense in natural language](#). *Preprint*,  
512 arXiv:1911.11641.

513 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie  
514 Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind  
515 Neelakantan, Pranav Shyam, Girish Sastry, Amanda  
516 Askell, Sandhini Agarwal, Ariel Herbert-Voss,  
517 Gretchen Krueger, Tom Henighan, Rewon Child,  
518 Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,  
519 Clemens Winter, Christopher Hesse, Mark Chen,  
520 Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin  
521 Chess, Jack Clark, Christopher Berner, Sam Mc-  
522 Candlish, Alec Radford, Ilya Sutskever, and Dario

Amodei. 2020. [Language models are few-shot learn-  
ers](#). *Preprint*, arXiv:2005.14165. 523 524

Zefan Cai, Yichi Zhang, Bofei Gao, Yuliang Liu,  
Yucheng Li, Tianyu Liu, Keming Lu, Wayne Xiong,  
Yue Dong, Junjie Hu, and Wen Xiao. 2025. [Pyra-  
midkv: Dynamic kv cache compression based  
on pyramidal information funneling](#). *Preprint*,  
arXiv:2406.02069. 525 526 527 528 529 530

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming  
Yuan, Henrique Ponde de Oliveira Pinto, Jared Kap-  
lan, Harri Edwards, Yuri Burda, Nicholas Joseph,  
Greg Brockman, Alex Ray, Raul Puri, Gretchen  
Krueger, Michael Petrov, Heidy Khlaaf, Girish Sas-  
try, Pamela Mishkin, Brooke Chan, Scott Gray,  
Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz  
Kaiser, Mohammad Bavarian, Clemens Winter,  
Philippe Tillet, Felipe Petroski Such, Dave Cum-  
mings, Matthias Plappert, Fotios Chantzis, Eliza-  
beth Barnes, Ariel Herbert-Voss, William Hebg-  
Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie  
Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain,  
William Saunders, Christopher Hesse, Andrew N.  
Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan  
Morikawa, Alec Radford, Matthew Knight, Miles  
Brundage, Mira Murati, Katie Mayer, Peter Welinder,  
Bob McGrew, Dario Amodei, Sam McCandlish, Ilya  
Sutskever, and Wojciech Zaremba. 2021. [Evaluat-  
ing large language models trained on code](#). *Preprint*,  
arXiv:2107.03374. 531 532 533 534 535 536 537 538 539 540 541 542 543 544 545 546 547 548 549 550 551

Yu Cheng, Duo Wang, Pan Zhou, and Tao Zhang. 2020.  
[A survey of model compression and acceleration for  
deep neural networks](#). *Preprint*, arXiv:1710.09282. 552 553 554

Arjun Choudhry, Chang Liu, Nina Żukowska, Yifu Cai,  
Mononito Goswami, and Artur Dubrawski. 2025.  
[Layermerge: Modality-agnostic depth pruning for  
efficient foundation model deployment](#). In *NeurIPS  
2025 Workshop on Efficient Reasoning*. 555 556 557 558 559

Kevin Clark, Urvashi Khandelwal, Omer Levy, and  
Christopher D. Manning. 2019. [What does bert  
look at? an analysis of bert’s attention](#). *Preprint*,  
arXiv:1906.04341. 560 561 562 563

Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot,  
Ashish Sabharwal, Carissa Schoenick, and Oyvind  
Tafjord. 2018. [Think you have solved question  
answering? try arc, the ai2 reasoning challenge](#).  
*Preprint*, arXiv:1803.05457. 564 565 566 567 568

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,  
Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias  
Plappert, Jerry Tworek, Jacob Hilton, Reiichiro  
Nakano, Christopher Hesse, and John Schulman.  
2021. [Training verifiers to solve math word prob-  
lems](#). *Preprint*, arXiv:2110.14168. 569 570 571 572 573 574

Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and  
Luke Zettlemoyer. 2023. [Qlora: Efficient finetuning  
of quantized llms](#). *Preprint*, arXiv:2305.14314. 575 576 577

Elias Frantar and Dan Alistarh. 2023. [Sparsegpt: Mas-  
sive language models can be accurately pruned in  
one-shot](#). *Preprint*, arXiv:2301.00774. 578 579 580

581	Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. 2024. <a href="#">The language model evaluation harness</a> .	
582		
583		
584		
585		
586		
587		
588		
589	Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. 2021. <a href="#">Transformer feed-forward layers are key-value memories</a> . <i>Preprint</i> , arXiv:2012.14913.	
590		
591		
592	Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. 2024. <a href="#">The llama 3 herd of models</a> . <i>arXiv preprint arXiv:2407.21783</i> .	
593		
594		
595		
596		
597	Andrey Gromov, Kushal Tirumala, Hassan Shapourian, Paolo Glorioso, and Daniel A. Roberts. 2025. <a href="#">The unreasonable ineffectiveness of the deeper layers</a> . <i>Preprint</i> , arXiv:2403.17887.	
598		
599		
600		
601	Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang. 2024. <a href="#">MiniLLM: Knowledge distillation of large language models</a> . In <i>The Twelfth International Conference on Learning Representations</i> .	
602		
603		
604		
605	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. <a href="#">Measuring massive multitask language understanding</a> . <i>Preprint</i> , arXiv:2009.03300.	
606		
607		
608		
609	Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. <a href="#">Distilling the knowledge in a neural network</a> . <i>Preprint</i> , arXiv:1503.02531.	
610		
611		
612	Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. <a href="#">Lora: Low-rank adaptation of large language models</a> . <i>Preprint</i> , arXiv:2106.09685.	
613		
614		
615		
616	Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. 2024. <a href="#">Swe-bench: Can language models resolve real-world github issues?</a> <i>Preprint</i> , arXiv:2310.06770.	
617		
618		
619		
620		
621	Bo-Kyeong Kim, Geonmin Kim, Tae-Ho Kim, Thibault Castells, Shinkook Choi, Junho Shin, and Hyounghyung Song. 2024. <a href="#">Shortened llama: Depth pruning for large language models with comparison of retraining methods</a> . <i>Preprint</i> , arXiv:2402.02834.	
622		
623		
624		
625		
626	Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. 2019. <a href="#">Similarity of neural network representations revisited</a> . <i>Preprint</i> , arXiv:1905.00414.	
627		
628		
629		
630	Grigory Kovalev and Mikhail Tikhomirov. 2025. <a href="#">Iterative layer-wise distillation for efficient compression of large language models</a> . <i>Preprint</i> , arXiv:2511.05085.	
631		
632		
633		
	Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. 2017. <a href="#">Race: Large-scale reading comprehension dataset from examinations</a> . <i>Preprint</i> , arXiv:1704.04683.	634
		635
		636
		637
	Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. <a href="#">Let’s verify step by step</a> . <i>Preprint</i> , arXiv:2305.20050.	638
		639
		640
		641
		642
	Gui Ling, Ziyang Wang, Yuliang Yan, and Qingwen Liu. 2024. <a href="#">Slingpt: Layer-wise structured pruning for large language models</a> . <i>Preprint</i> , arXiv:2412.18110.	643
		644
		645
	Deyuan Liu, Zhanyue Qin, Hairu Wang, Zhao Yang, Zecheng Wang, Fangying Rong, Qingbin Liu, Yanchao Hao, Xi Chen, Cunhang Fan, Zhao Lv, Zhiying Tu, Dianhui Chu, Bo Li, and Dianbo Sui. 2025. <a href="#">Pruning via merging: Compressing llms via manifold alignment based layer merging</a> . <i>Preprint</i> , arXiv:2406.16330.	646
		647
		648
		649
		650
		651
		652
	Zechun Liu, Barlas Oguz, Changsheng Zhao, Ernie Chang, Pierre Stock, Yashar Mehdad, Yangyang Shi, Raghuraman Krishnamoorthi, and Vikas Chandr. 2023. <a href="#">Llm-qat: Data-free quantization aware training for large language models</a> . <i>Preprint</i> , arXiv:2305.17888.	653
		654
		655
		656
		657
		658
	Yao Lu, Yuqi Li, Wenbin Xie, Shanqing Yu, Qi Xuan, Zhaowei Zhu, and Shiping Wen. 2025. <a href="#">The structural scalpel: Automated contiguous layer pruning for large language models</a> . <i>Preprint</i> , arXiv:2510.23652.	659
		660
		661
		662
	Xinyin Ma, Gongfan Fang, and Xinchao Wang. 2023. <a href="#">Llm-pruner: On the structural pruning of large language models</a> . <i>Preprint</i> , arXiv:2305.11627.	663
		664
		665
	Zehong Ma, Shiliang Zhang, Longhui Wei, and Qi Tian. 2025. <a href="#">Efficient multi-modal long context learning for training-free adaptation</a> . <i>Preprint</i> , arXiv:2505.19812.	666
		667
		668
		669
	Xin Men, Mingyu Xu, Qingyu Zhang, Bingning Wang, Hongyu Lin, Yaojie Lu, Xianpei Han, and Weipeng Chen. 2024. <a href="#">Shortgpt: Layers in large language models are more redundant than you expect</a> . <i>Preprint</i> , arXiv:2403.03853.	670
		671
		672
		673
		674
	Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2023. <a href="#">Locating and editing factual associations in gpt</a> . <i>Preprint</i> , arXiv:2202.05262.	675
		676
		677
	Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. 2016. <a href="#">Pointer sentinel mixture models</a> . <i>Preprint</i> , arXiv:1609.07843.	678
		679
		680
	Zehua Pei, Hui-Ling Zhen, Xianzhi Yu, Sinno Jialin Pan, Mingxuan Yuan, and Bei Yu. 2025. <a href="#">From pruning to grafting: Dynamic knowledge redistribution via learnable layer fusion</a> . <i>Preprint</i> , arXiv:2411.14507.	681
		682
		683
		684
	Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin,	685
		686
		687

688	Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang,	Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali	745
689	Jiayi Yang, Jingren Zhou, Junyang Lin, Kai Dang,	Farhadi, and Yejin Choi. 2019. <a href="#">Hellswag: Can</a>	746
690	Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li,	<a href="#">a machine really finish your sentence?</a> <i>Preprint</i> ,	747
691	Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji	arXiv:1905.07830.	748
692	Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang		
693	Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang	Yang Zhang, Yawei Li, Xinpeng Wang, Qianli Shen,	749
694	Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru	Barbara Plank, Bernd Bischl, Mina Rezaei, and Kenji	750
695	Zhang, and Zihan Qiu. 2025. <a href="#">Qwen2.5 technical</a>	Kawaguchi. 2024. <a href="#">Finercut: Finer-grained inter-</a>	751
696	<a href="#">report</a> . <i>Preprint</i> , arXiv:2412.15115.	<a href="#">pretable layer pruning for large language models.</a>	752
		<i>Preprint</i> , arXiv:2405.18218.	753
697	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine		
698	Lee, Sharan Narang, Michael Matena, Yanqi Zhou,	Longguang Zhong, Fanqi Wan, Ruijun Chen, Xiaojun	754
699	Wei Li, and Peter J. Liu. 2023. <a href="#">Exploring the limits</a>	Quan, and Liangzhi Li. 2025. <a href="#">Blockpruner: Fine-</a>	755
700	<a href="#">of transfer learning with a unified text-to-text trans-</a>	<a href="#">grained pruning for large language models.</a> <i>Preprint</i> ,	756
701	<a href="#">former</a> . <i>Preprint</i> , arXiv:1910.10683.	arXiv:2406.10594.	757
702	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavat-	Xunyu Zhu, Jian Li, Yong Liu, Can Ma, and Weiping	758
703	ula, and Yejin Choi. 2019. <a href="#">Winogrande: An adver-</a>	Wang. 2024. <a href="#">A survey on model compression for</a>	759
704	<a href="#">sarial winograd schema challenge at scale.</a> <i>Preprint</i> ,	<a href="#">large language models.</a> <i>Preprint</i> , arXiv:2308.07633.	760
705	arXiv:1907.10641.		
706	Jiwon Song, Kyungseok Oh, Taesu Kim, Hyungjun		
707	Kim, Yulhwa Kim, and Jae-Joon Kim. 2024. <a href="#">Sleb:</a>		
708	<a href="#">Streamlining llms through redundancy verification</a>		
709	<a href="#">and elimination of transformer blocks.</a> <i>Preprint</i> ,		
710	arXiv:2402.09025.		
711	Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann		
712	Dubois, Xuechen Li, Carlos Guestrin, Percy Liang,		
713	and Tatsunori B. Hashimoto. 2023. <a href="#">Stanford alpaca:</a>		
714	<a href="#">An instruction-following llama model.</a> <a href="https://github.com/tatsu-lab/stanford_alpaca">https://</a>		
715	<a href="https://github.com/tatsu-lab/stanford_alpaca">github.com/tatsu-lab/stanford_alpaca</a> .		
716	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-		
717	bert, Amjad Almahairi, Yasmine Babaei, Nikolay		
718	Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti		
719	Bhosale, et al. 2023. <a href="#">Llama 2: Open founda-</a>		
720	<a href="#">tion and fine-tuned chat models.</a> <i>arXiv preprint</i>		
721	<i>arXiv:2307.09288</i> .		
722	Wenxiao Wang, Wei Chen, Yicong Luo, Yongliu Long,		
723	Zhengkai Lin, Liye Zhang, Binbin Lin, Deng Cai,		
724	and Xiaofei He. 2024. <a href="#">Model compression and effi-</a>		
725	<a href="#">cient inference for large language models: A survey.</a>		
726	<i>Preprint</i> , arXiv:2402.09748.		
727	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten		
728	Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and		
729	Denny Zhou. 2023. <a href="#">Chain-of-thought prompting elic-</a>		
730	<a href="#">its reasoning in large language models.</a> <i>Preprint</i> ,		
731	arXiv:2201.11903.		
732	Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song		
733	Han, and Mike Lewis. 2024. <a href="#">Efficient streaming lan-</a>		
734	<a href="#">guage models with attention sinks.</a> In <i>The Twelfth</i>		
735	<i>International Conference on Learning Representa-</i>		
736	<i>tions</i> .		
737	Yifei Yang, Zouying Cao, and Hai Zhao. 2024. <a href="#">Laco:</a>		
738	<a href="#">Large language model pruning via layer collapse.</a>		
739	<i>Preprint</i> , arXiv:2402.11187.		
740	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran,		
741	Thomas L. Griffiths, Yuan Cao, and Karthik		
742	Narasimhan. 2023. <a href="#">Tree of thoughts: Deliber-</a>		
743	<a href="#">ate problem solving with large language models.</a>		
744	<i>Preprint</i> , arXiv:2305.10601.		

## A Appendix

### A.1 Details of Implementations

In this section, we elaborate our experimental setup. We utilize the C4 dataset for the initial phases, randomly sampling 500 samples for CKA similarity matrix computation and 10,000 samples for learning the merging coefficients, with the random seed fixed at 42. Regarding the merging strategy, we constrain the block size between 4 and 6 for most models, while adjusting this range to 5 and 7 for Llama2-13B. We set the similarity threshold to 0.85 and  $\alpha$  to 1.5 in Equation 5, and  $\beta$  to 0.3 in Equation 6. The merging start depth is initialized at layer 17 for Llama3-8B and Llama2-7B, and at layer 22 for Llama2-13B. Conversely, we set the start depth to 0 for Qwen models to accommodate their distinct layer-wise redundancy patterns observed in the CKA analysis. During the coefficients learning process, we train for 2 epochs with a batch size of 4, gradient accumulation of 4, a learning rate of 0.1, and a temperature of 1. For the post-pruning recovery phase, we fine-tune the model on the Alpaca dataset for 2 epochs using LoRA. The hyperparameters include a learning rate of  $1 \times 10^{-4}$ , LoRA rank  $r = 8$ ,  $\alpha = 16$ , and a dropout of 0.05, with a batch size of 4 and gradient accumulation of 16. We adhere to the official configurations for all baseline methods with two exceptions. For LLM-Pruner, we adopt its layer pruning configuration; for LaCo, we terminate the iteration early once the desired pruning ratio is achieved.

### A.2 Details of Datasets

#### A.2.1 Pruning Datasets

**C4** (Colossal Clean Crawled Corpus) (Raffel et al., 2023) is a massive, cleaned version of the Common Crawl web corpus. A subset is sampled to compute the CKA similarity matrix, while a separate subset is employed to optimize the merging coefficients. This dataset serves as a standard benchmark for measuring perplexity on continuous text streams.

#### A.2.2 Post Pruning Recovery Datasets

**Alpaca** (Taori et al., 2023) To restore model performance after pruning, we employ the Alpaca dataset for instruction tuning. This dataset consists of synthetic instruction-response pairs designed to enhance the capability of the model to follow user commands and maintain response quality.

---

### Algorithm 1 Weight-Aware Learnable Merging Framework

---

**Input:** Pre-trained LLM  $\mathcal{M}$  with  $L$  layers, Calibration Dataset  $\mathcal{D}$ , Pruning Target  $K$

**Output:** Compressed Model  $\mathcal{M}'$

#### # Phase 1: Identifying Redundant Layers

- 1: Sample small subset  $\mathcal{D}_{sim} \subset \mathcal{D}$  where  $|\mathcal{D}_{sim}| = 500$  {See Appendix A.1}
- 2: Extract hidden states  $\mathbf{X}^{(l)}$  for all layers using  $\mathcal{D}_{sim}$
- 3: **for** each pair  $(l, k)$  in layers **do**
- 4:   Compute Gram matrices  $\mathbf{K}^{(l)}, \mathbf{K}^{(k)}$  based on  $\mathcal{D}_{sim}$
- 5:   Compute CKA similarity  $s_{lk}$  via HSIC in Eq. 4
- 6: **end for**
- 7: Construct similarity matrix  $\mathbf{S} \in \mathbb{R}^{L \times L}$
- 8: Initialize DP table  $V$  and solve for optimal blocks  $\mathbb{B}$  using Eq. 7

9: *Output:* Set of non-overlapping blocks  $\mathbb{B}$

#### # Phase 2: Learning Merging Coefficients

- 10: Sample larger subset  $\mathcal{D}_{train} \subset \mathcal{D}$  where  $|\mathcal{D}_{train}| = 10,000$  {See Appendix A.1}
  - 11: Initialize merging coefficients  $\mathbf{F}_\tau \in \mathbb{R}^{|\mathcal{T}| \times B}$  for blocks in  $\mathbb{B}$
  - 12: Freeze original LLM parameters
  - 13: **while** not converged **do**
  - 14:   Sample batch  $x$  from  $\mathcal{D}_{train}$
  - 15:   **for** each block in  $\mathbb{B}$  **do**
  - 16:     Compute merging coefficients via Eq. 8:  
     $\alpha_{\tau,l} = \text{Softmax}(F_{\tau,l})$
  - 17:     Merge weights via Eq. 9:  $\tilde{W}_\tau = \sum_{j=1}^B \alpha_{\tau,j} \cdot W_\tau^{(j)}$
  - 18:   **end for**
  - 19:   Forward pass to get logits  $P_{\mathcal{M}}(x)$  and  $P_{\mathcal{M}'}(x)$
  - 20:   Compute Loss in Eq. 10:  $\mathcal{L} = \mathcal{D}_{KL}(P_{\mathcal{M}} || P_{\mathcal{M}'})$
  - 21:   Update  $\mathbf{F}$  via gradient descent
  - 22: **end while**
  - 23: Apply learned  $\alpha$  to permanently merge layers
  - 24: Fine-tune  $\mathcal{M}'$  using LoRA on instruction data
  - 25: **return**  $\mathcal{M}'$
- 

### A.2.3 Evaluation Datasets

We implement the LM Evaluation Harness (Gao et al., 2024) with default configurations to assess performance across various downstream tasks.

**WikiText2** (Merity et al., 2016) This dataset evaluates language modeling capabilities on high-quality encyclopedic content. It retains original

808

809

810

811

812

813

814

Model	Method	Pruning Rate	PPL ( $\downarrow$ )	HellaSwag	PIQA	Winogrande	MMLU	ARC-e	ARC-c	RACE	GSM8K	Avg. Score
Llama2-7B	Dense	0.00%	8.76	76.19	78.73	69.46	41.71	73.82	44.97	40.10	13.42	54.80
	MKA	18.75%	1008.68	54.59	66.81	59.19	33.29	55.81	38.05	30.24	0.76	42.34
	ShortGPT	18.75%	27.04	65.18	72.03	<u>64.64</u>	25.90	57.45	38.31	33.49	1.74	44.84
	SLEB	18.75%	<u>14.77</u>	64.89	<u>74.48</u>	59.83	26.48	59.68	34.64	33.78	2.05	44.48
	LaCo	18.75%	100.73	58.21	70.46	62.43	31.68	57.95	37.71	31.77	0.45	43.83
	LLMPruner	18.75%	30.36	<u>68.27</u>	74.16	63.22	<u>34.49</u>	<u>65.70</u>	<u>41.21</u>	<u>35.60</u>	<b>6.37</b>	<u>48.63</u>
	Ours	18.75%	<b>14.69</b>	<b>71.52</b>	<b>75.79</b>	<b>67.25</b>	<b>40.14</b>	<b>66.50</b>	<b>41.89</b>	<b>39.43</b>	<u>5.00</u>	<b>50.94</b>
	MKA	25.00%	1842.92	49.81	65.72	57.70	25.06	52.10	36.26	30.05	0.38	39.64
	ShortGPT	25.00%	47.19	59.53	68.44	<u>66.06</u>	23.49	52.78	34.04	32.06	1.52	42.24
	SLEB	25.00%	<u>20.33</u>	60.91	71.60	58.56	27.66	54.21	32.08	31.67	1.52	42.28
	LaCo	25.00%	145.97	51.37	65.51	59.12	<u>31.33</u>	52.99	34.22	31.29	0.30	40.77
	LLMPruner	25.00%	40.57	<u>65.63</u>	<b>72.42</b>	62.35	25.14	61.83	<b>39.33</b>	<u>35.12</u>	<b>4.32</b>	<u>45.77</u>
	Ours	25.00%	<b>20.28</b>	<b>67.06</b>	<u>72.09</u>	<b>68.19</b>	<b>38.72</b>	<b>62.25</b>	<u>38.57</u>	<b>38.56</b>	<u>3.26</u>	<b>48.59</b>
	MKA	31.25%	3569.59	46.11	61.53	58.56	30.44	48.40	32.42	31.39	0.30	38.64
	ShortGPT	31.25%	119.99	50.84	63.11	<u>63.54</u>	<u>37.22</u>	45.03	32.17	32.63	0.53	40.63
	SLEB	31.25%	<u>28.78</u>	54.27	68.50	<u>54.22</u>	25.64	50.25	31.06	30.62	1.44	39.50
	LaCo	31.25%	117.32	46.32	61.10	60.54	34.39	37.92	32.08	31.48	1.06	38.11
	LLMPruner	31.25%	51.49	<u>61.89</u>	<u>70.02</u>	61.64	28.31	<b>56.82</b>	<b>35.92</b>	<u>35.02</u>	<u>2.12</u>	<u>43.97</u>
Ours	31.25%	<b>26.88</b>	<b>62.04</b>	<b>70.18</b>	<b>66.22</b>	<b>37.91</b>	<u>54.88</u>	<u>34.39</u>	<b>37.03</b>	<b>2.20</b>	<b>45.61</b>	

Table 5: Zero-shot performance of Llama2-7B model with distinct pruning strategies. “Dense” denotes the unpruned models and “PPL” denotes the perplexity evaluated on Wikitext2. **Bold** fonts indicate the best results, while underlined values represent the sub-optimal performance.

article structures to test the ability of the model to handle long-term dependencies.

**HellaSwag** (Zellers et al., 2019) As a benchmark for commonsense natural language inference, HellaSwag challenges models to select the most plausible continuation of a sentence. It employs adversarial filtering to ensure the task requires genuine reasoning rather than simple statistical matching.

**PIQA** (Bisk et al., 2019) Focusing on physical commonsense, PIQA presents questions regarding the interaction of physical objects. The model must identify the plausible solution to a specified goal, testing its grounding in real-world physics.

**Winogrande** (Sakaguchi et al., 2019) This dataset is designed to evaluate robust commonsense reasoning and pronoun resolution. It filters out biases prevalent in earlier datasets to ensure performance relies on reasoning capabilities rather than algorithmic artifacts.

**MMLU** (Hendrycks et al., 2021) MMLU serves as a comprehensive test of general world knowledge and problem-solving ability. It covers a wide array of subjects across STEM, the humanities, and social sciences, ranging from elementary to professional levels.

**ARC** (Clark et al., 2018) The AI2 Reasoning Challenge (ARC) evaluates grade-school scientific knowledge. We report results on both the Easy subset (ARC-e), which tests basic fact retrieval, and the Challenge subset (ARC-c), which requires deep reasoning to solve complex scientific scenarios.

**RACE** (Lai et al., 2017) Collected from English

reading comprehension examinations, RACE requires diverse reasoning skills including summarization, inference, and detail extraction to select correct answers based on provided passages.

**GSM8K** (Cobbe et al., 2021) This benchmark assesses multi-step mathematical reasoning. It consists of high-quality grade school math word problems that require the model to generate a chain of thought to reach the correct numerical solution.

### A.3 Extended Evaluation on LLaMA2-7B

Table 5 presents the performance comparison between WaLeM and baseline methods on Llama2-7B across various pruning rates. Our proposed WaLeM consistently outperforms existing baselines on the majority of tasks. We observe particularly significant improvements on HellaSwag, Winogrande, MMLU and RACE. Notably, WaLeM retains over 90% of the original performance on MMLU and RACE even with 31.25% of layers pruned.

### A.4 Detailed Results for Sensitivity to Merging Granularity

Table 6 details the numerical results corresponding to Figure 4. We explored various block size constraints to determine the optimal trade-off between compression rate and model performance. Our experiments reveals an inverted-U performance trajectory. Small blocks (2 to 4 layers) fail to fully exploit redundancy between consecutive layers, leading to an increased number of blocks and merged layers. This fragmentation exacerbates error accumulation

Model	Size Range	PPL ( $\downarrow$ )	Avg. Score	Block Indices
Llama3-8B	2-4	13.67	58.93	{18-19}, {21-24}, {27-28}, {29-30}
	3-5	13.62	60.20	{21-24}, {27-30}
	4-6	13.62	60.20	{21-24}, {27-30}
	5-7	14.05	55.12	{18-24}
	6-8	14.05	55.12	{18-24}
Llama2-13B	2-4	12.15	55.78	{26-29}, {30-31}, {32-33}, {34-36}, {37-38}
	3-5	11.63	56.73	{27-31}, {32-36}
	4-6	11.66	56.58	{26-29}, {30-35}
	5-7	11.63	56.73	{27-31}, {32-36}
Qwen2.5-14B	2-4	170.80	57.29	{1-2}, {4-6}, {11-13}, {20-21}, {23-24}, {26-27}, {28-29}, {42-43}
	3-5	315.61	54.80	{4-6}, {11-13}, {17-19}, {20-22}, {23-25}
	4-6	11.95	55.65	{11-14}, {20-24}, {26-29}
	5-7	11.96	55.58	{11-15}, {20-26}
	6-8	14.21	56.53	{20-25}, {36-41}

Table 6: Performance comparison of pruned models with different merging block sizes at a 20% pruning ratio. "Block Indices" column details the specific layer indices contained within each merged block.

during forward propagation. Conversely, larger blocks (6 to 8 layers) degrade performance by over-compressing distinct feature representations. Empirically, a block size range of 3 to 6 offers the most effective granularity, yielding the highest average accuracy and lowest perplexity on Llama3-8B and Llama2-13B.

We formally demonstrate that small block sizes lead to higher error accumulation. Consider a network interval of 4 layers. We compare the Small Block Strategy (merging two 2-layer groups sequentially) against the Large Block Strategy (merging 4 layers directly).

Let  $\epsilon$  denote the forward propagation error arising from merging 2 layers. Assuming the fusion error scales linearly with the number of layers, the error for a 4-layer merge is  $2\epsilon$ . Let  $\lambda$  represent the feature amplification factor of the network layers.

For the Small Block Strategy, the error from the first merged group is amplified by  $\lambda$  during forward propagation and accumulates with the error from the second group:

$$E_{small} = \lambda \cdot \epsilon + \epsilon = \epsilon(\lambda + 1) \quad (11)$$

In contrast, the Large Block Strategy introduces the fusion error only once:

$$E_{large} = 2\epsilon \quad (12)$$

Comparing the two strategies:

$$\frac{E_{small}}{E_{large}} = \frac{\epsilon(\lambda + 1)}{2\epsilon} = \frac{\lambda + 1}{2} \quad (13)$$

Since deep networks inevitably amplify input perturbations ( $\lambda > 1$ ), we derive:

$$\frac{\lambda + 1}{2} > 1 \implies E_{small} > E_{large} \quad (14)$$

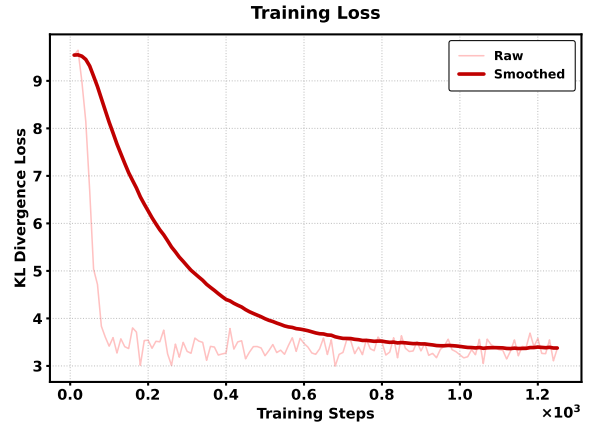


Figure 6: KL divergence loss during the merging coefficients learning phase on Llama2-13B with a 30% pruning rate.

This inequality theoretically validates that merging of small blocks results in greater total error compared to direct merging with larger blocks.

## A.5 Pseudocode of WaLeM framework

We provide the comprehensive pseudocode for the WaLeM framework in Algorithm 1. This process encompasses the three distinct phases: global redundancy analysis using CKA, learnable merging training via knowledge distillation, and the final post-fusion recovery.

## A.6 Details of Training Process

We detail the optimization process of merging coefficients for Llama2-13B with a 30% pruning rate. Figure 6 illustrates the KL divergence loss during this phase. Given that the merging strategy comprises three distinct blocks, Figure 7–9 depict

926 the evolution of merging coefficients for specific  
927 weight types across different layers.

### 928 **A.7 Details of Software Environment**

929 All experiments were implemented using Python  
930 3.11.14 in a Conda environment. The major soft-  
931 ware components and their corresponding versions  
932 are detailed as follows:

<b>Software</b>	<b>Version</b>
PyTorch	2.9.1
Transformers	4.57.1
CUDA Libraries	12.8
PEFT	0.18.0
LM Eval	0.4.9.1
Accelerate	1.11.0
Datasets	4.4.1
Numpy	2.3.4
Pandas	2.3.3

933 Additional dependencies include scikit-learn  
934 (1.7.2), wandb (0.23.0), and various utilities for  
935 data processing and model optimization.

### 936 **A.8 Case Study**

937 Table 7 presents qualitative samples generated by  
938 Llama3-8B at a 31.25% pruning rate. We observe  
939 that our WaLeM method effectively mitigates the  
940 redundancy that typically leads to semantic drift  
941 or tangential digression in the original model. As  
942 a result, the pruned model generates text with im-  
943 proved discourse structure and logical coherence,  
944 even at high compression rate. Although minor  
945 numerical hallucinations persist, the pruning pro-  
946 cess demonstrates a significant de-noising effect:  
947 it appears to filter out low quality corpus patters,  
948 such as web crawl artifacts, and yields cleaner and  
949 more focused outputs.

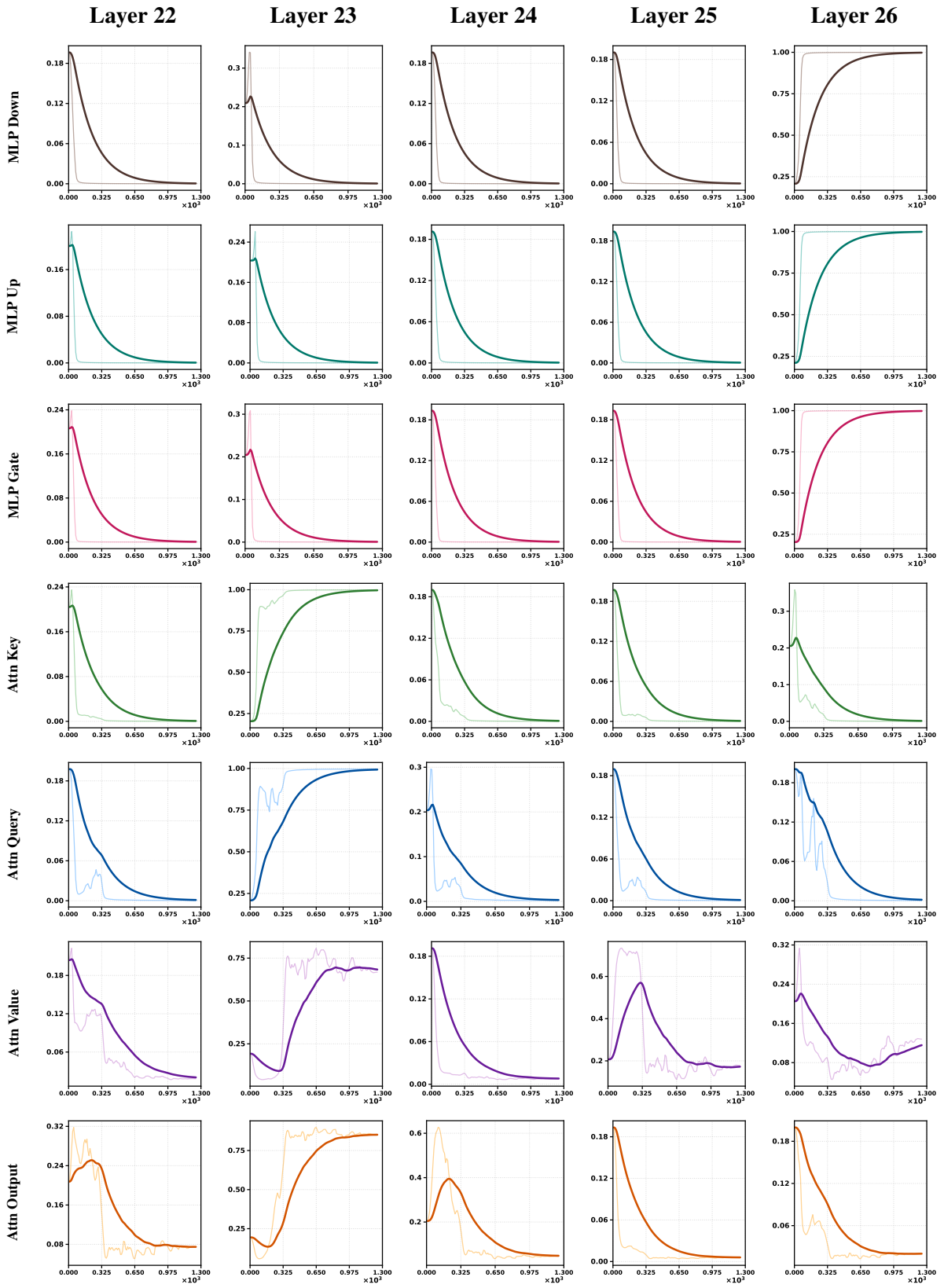


Figure 7: Visualization of the learning process of merging coefficients across layers within Block 1. **Left Axis:** Weight Types. **Top Axis:** Layer Index. Light lines represent raw data; bold lines represent smoothed trends.

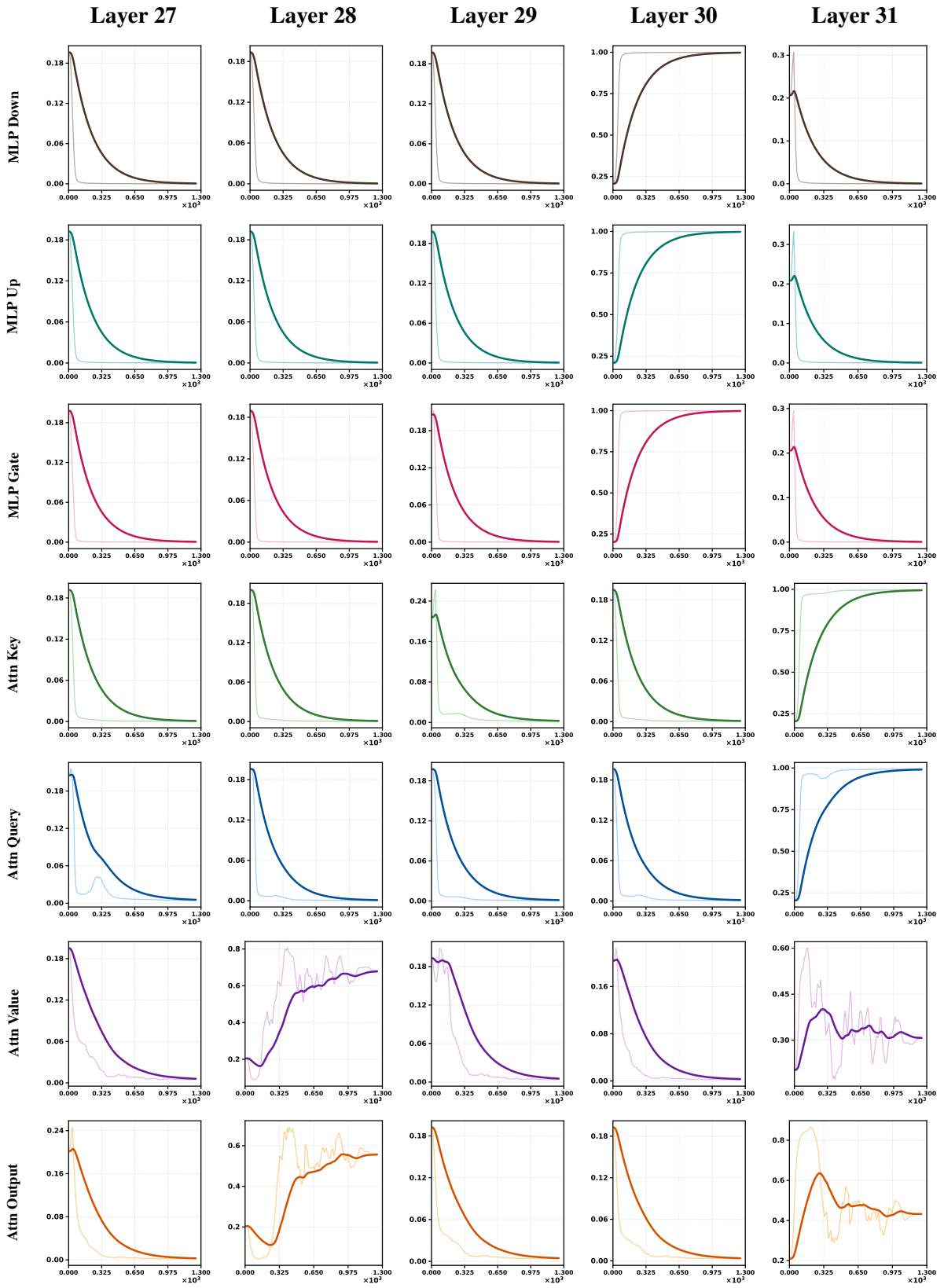


Figure 8: Visualization of the learning process of merging coefficients across layers within Block 2. **Left Axis:** Weight Types. **Top Axis:** Layer Index. Light lines represent raw data; bold lines represent smoothed trends.

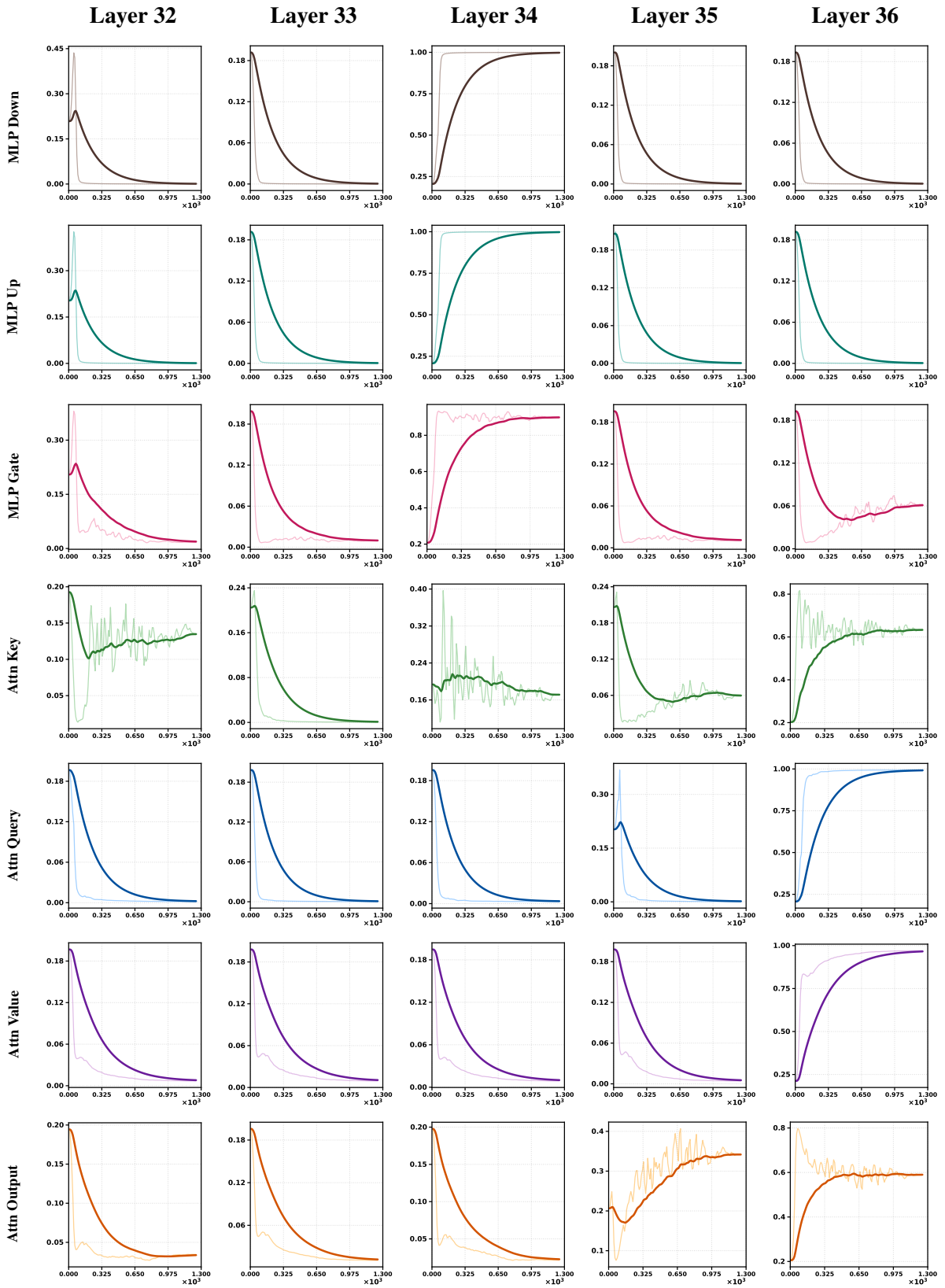


Figure 9: Visualization of the learning process of merging coefficients across layers within Block 3. **Left Axis:** Weight Types. **Top Axis:** Layer Index. Light lines represent raw data; bold lines represent smoothed trends.

Original	<p>Creating a strong brand identity is essential for any modern business. While a logo is a good start, the most successful companies understand that branding is actually about more than just what you see and hear. The best brands also consider how they are perceived by their customers.</p> <p>When it comes to establishing trust in your audience, there’s nothing like having an expert speak on behalf of your company. So why not use the opportunity?</p> <p>Hiring spokespersons can be expensive and time-consuming – but these days, video content has become increasingly important as well!</p> <p>In this article we’ll discuss 5 benefits associated with hiring experts who will represent yourself or organization through various means including social media campaigns such as YouTube videos which offer great exposure at little cost when compared against traditional advertising methods (TV spots).</p> <p>1. Expertise: Hiring someone else speaks volumes</p> <p>If you’re looking for something specific, then chances are other people have found success already. For example if I want some help figuring out my taxes because tax season is coming up soon – would anyone think twice before going online searching “tax accountant near me”? Well no one really needs to since all sorts of websites exist dedicated solely towards finding answers from professionals within minutes rather than hours waiting around while trying things themselves without guidance first hand experience etcetera...</p>
Pruned	<p>Creating a strong brand identity is essential for any modern business. While a logo is a good start, the most successful companies understand that branding is actually about creating unique and memorable experiences for their customers.</p> <p>Here’s how to create an effective brand image:</p> <ol style="list-style-type: none"> <li>1. Identify your target audience: Before you can create anything from scratch it’s important to know who you’re targeting in terms of demographics such as age group or job titles etc</li> <li>2. Define your value proposition: To have a clear idea which makes us different, we need to define our value proposition - what exactly do they pay for?</li> <li>3. Create visual elements: After identifying your target audiences and defining values it becomes easy to choose colors, fonts, images &amp; design style that resonates with them</li> <li>4. Develop your digital presence: With more than 9 billion users worldwide on social media platforms and millions every second considering whether there is enough room left between two sides? Social media has become one of the easiest ways to reach out to potential buyers. So having presence at all possible platforms like Facebook, Instagram, Twitter, LinkedIn, YouTube helps make a short attention span user engaged for longer time longer period before turning away by providing engaging and informative information!</li> <li>5. Engage with customer feedback: When you want people to buy your</li> </ol>
Original	<p>Explain the concept of ‘Photosynthesis’ to a 5-year-old child using a simple analogy. Photosynthesis is like cooking. Plants use sunlight as energy, take in carbon dioxide from air and water from soil, make their own food (glucose) which they need for growth, just like we cook our meals at home.</p> <p>You might also want to ask your kid how plants grow bigger every day? It’s because it needs glucose or sugar that comes out of photosynthesis! So next time you are outside playing with them, tell them about this amazing process called Photosynthesis!</p> <p>Do share any other good examples if you have heard so far!</p> <p>- <a href="https://www.facebook.com/sharer/sharer.php?u=https://aakashdigitalsrv1.meritnation.com/ask-answer/qu">https://www.facebook.com/sharer/sharer.php?u=https://aakashdigitalsrv1.meritnation.com/ask-answer/qu</a></p> <p>- <a href="/board/detail/board.do?gnbrd=&amp;cmd=attach&amp;atch_file_id=17007054">/board/detail/board.do?gnbrd=&amp;cmd=attach&amp;atch_file_id=17007054</a></p> <p>The following reaction represents photosynthesis:</p> $6\text{CO}_2 + 12\text{H}_2\text{O} \rightarrow \text{C}_6\text{H}_{12}\text{O}_6 + 6\text{O}_2$ <p>We can represent it by considering each molecule on its side</p> $\text{C O} \rightarrow \text{C O CO CO CO OH HO}$
Pruned	<p>Explain the concept of ‘Photosynthesis’ to a 5-year-old child using a simple analogy. Photosynthesis is like cooking. Plants use sunlight as an ingredient and feed it into their bodies, then produce oxygen which can be used by humans or animals that breathe it out at all time while we eat plants such as fruits, vegetables, rice, potatoes etc. So instead of us cooks directly from our food sources, plants act as intermediaries for their own nutritional needs. It’s just like how people cook and consume different foods depending on what they taste or have preferences with regard to healthy diets and nutrition. The same applies when photosynthesis occurs within the plant itself. Like how things are put in order through the process of cooking so too do photosynthesis involve intricate processes involving light energy (sun), carbon dioxide air, water molecules being converted into glucose (sugar) and oxygen, essential nutrients for both human and plants alike.</p>

Table 7: Generated samples from original Llama3-8B and the model pruned via WaLeM with a 31.25% pruning rate.