PLUG-AND-FOLD: WEIGHT-PRESERVING STRUCTURED COMPRESSION FOR LARGE LANGUAGE MODELS

Anonymous authors

000

001

002

004

006

008 009 010

011 012

013

014

015

016

017

018

019

021

023

025

026

027

028

029

031

032

034

037

038

040

041 042

043

044

046

047

048

051

052

Paper under double-blind review

ABSTRACT

Large Language Models (LLMs) have achieved remarkable performance across a wide range of tasks, but their growing size poses significant challenges for deployment and efficiency. Among existing model compression methods, structured pruning has emerged as a popular approach for reducing model size. However, pruning removes structural components such as layers, heads, or channels, which can disrupt pre-trained weights and lead to fragile recovery fine-tuning process. In this work, we propose Plug-and-Fold (PnF), a weight-preserving yet structurally effective compression method. Rather than removing weights or modifying the model architecture, PnF introduces lightweight, learnable adapter modules into the projection layers of attention and feed-forward networks. These adapters are trained while keeping the original weights frozen, and are later folded into the base weights via simple matrix multiplications. This process yields a compressed model that is structurally identical to the original and incurs no additional runtime overhead. We evaluate PnF across a variety of benchmarks and model scales, demonstrating consistent improvements over recent state-of-the-art structured compression baselines. Our results highlight that preserving the integrity of pretrained weights not only simplifies the compression pipeline, but also improves generalization and performance recovery in compressed LLMs.

1 Introduction

Large language models (LLMs) based on the Transformer (Vaswani et al., 2017) have achieved remarkable progress across various domains, including natural language processing (Zhao et al., 2023; Jiang et al., 2024a; Radford et al., 2018), code generation (Jiang et al., 2024b), computer vision (Liu et al., 2023a; Hamadi, 2023), and scientific applications (Zhang et al., 2025; Lin et al., 2023). This progress is attributable to two factors: (1) scaling model size to billions to trillions of parameters (Team et al., 2024; Islam & Moushi, 2025; Team et al., 2025; Zhang & Sennrich, 2019) and (2) pre-training on massive, diverse corpora (Langlais et al., 2025; Liu et al., 2024). Together, these endow LLMs with deep language understanding and ability to generate high-quality code, text, and multi-modal contents.

Despite these successes, their massive arameter sizes pose critical challenges: they require large storage, memory footprints, increase inference latency, and substantial computation for training and deployment, especially in resource-constrained settings. To address these practical limitations, a substantial body of research has focused on model compression techniques that shrink the footprint while preserving performance. These methods can be grouped into three principal categories: (1) knowledge distillation, which transfers capabilities from a large teacher to a smaller student (Hinton, 2014; Ojha et al., 2023; Agarwal et al., 2023; Bing et al., 2025; Cui et al., 2025); (2) quantization, which lowers numerical precision to save memory and accelerate inference (Liu et al., 2023b; Li et al., 2024b; Shang et al., 2023; Hu et al., 2025; An et al., 2025); and (3) pruning, a structured approach that removes redundant channels, heads, or layers (Voita et al., 2019; Gao et al., 2024b; Ma et al., 2023; Ashkboos et al., 2024; Men et al., 2024; Mugnaini et al., 2025; Yang et al., 2024).

Pruning gained a lot of attention since it leverages the pre-trained weights of the original model and typically does not require to training a new network from the ground up. Moreover, once

the unnecessary components have been eliminated, the resulting model can be further compressed through quantization, yielding additional reductions in memory consumption and inference latency. In the context of LLMs, most prior works on pruning focuses on three kinds of structural reductions: (i) deleting channels from the projection weights in attention and feed-forward network (Ashkboos et al., 2024; Gao et al., 2024b; Ma et al., 2023), (ii) removing heads in the multi-head attention (Voita et al., 2019; Mugnaini et al., 2025), and (iii) pruning whole transformer layers (Yang et al., 2024; Men et al., 2024). The selection of components to prune is guided by metrics that estimate the impact of removal, such as the magnitude of weight and activation (Sun et al.), cosine similarity (Men et al., 2024), or the L2-norm (Ashkboos et al., 2024). Although pruning leaves the overall transformer architecture intact, it disrupts parameters that were carefully tuned during large-scale pretraining, leading to inevitable performance loss. Consequently, many approaches incorporate a recovery finetuning (RFT) stage to restore accuracy, often employing the lightweight adapter like LoRA (Voita et al., 2019; Gao et al., 2024b; Ma et al., 2023; Ashkboos et al., 2024; Men et al., 2024; Mugnaini et al., 2025; Yang et al., 2024b. However, the recovery process can be fragile: even extensive RFT often fails to fully restore the performance of precisely optimized foundation models.

To overcome these limitations, we propose a **weight-preserving** structured compression that retain the integrity of pretrained weight while still achieving substantial efficiency gains. Our method, Plug-and-Fold (PnF), inserts lightweight, learnable adapter modules into the original projection matrices of the attention and feed-forward sub-layers rather than removing heads, channels, or layers. During training, only the adapters are updated while the pretrained weights remains completely frozen, preserving the expressivity and knowledge encoded in the original model. Once training is complete, the adapters are folded back into the base weights by simple matrix multiplications, resulting in a compressed model that is structurally identical to the original. Because no architectural modification is introduced and no extra operations are required during inference, PnF can be integrated seamlessly into existing serving frameworks and hardware accelerators.

We evaluate PnF with extensive experiments covering a broad spectrum of model sizes and compression rates. To validate its effectiveness, PnF is benchmarked against the latest state-of-the-art structured-compression baselines on a diverse set of tasks that demand varied domain knowledge and comprehensive capabilities. Across all settings, PnF consistently surpasses existing methods, delivering notable gains in downstream performance. These results show that preserving the integrity of pretrained weights not only yields a simpler and more scalable compression pipeline, but also enhances the recovery of accuracy and the generalization ability of the compressed models.

The main contributions of our paper are summarized as follows:

- We propose Plug-and-Fold (PnF), a novel weight preserving structured compression method that inserts lightweight, learnable adapter modules into the original projection layers without modifying the model architecture.
- After training, the adapters are folded into the base weights via simple matrix multiplications, resulting in a compressed model that is structurally identical to the original model and reduces runtime effectively.
- Extensive experiments demonstrate that PnF outperforms recent state-of-the-art structured-compression baselines across a wide range of model scales and benchmark tasks, confirming its effectiveness and scalability.

2 BACKGROUND

2.1 Decoder-based Transformer Architecture

Large Language Models (LLMs) primarily leverage a decoder-based Transformer architecture composed of stacked decoder blocks. These blocks consist of two core components: the Multi-Head Self-Attention (MHSA) mechanism and the Feed Forward Network (FFN). These components form the core layers of decoder blocks, enabling sequential data processing and contextual understanding.

2.1.1 Multi-Head Self-Attention (MHSA)

The MHSA mechanism enables the model to dynamically weight and aggregate contextual information from different positions in the input sequence by utilizing attention heads. Formally, let the l-th

decoder block takes input hidden state $X^{(l-1)} \in \mathbb{R}^{n \times d_{\text{embed}}}$, where n and d_{embed} is the length and the dimension of the input, respectively. For the i-th attention head, $i \in \{1, \cdots, n_h\}$, the MHSA mechanism computes the query vectors $Q_i^{(l)} \in \mathbb{R}^{n \times d_{\text{head}}}$, key vectors $K_i^{(l)} \in \mathbb{R}^{n \times d_{\text{head}}}$, and value vectors $V_i^{(l)} \in \mathbb{R}^{n \times d_{\text{head}}}$ as follows:

$$Q_i^{(l)} = X^{(l-1)} W_{Q_i^{(l)}}, \quad K_i^{(l)} = X^{(l-1)} W_{K_i^{(l)}}, \quad V_i^{(l)} = X^{(l-1)} W_{V_i^{(l)}}, \tag{1}$$

where $W_{Q_i^{(l)}},~W_{K_i^{(l)}},~W_{V_i^{(l)}} \in \mathbb{R}^{d_{\mathrm{embed}} \times d_{\mathrm{head}}}$ are the learned weight parameters for query, key, and value projections, and d_{head} is the dimension of the head (often $d_{\mathrm{head}} = \frac{d_{\mathrm{embed}}}{n_h}$). Then, the self-attention operation is applied to each triple $(Q_i^{(l)}, K_i^{(l)}, V_i^{(l)})$ and computes the attention output of the i-th head $Z_i^{(l)}$ as follows:

$$Z_{i}^{(l)} = \text{Attention}(Q_{i}^{(l)}, K_{i}^{(l)}, V_{i}^{(l)}) = \text{Softmax}\Big(\frac{Q_{i}^{(l)} \left(K_{i}^{(l)}\right)^{\top}}{\sqrt{d_{k}}}\Big) V_{i}^{(l)}, \tag{2}$$

where $\sqrt{d_k}$ is a scaling factor applied to ensure numerical stability. To represent comprehensive contextual information, these outputs from individual heads are concatenated and transformed as follows:

$$Z^{(l)} = \operatorname{Concat}(Z_1^{(l)}, \dots, Z_h^{(l)}) W_{O^{(l)}} \in \mathbb{R}^{n \times d_{\text{embed}}}, \tag{3}$$

where $\operatorname{Concat}(\cdot)$ is the concatenation operation and $W_{O^{(l)}} \in \mathbb{R}^{(hd_{\operatorname{head}}) \times d_{\operatorname{embed}}}$ is learned weight parameters for output.

2.1.2 FEED-FORWARD NETWORK (FFN)

Following the MHSA mechanism, the output is passed through a Feed Forward Network (FFN) to enhance the model's capacity to process through non-linear transformations and increased number of parameters. The FFN is often applies linear transformations separated by a nonlinear activation function $\sigma(\cdot)$ (e.g., SiLU(Elfwing et al., 2018)). For example, SwiGLU (Shazeer, 2020) module is defined as follows:

$$SwiGLU(Z^{(l)}) = \left(\sigma(Z^{(l)}W_{gate^{(l)}}) \odot Z^{(l)}W_{up^{(l)}}\right)W_{down^{(l)}}$$

$$\tag{4}$$

where σ is the Swish activation function (Ramachandran et al., 2018), and $W_{\text{gate}^{(l)}}, W_{\text{up}^{(l)}} \in \mathbb{R}^{d_{\text{embed}} \times d_{\text{inter}}}$, and $W_{\text{up}^{(l)}} \in \mathbb{R}^{d_{\text{inter}} \times d_{\text{embed}}}$ are learnable parameters with the intermediate dimension d_{inter} .

3 METHOD

In this section, we present Plug-and-Fold (PnF) compression, a straightforward yet effective compression method for large language models, whose complete workflow is illustrated in Figure 1¹. The main objective of this method is to preserve the original projection weight during training while reducing their dimensionality, yeidling a compact model that maintains the original signal.

Section 3.1 introduce the PnF adapter, a foldable compression module plugged into the original projection weights and trained to induce low-dimensional projection while preserving the original signal. Section 3.2 presents training schemes used to train these adapters effectively. Finally, Section 3.3 describes how the trained PnF adapters are folded into low-dimensional projection weights, producing a compact model that is computationally efficient while preserving performance suitable for deployment.

3.1 PLUG-AND-FOLD (PNF) COMPRESSION

3.1.1 PLUG-AND-FOLD (PNF) ADAPTER

In order to preserve the original signal while training, Plug-and-Fold adapters are plugged into the pre-trained model. Given a pre-trained linear weight $W \in \mathbb{R}^{m \times n}$, we define the PnF adapter as a linear projection:

$$P \in \mathbb{R}^{n \times r},\tag{5}$$

¹Snowflake and Fire icons created by Freepik – Flaticon

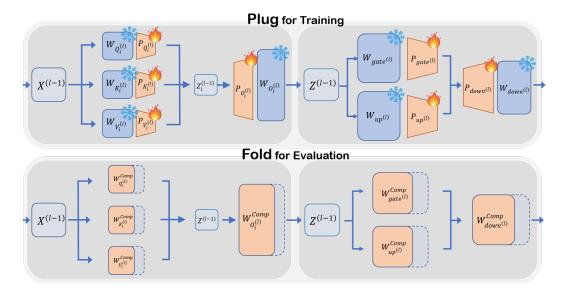


Figure 1: **Visualization of Plug-and-Fold framework**. The top half illustrates the training phase: lightweight PnF adapters are *plugged* into the pretrained linear layers and project to a reduced-dimensional space; the backbone weights remain frozen (shown as snowflakes), while the adapters are the only trainable components (shown as fire), enabling them to fully leverage the already-optimized structure. The bottom half shows the evaluation phase: after training, each adapter is *folded* back into its corresponding weight matrix via a simple matrix multiplication, yielding a compressed model that preserves the original architecture, interface, and performance.

where r < n. The adapter is applied to W and subsequently trained to recover the performance of the original model. Formally, our aim is to find an adapter P that satisfies:

$$\mathcal{P}(W) \approx \mathcal{P}(WP),$$
 (6)

where $\mathcal{P}(\cdot)$ denotes the performance measures on various tasks induced by the corresponding weight. Consequently, projecting the weights through the trained adapter P that satisfies Eq. (6) yields output representations in the reduced-dimensional space (r-dimension), while preserving a quality comparable to that of the full-size model. i.e., this projection yields compact representations that preserve the fidelity of the original weight matrix, allowing highly efficient deployment across a broad range of downstream tasks.

3.1.2 PNF ADAPTER FOR MHSA

We now explain how PnF adapter is integrated into the MHSA layer of an LLM. Let the projection weights for queries, keys, values and the output at layer l be $W_{Q_i^{(l)}}, W_{K_i^{(l)}}, W_{V_i^{(l)}} \in \mathbb{R}^{d_{\text{embed}} \times d_{\text{head}}}$, and $W_{O^{(l)}} \in \mathbb{R}^{(n_h d_{\text{head}}) \times d_{\text{embed}}}$, where n_h is the number of attention heads. For each of these matrices, we plug in a corresponding PnF adapter with dimension $r_{\text{head}} < d_{\text{head}}$:

$$P_{Q_i^{(l)}}, \ P_{K_i^{(l)}}, \ P_{V_i^{(l)}} \in \mathbb{R}^{d_{\text{head}} \times r_{\text{head}}^{(l)}}, \ \text{and} \ P_{O^{(l)}} \in \mathbb{R}^{(n_h r_{\text{head}}^{(l)}) \times (n_h d_{\text{head}})} \tag{7}$$

These adapters, multiplied with the original weights, produce lower-dimensional projections:

$$\begin{split} W_{Q_i^{(l)}} P_{Q_i^{(l)}} &\in \mathbb{R}^{d_{\text{embed}} \times r_{\text{head}}^{(l)}} \\ W_{K_i^{(l)}} P_{K_i^{(l)}} &\in \mathbb{R}^{d_{\text{embed}} \times r_{\text{head}}^{(l)}} \\ W_{V_i^{(l)}} P_{V_i^{(l)}} &\in \mathbb{R}^{d_{\text{embed}} \times r_{\text{head}}^{(l)}} \\ \end{pmatrix} \tag{8}$$

Thus, each attention projection incorporates a learnable low-rank adapter. After training, folding the adapter into the original weight via matrix multiplication gives substantial reduction in both memory usage and computational overhead, while maintaining output quality of the uncompressed model.

3.1.3 PNF ADAPTER FOR FEED FORWARD NETWORK

Next, we present the applicaiton of PnF adapters to the FFN. Let the gate, up-projection, and down-projection at layer l be $W_{\text{gate}^{(l)}}, W_{\text{up}^{(l)}} \in \mathbb{R}^{d_{\text{embed}} \times d_{\text{inter}}}$, and $W_{\text{down}^{(l)}} \in \mathbb{R}^{d_{\text{inter}} \times d_{\text{embed}}}$, respectively. For these matrices, we introduce the corresponding PnF adapters:

$$P_{\text{gate}^{(l)}}, P_{\text{up}^{(l)}} \in \mathbb{R}^{d_{\text{inter}} \times r_{\text{inter}}^{(l)}}, \text{ and } P_{\text{down}^{(l)}} \in \mathbb{R}^{r_{\text{inter}}^{(l)} \times d_{\text{inter}}}$$
 (9)

where $r_{\text{inter}}^{(l)} < d_{\text{inter}}^{(l)}$. Multiplying these adapter with the original weights yields the compressed projections:

$$\begin{aligned} W_{\text{gate}^{(l)}} P_{\text{gate}^{(l)}} &\in \mathbb{R}^{d_{\text{embed}} \times r_{\text{inter}}^{(l)}} \\ W_{\text{up}^{(l)}} P_{\text{up}^{(l)}} &\in \mathbb{R}^{d_{\text{embed}} \times r_{\text{inter}}^{(l)}} \\ P_{\text{down}^{(l)}} W_{\text{down}^{(l)}} &\in \mathbb{R}^{r_{\text{inter}}^{(l)} \times d_{\text{embed}}} \end{aligned} \tag{10}$$

Therefore, similar to that of the attention mechanism with PnF adapters above, each FFN layer is equipped with a learnable low-rank adapter. Because the feed-forward network (FFN) comprises the majority of a transformer's parameters, folding the adapters into the original weights provides substantial savings in both memory and computation.

3.2 Training Pipeline for PnF adapter

To obtain PnF adapters with high fidelity, we propose a three-stage training pipeline: (i) **Compression Planning** that determines the per-layer degree of dimensionality reduction, (ii) **Group-wise Sequential Training** that stabilizes optimization by sequentially training a small, isolated set of adapters, and (iii) **KL-divergence Distillation Loss** that aligns the compressed model's output distribution with the original model's distribution.

Stage 1: Compression Planning Based on desired compression ratio (e.g., 20%), we first determine the degree of reduction of dimensionality (i.e., $r_{\rm head}^{(l)}$ and $r_{\rm inter}^{(l)}$) for each layer l. While the allocation of reductions can be flexible, we recommend a pyramidal schedule where deeper layers (closer to the language modeling head) are compressed more aggressively, and earlier layers receive milder reductions. Prior work on layer pruning Men et al. (2024); Gromov et al. (2024) shows that later (upper) layers can often be removed with little impact on downstream performance, indicating that they contribute less to the model's expressivity. Based on this finding, we allocate a larger portion of the compression budget to the top of the model.

Because the reduction ratio can be explicitly set, the approach is highly flexible and can be tailored to meet a user's requirements. Our empirical studies reveal that applying a higher compression rate to the FFN yields considerably better results than compressing the MHSA modules, and a concrete example of this planning is provided in the Appendix A.

Stage 2: Group-wise Sequential Training Plugging all adapters at once might perturb the original model's signal at the beginning of training, inducing covariate shift and misleading gradients. Alternatively, training a single adapter at a time preserves this signal but is prohibitively slow. To address this issue, we introduce Group-wise Sequential Training. This training scheme trains small groups of adapters in turn, retaining most of the signal preservation benefits while substantially reducing training time and stabilizing convergence, which is further discussed in Section 4.3.1. Formally, we first partition the L transformer layers into disjoint groups of size N, starting from the top of the model (output side) and moving downward. The k-th group is defined as:

$$G_k = \{L - kN + 1, \dots, L - (k-1)N\}, \quad k = 1, 2, \dots, n_g,$$
 (11)

where $n_g = \lfloor L/N \rfloor$ is the number of groups. Given the compression plan that specifies per-layer reductions (i.e., $r_{\rm head}^{(l)}$ and $r_{\rm inter}^{(l)}$), we first identify which group contain layers slated for compression. Then training proceeds sequentially from \mathcal{G}_1 towards \mathcal{G}_{n_g} .

At step k, if \mathcal{G}_k includes layers selected by the compression plan, we insert adapters only into those layers and train them, while keeping the adapters trained in previous groups $(\mathcal{G}1,\ldots,\mathcal{G}_{k-1})$ frozen. During this phase, only the parameters of current group are updated; all previous groups remain frozen with their trained adapters, while remaining groups $(\mathcal{G}_{k+1},\cdots,\mathcal{G}_{n_g})$ remain frozen without adapters (i.e., in their original state).

An instance of group-wise sequential training is illustrated in Figure 2, given L=36 and N=4, the compression plan targeting layers 13 - 36 covers six groups $(\mathcal{G}_1,\cdots,\mathcal{G}_6)$. We train these six groups sequentially from the output side toward the input (i.e., $\mathcal{G}_1\to\cdots\to\mathcal{G}_6$) while the lower 12 layers remain uncompressed. By activating one small group per step and keeping the remaining group fixed, this approach preserves the backbone signal and improves optimization stability.

Stage 3: KL-divergence Distillation Loss During the group-wise sequential training for the adapters, we adopt a Kullback-Leibler (KL) divergence loss. Specifically, the logits of the PnF-plugged model are aligned with those of the frozen backbone model by minimizing:

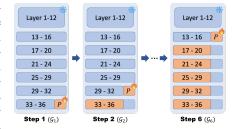


Figure 2: Visualization of group-wise sequential training. Training proceeds group by group, beginning with the output side. At any step, only the current group G_i is updated while all other groups stay frozen, which preserves the backbone signal and enhances optimization stability.

$$\mathcal{L}_{KL} = KL(p_W||p_{WP}) \tag{12}$$

where p_W and p_{WP} denote the predictive distribution of the backbone and the PnF-plugged models, respectively.

We adopt a KL-divergence distillation loss for two reasons. First, the goal of compression is to produce a smaller model that reproduces the original model's behavior. The KL-divergence can achieve this by aligning the predictive distribution of student (PnF-plugged model) with the teacher (original model). Second, recent studies (Bercovich et al., 2024; Muralidharan et al., 2024; Li et al., 2024a) report that KL-based distillation often outperforms cross-entropy, yielding better downstream performance.

3.3 Deployment for Inference

After the adapters are fully trained leveraging unhindered pre-trained weights, they can be seam-lessly integrated into the backbone model. In MHSA, for example, each adapter is folded into its corresponding pre-trained weight matrix via matrix multiplication:

$$\begin{split} W_{Q_{i}^{(l)}}P_{Q_{i}^{(l)}} &\to W_{Q_{i}^{(l)}}^{\text{Comp}} \\ W_{K_{i}^{(l)}}P_{K_{i}^{(l)}} &\to W_{K_{i}^{(l)}}^{\text{Comp}} \\ W_{V_{i}^{(l)}}P_{V_{i}^{(l)}} &\to W_{V_{i}^{(l)}}^{\text{Comp}} \\ \end{split} \tag{13}$$

A similar folding procedure applies to FFN, where each adapter is integrated into its corresponding weight matrix:

$$\begin{split} W_{\text{gate}^{(l)}} P_{\text{gate}^{(l)}} &\to W_{\text{gate}^{(l)}}^{\text{Comp}} \\ W_{\text{up}^{(l)}} P_{\text{up}^{(l)}} &\to W_{\text{up}^{(l)}}^{\text{Comp}} \\ P_{\text{down}^{(l)}} W_{\text{down}^{(l)}} &\to W_{\text{down}^{(l)}}^{\text{Comp}} \end{split} \tag{14}$$

The resulting weights directly replace the original model, reducing parameter counts and computational costs while preserving the model's architectural structure and inference pipeline. This fold-in operation has two key benefits. First, deployment is simple: the trained PnF adapters are folded into the original weights via plain matrix multiplications—no auxiliary metrics, graph edits, or specialized operators. Second, it ensures that the deployed model remains identical structure and interface to the original model, which facilitates compatibility with existing serving frameworks and hardware accelerators.

4 EXPERIMENTS

In this section, we first evaluate the PnF Compression method against several widely-used compression methods across different compression rates and original model sizes, demonstrating its

Table 1: Performance of the various compression methods on Qwen-3-8B-Base. The pretrained backbone model and its compressed variants are evaluated across multiple benchmarks at several compression rates. The best and second-best results at each compression rate are highlighted with **boldface** and <u>underline</u>, respectively.

Method	CR	PIQA	HS	WG	CSQA	ARC-e	ARC-c	OBQA	boolq	SIQA	mmlu	ld	Avg	RP
Baseline	0% (8.19B)	0.793	0.786	0.724	0.860	0.801	0.573	0.410	0.830	0.547	0.747	0.709	0.707	1.000
	20% (6.52B)	0.716	0.617	0.665	0.195	0.644	0.401	0.376	0.749	0.418	0.247	0.571	0.509	0.720
Slice GPT	30% (5.71B)	0.667	0.544	0.624	0.199	0.511	0.317	0.362	0.601	0.404	0.231	0.505	0.451	0.638
	40% (4.91B)	0.618	0.447	0.586	0.194	0.405	0.263	0.332	0.523	0.392	0.230	0.422	0.401	0.567
	20% (6.65B)	0.733	0.645	0.658	0.627	0.665	0.422	0.382	0.673	0.453	0.560	0.587	0.582	0.824
LaCo	30% (5.88B)	0.687	0.524	0.589	0.405	0.561	0.337	0.320	0.722	0.425	0.362	0.522	0.496	0.701
	40% (5.10B)	0.614	0.398	0.554	0.205	0.423	0.277	0.292	0.501	0.387	0.242	0.305	0.382	0.540
	20% (6.65B)	0.757	0.612	0.559	0.211	0.647	0.375	0.400	0.618	0.441	0.255	0.508	0.489	0.692
LLM-Streamline	30% (5.88B)	0.717	0.501	0.534	0.192	0.524	0.303	0.348	0.617	0.393	0.229	0.358	0.429	0.606
	40% (5.10B)	0.589	0.362	0.571	0.196	0.356	0.264	0.286	0.430	0.376	0.230	0.017	0.334	0.473
	20% (6.65B)	0.632	0.362	0.513	0.195	0.439	0.261	0.300	0.553	0.368	0.247	0.070	0.358	0.506
Short GPT	30% (5.88B)	0.608	0.326	0.507	0.187	0.416	0.238	0.286	0.462	0.356	0.231	0.059	0.334	0.473
	40% (5.10B)	0.572	0.287	0.526	0.185	0.367	0.214	0.262	0.440	0.347	0.229	0.021	0.314	0.444
	20% (6.55B)	0.774	0.714	0.709	0.757	0.773	0.479	0.410	0.818	0.521	0.645	0.677	0.661	0.935
Ours	30% (5.74B)	0.749	0.651	0.658	0.553	0.687	0.412	0.372	0.776	0.483	0.501	0.629	0.588	0.832
	40% (4.91B)	0.719	0.587	0.626	0.476	0.655	0.378	0.358	0.749	0.427	0.398	0.538	0.545	0.771

effectiveness (Section 4.2). We then examine the impact of our weight-preserving mechanism and training strategies through an ablation study (Section 4.3).

4.1 EXPERIMENTAL SETUP

All experiments were conducted to systematically compare the effectiveness of various large language model (LLM) compression techniques across a suite of widely-used benchmark tasks. We evaluated each method Slice-GPT (Ashkboos et al., 2024), LaCo (Yang et al., 2024), ShortGPT (Men et al., 2024), LLM-Streamline (Chen et al., 2025), and our proposed method in three target compression rates (approximately 20%, 30%, and 40%) relative to the original model size. The baselines consist of the uncompressed models: Qwen3-4B-Base, Qwen3-8B-Base, OPT 2.7B, and OPT 6.7B.

The evaluation benchmarks include: PIQA (physical commonsense reasoning), HellaSwag (commonsense inference), WinoGrande (pronoun resolution), CSQA (commonsense QA), ARC-e/ARC-c (science questions), OpenBookQA, BoolQ (boolean QA), Social IQA (multiple-choice), MMLU (multi-task language understanding), and Lambda OpenAI (factual QA). Each model's performance is measured using task-specific accuracy, or accuracy norm if available, reported per dataset. For each compression approach and setting, we tabulate the compression rate (CR) and all benchmark scores, along with the average performance (AVG) across tasks and relative performance rate (RP).

For a fair comparison, all compressed models underwent a performance recovery phase following the respective compression procedure. Specifically, our approach utilizes adapter training for post-compression recovery; the Streamline baseline employs light layer training; and other methods adopt LoRA (Hu et al., 2022) training as their recovery protocol. All recovery procedures leveraged the SlimPajama dataset (Soboleva et al., 2023), sampling 600,000 training instances, each with a sequence length of 1,024 tokens, to ensure consistency and robustness in recovered performance across all benchmarks. Comprehensive implementation and experimental details are provided in Appendix A.

4.2 RESULTS

We evaluate the proposed compression method on two base LLMs, Qwen3-4B-Base and Qwen3-8B-Base, under compression rates of approximately 20%, 30%, and 40% relative to their original parameter counts. All models were assessed in a zero-shot setting using the LLM evaluation library (Gao et al., 2024a). Additional experiments, including evaluations on other LLM variants and in five-shot settings, are reported in Appendix B.

Tables 2 and 1 summarize the results on compressing Qwen3-4B-Base and Qwen3-8B-Base, respectively. Across all compression rates, our method consistently outperforms competing approaches on most benchmarks, while preserving performance close to that of the uncompressed models. The advantage is most evident on knowledge-intensive tasks such as CSQA, MMLU, and ARC, which rely heavily on retrieving and applying pretrained knowledge. On benchmarks emphasizing common-

Table 2: Performance of the different compression methods on Qwen3-4B-Base. The pretrained backbone and its compressed variants are evaluated on the same set of benchmarks and compression rates as in Table 1. For each compression rate, the best result is shown in **boldface** and the second-best in <u>underlined</u> text.

Method	CR	PIQA	HS	WG	CSQA	ARC-e	ARC-c	OBQA	boolq	SIQA	mmlu	ld	Avg	RP
Baseline	0% (4.02B)	0.779	0.736	0.703	0.827	0.760	0.516	0.412	0.830	0.502	0.713	0.690	0.679	1.000
	20% (3.53B)	0.688	0.554	0.628	0.197	0.546	0.346	0.338	0.723	0.411	0.236	0.528	0.472	0.696
Slice GPT	30% (3.06B)	0.633	0.462	0.599	0.193	0.431	0.260	0.308	0.680	0.386	0.231	0.441	0.420	0.619
	40% (2.65B)	0.584	0.384	0.553	0.197	0.348	0.251	0.276	0.602	0.371	0.230	0.359	0.378	0.556
	20% (3.22B)	0.715	0.578	0.631	0.586	0.634	0.387	0.358	0.738	0.434	0.584	0.502	0.559	0.823
LaCo	30% (2.81B)	0.644	0.470	0.589	0.306	0.517	0.317	0.282	0.651	0.404	0.335	0.359	0.443	0.653
	40% (2.41B)	0.630	0.416	0.562	0.195	0.453	0.273	0.284	0.606	0.388	0.234	0.341	0.398	0.587
	20% (3.22B)	0.739	0.559	0.556	0.196	0.619	0.369	0.378	0.558	0.417	0.235	0.448	0.461	0.679
LLM-Streamline	30% (2.81B)	0.678	0.443	0.530	0.195	0.498	0.272	0.336	0.586	0.395	0.229	0.330	0.408	0.601
	40% (2.41B)	0.581	0.351	0.556	0.196	0.352	0.274	0.290	0.426	0.378	0.230	0.006	0.331	0.488
	20% (3.22B)	0.694	0.557	0.589	0.561	0.645	0.411	0.344	0.684	0.417	0.487	0.529	0.538	0.792
Short GPT	30% (2.81B)	0.654	0.386	0.551	0.185	0.492	0.308	0.312	0.588	0.372	0.245	0.253	0.395	0.582
	40% (2.41B)	0.548	0.274	0.519	0.222	0.319	0.226	0.238	0.538	0.350	0.244	0.029	0.319	0.469
	20% (3.22B)	0.736	0.662	0.669	0.779	0.704	0.436	0.382	0.784	0.501	0.657	0.651	0.633	0.932
Ours	30% (2.82B)	0.712	0.588	0.618	0.628	0.665	0.380	0.362	0.749	0.464	0.524	0.595	0.571	0.842
	40% (2.41B)	0.702	0.513	0.587	0.420	0.552	0.310	0.342	0.685	0.421	0.395	0.542	0.497	0.732

sense reasoning and general language understanding (e.g., HellaSwag, WinoGrande), the performance gap between methods is smaller, yet our approach still achieves the best overall balance across tasks.

When comparing Qwen3-4B-Base and Qwen3-8B-Base, we observe that the larger base model retains higher absolute accuracy across all compression methods and rates, reflecting its greater capacity. However, the relative performance preservation (RP) of our method remains consistently strong for both model scales, demonstrating its robustness. Notably, the 8B model shows slightly smaller performance degradation under compression, suggesting that larger models may provide more redundancy that can be better exploited during parameter reduction. This trend highlights that while scaling up improves baseline performance, an effective compression strategy is crucial. Overall, our method achieves stable gains across both model sizes, indicating strong generalizability of the approach.

Discussion. These findings suggest that updating adapter weights while preserving core model parameters is critical for effective LLM compression. Retaining the pretrained weight structure allows the compressed models to maintain essential knowledge and reasoning capabilities needed for complex tasks. In contrast, methods that aggressively modify core parameters tend to incur larger performance degradation, particularly on knowledge-demanding benchmarks.

4.3 ABLATIONS

4.3.1 Training Strategy

To understand how the size of the adapter groups influences effectiveness and efficiency, we performed an ablation study in which the group size N was varied while keeping all other hyper parameters, compression plan, learning rate schedule, and total training epochs identical to the default configuration described in Appendix A. The experiments, summarized in Table 3, were conducted on Qwen-3-4B-Base compressed to a 20% reduction rate.

When N=36 every adapter is inserted and trained at once, which minimizes the number of training phases but perturbs the entire backbone simultaneously. This large covariate shift leads to unstable gradients and a noticeable drop in downstream performance, as reflected by an average score of 0.6182. At the opposite extreme, N=1 updates one adapter at a time, moving sequentially through the 36 layers. Because only a single component is altered during each step, the original signal is largely preserved, resulting in the highest average performance. However, the training iteration grows roughly linearly with the number of groups, making this setting impractical for larger models.

Our default configuration adopts ${\cal N}=4$, grouping four consecutive layers together. This approach retains most of the stability advantages

Table 3: Ablation of the group size N used in the group wise sequential training scheme. The table reports the average downstream score.

Group size	Avg
N=36 (all)	0.6182
N=1	0.6346
N=4 (ours)	0.6329

of the single-adapter regime while dramatically reducing the total number of training phases. The

resulting average score (0.6329) is only marginally below the optimal N=1 setting, yet the computational cost is comparable to the "all-at-once" baseline. Consequently, we select N=4 as the standard group size for all subsequent experiments.

4.3.2 IMPACT OF RECOVERY-TRAINING SET SIZE

Table 4: Effect of recovery-training set size on the performance of our 20% compressed Qwen-3-4B-Base. Results are reported for four different sample budgets (300k, 600k, 1M, and 2M) on a range of downstream benchmarks.

Method	CR	Samples	PIQA	HS	WG	CSQA	ARC-e	ARC-c	OBQA	boolq	SIQA	mmlu	ld	Avg	RP
Baseline	-	-	0.7786	0.7364	0.7032	0.8272	0.7597	0.5162	0.4120	0.8299	0.5015	0.7131	0.6898	0.6789	1.0000
		300K	0.7163	0.6433	0.6630	0.7802	0.7046	0.4181	0.3720	0.7976	0.4928	0.6476	0.6418	0.6252	0.9209
O	20%	600K	0.7363	0.6622	0.6690	0.7790	0.7044	0.4358	0.3820	0.7837	0.5013	0.6573	0.6514	0.6329	0.9322
Ours	20%	1M	0.7350	0.6757	0.6788	0.8354	0.7022	0.4488	0.3720	0.7985	0.4923	0.7084	0.6693	0.6469	0.9528
		2M	0.7679	0.7230	0.6890	0.8215	0.7513	0.5060	0.4020	0.8315	0.4908	0.7076	0.6804	0.6701	0.9870

In this section we evaluate how the size of the recovery-training set influences the effectiveness of our compression pipeline. Table 4 reports results for four different sample budgets (300K, 600K, 1M, and 2M) under a fixed compression rate of 20%. As the number of training instances grows, downstream performance improves consistently across virtually all benchmarks: the average score rises from 0.6252 (300K samples) to 0.6701 (2M samples), and the relative performance (RP) climbs from 0.9209 to 0.9870, narrowing the gap with the uncompressed baseline (Avg=0.6789). For most tasks the improvement is gradual, but a few—namely HS, BoolQ, OBQA, and ARC—show a different pattern. With only 300K–1M samples their scores increase only marginally, reflecting the limited signal provided by a small recovery set. Once the sample count reaches over 1M, the gains accelerate sharply; at 2M samples these tasks almost match the baseline performance (e.g., HS jumps from 0.6757 to 0.7230, BoolQ from 0.7985 to 0.8315, OBQA from 0.3720 to 0.4020, ARC-e from 0.7022 to 0.7513). This behavior suggests that certain evaluation sets require a richer recovery signal before the compressed model can fully exploit the knowledge retained in the frozen backbone.

Overall, the results show that even a modest recovery set captures more than 90% of the attainable relative performance (RP). When the recovery data are scaled to a few million examples, the compressed model nearly matches the uncompressed baseline, incurring less than a 2% performance drop while preserving the 20% compression ratio.

5 LIMITATION

A possible limitation of our approach is that the first stage of the pipeline is deliberately empirical. Although this stage grants users freedom to design compression plans, it also places a burden on the practitioner to possess a priori knowledge about the model's relative importance of its components for the tasks of interest. In practice, an uninformed choice of reduction rates can lead to sub-optimal performance or unnecessary training overhead. On the other hand, this very flexibility makes the stage a useful diagnostic tool: by systematically varying the groups that are compressed, users can probe which parts of an LLM are most critical for specific linguistic or reasoning abilities. Future work could therefore focus on automated or data-driven heuristics (e.g., sensitivity analyses, reinforcement-learning controllers) that suggest compression configurations with minimal user intervention, while still preserving the analytical benefits of the current empirical design.

6 Conclusion

In this work, we introduce a novel framework Plug-and-Fold (PnF), a compression framework that preserves both weights and structure of the pretrained LLM. In our workflow, lightweight PnF adapters are first plugged into a pretrained LLM's weight matrices. After going through adaption phase, adapters are folded back into the base model via simple matrix multiplication. The resulting model is structurally identical to the original backbone yet enjoys substantial reductions in parameters with unimpaired performance. Extensive experiments on four backbones and three compression rates show PnF consistently outperforms strong baselines, highlighting the benefit of retaining pretrained weights. Ablation studies on training strategies confirm the effectiveness of our workflow, while experiments on recovery-training set size demonstrate that with sufficient data PnF can nearly match the original model's performance. In summary, Plug-and-Fold provides an efficient, scalable, architecture-preserving compression pipeline that maintains the expressive power of large pretrained LLMs, enabling deployment on resource-constrained hardware without performance loss.

REFERENCES

- Rishabh Agarwal, Nino Vieillard, Piotr Stanczyk, Sabela Ramos, Matthieu Geist, and Olivier Bachem. Gkd: Generalized knowledge distillation for auto-regressive sequence models. <u>CoRR</u>, 2023.
- Yongqi An, Xu Zhao, Tao Yu, Ming Tang, and Jinqiao Wang. Systematic outliers in large language models. arXiv preprint arXiv:2502.06415, 2025.
- Saleh Ashkboos, Maximilian L Croci, Marcelo Gennari do Nascimento, Torsten Hoefler, and James Hensman. Slicegpt: Compress large language models by deleting rows and columns. <u>arXiv:2401.15024</u>, 2024.
- Akhiad Bercovich, Tomer Ronen, Talor Abramovich, Nir Ailon, Nave Assaf, Mohammad Dabbah, Ido Galil, Amnon Geifman, Yonatan Geifman, Izhak Golan, et al. Puzzle: Distillation-based nas for inference-optimized llms. arXiv preprint arXiv:2411.19146, 2024.
- Zhaodong Bing, Linze Li, and Jiajun Liang. Optimizing knowledge distillation in transformers: Enabling multi-head attention without alignment barriers. arXiv preprint arXiv:2502.07436, 2025.
- Xiaodong Chen, Yuxuan Hu, Jing Zhang, Yanling Wang, Cuiping Li, and Hong Chen. Streamlining redundant layers to compress large language models. In The Thirteenth International Conference on Learning Representations, 2025. URL https://openreview.net/forum?id=IC5RJvRoMp.
- Xiao Cui, Mo Zhu, Yulei Qin, Liang Xie, Wengang Zhou, and Houqiang Li. Multi-level optimal transport for universal cross-tokenizer knowledge distillation on language models. In <u>Proceedings</u> of the AAAI Conference on Artificial Intelligence, volume 39, pp. 23724–23732, 2025.
- Stefan Elfwing, Eiji Uchibe, and Kenji Doya. Sigmoid-weighted linear units for neural network function approximation in reinforcement learning. Neural networks, 107:3–11, 2018.
- 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. The language model evaluation harness, 07 2024a. URL https://zenodo.org/records/12608602.
- Shangqian Gao, Chi-Heng Lin, Ting Hua, Zheng Tang, Yilin Shen, Hongxia Jin, and Yen-Chang Hsu. Disp-llm: Dimension-independent structural pruning for large language models. <u>Advances in Neural Information Processing Systems</u>, 37:72219–72244, 2024b.
- Andrey Gromov, Kushal Tirumala, Hassan Shapourian, Paolo Glorioso, and Daniel A Roberts. The unreasonable ineffectiveness of the deeper layers. arXiv preprint arXiv:2403.17887, 2024.
- Raby Hamadi. Large language models meet computer vision: A brief survey. <u>arXiv preprint</u> arXiv:2311.16673, 2023.
- G Hinton. Distilling the knowledge in a neural network. In <u>Deep Learning and Representation</u> Learning Workshop in Conjunction with NIPS, 2014.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. <u>ICLR</u>, 1(2):3, 2022.
- Xing Hu, Yuan Cheng, Dawei Yang, Zukang Xu, Zhihang Yuan, Jiangyong Yu, Chen Xu, Zhe Jiang, and Sifan Zhou. Ostquant: Refining large language model quantization with orthogonal and scaling transformations for better distribution fitting. arXiv preprint arXiv:2501.13987, 2025.
- Raisa Islam and Owana Marzia Moushi. Gpt-40: The cutting-edge advancement in multimodal llm. In Intelligent Computing-Proceedings of the Computing Conference, pp. 47–60. Springer, 2025.

- Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, Lélio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mixtral of experts, 2024a. URL https://arxiv.org/abs/2401.04088.
 - Juyong Jiang, Fan Wang, Jiasi Shen, Sungju Kim, and Sunghun Kim. A survey on large language models for code generation. arXiv preprint arXiv:2406.00515, 2024b.
 - Pierre-Carl Langlais, Carlos Rosas Hinostroza, Mattia Nee, Catherine Arnett, Pavel Chizhov, Eliot Krzystof Jones, Irène Girard, David Mach, Anastasia Stasenko, and Ivan P Yamshchikov. Common corpus: The largest collection of ethical data for llm pre-training. <u>arXiv preprint</u> arXiv:2506.01732, 2025.
 - Shengrui Li, Junzhe Chen, Xueting Han, and Jing Bai. Nuteprune: Efficient progressive pruning with numerous teachers for large language models. arXiv preprint arXiv:2402.09773, 2024a.
 - Shiyao Li, Xuefei Ning, Luning Wang, Tengxuan Liu, Xiangsheng Shi, Shengen Yan, Guohao Dai, Huazhong Yang, and Yu Wang. Evaluating quantized large language models. In <u>Proceedings of the 41st International Conference on Machine Learning</u>, pp. 28480–28524, 2024b.
 - Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin, Robert Verkuil, Ori Kabeli, Yaniv Shmueli, et al. Evolutionary-scale prediction of atomic-level protein structure with a language model. Science, 379(6637):1123–1130, 2023.
 - Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. <u>Advances in neural information processing systems</u>, 36:34892–34916, 2023a.
 - Jing Liu, Ruihao Gong, Xiuying Wei, Zhiwei Dong, Jianfei Cai, and Bohan Zhuang. Qllm: Accurate and efficient low-bitwidth quantization for large language models. arXiv:2310.08041, 2023b.
 - Yang Liu, Jiahuan Cao, Chongyu Liu, Kai Ding, and Lianwen Jin. Datasets for large language models: A comprehensive survey. arXiv preprint arXiv:2402.18041, 2024.
 - Xinyin Ma, Gongfan Fang, and Xinchao Wang. Llm-pruner: On the structural pruning of large language models. <u>Advances in neural information processing systems</u>, 36:21702–21720, 2023.
 - Xin Men, Mingyu Xu, Qingyu Zhang, Bingning Wang, Hongyu Lin, Yaojie Lu, Xianpei Han, and Weipeng Chen. Shortgpt: Layers in large language models are more redundant than you expect. arXiv preprint arXiv:2403.03853, 2024.
 - Leandro Giusti Mugnaini, Bruno Lopes Yamamoto, Lucas Lauton de Alcantara, Victor Zacarias, Edson Bollis, Lucas Pellicer, Anna Helena Reali Costa, and Artur Jordao. Efficient llms with amp: Attention heads and mlp pruning. arXiv preprint arXiv:2504.21174, 2025.
 - Saurav Muralidharan, Sharath Turuvekere Sreenivas, Raviraj Joshi, Marcin Chochowski, Mostofa Patwary, Mohammad Shoeybi, Bryan Catanzaro, Jan Kautz, and Pavlo Molchanov. Compact language models via pruning and knowledge distillation. Advances in Neural Information Processing Systems, 37:41076–41102, 2024.
 - Utkarsh Ojha, Yuheng Li, Anirudh Sundara Rajan, Yingyu Liang, and Yong Jae Lee. What knowledge gets distilled in knowledge distillation? <u>Advances in Neural Information Processing Systems</u>, 36:11037–11048, 2023.
- Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. 2018.
 - Prajit Ramachandran, Barret Zoph, and Quoc V Le. Searching for activation functions. In International Conference on Learning Representations, 2018.
 - Yuzhang Shang, Zhihang Yuan, Qiang Wu, and Zhen Dong. Pb-llm: Partially binarized large language models. arXiv preprint arXiv:2310.00034, 2023.

602

603

604

605

606 607

608

609 610

612

613

614

615

616617

618

619

620

621 622

623

624

625

626

627

644

- Noam Shazeer. Glu variants improve transformer. <u>arXiv preprint arXiv:2002.05202</u>, 2020.
- 596 Daria Soboleva, Faisal Al-Khateeb, Robert Myers, Jacob R Steeves, Joel Hestness. and Nolan Dey. SlimPajama: A 627B token cleaned and dedu-597 plicated version of RedPajama. https://www.cerebras.net/blog/ 598 slimpajama-a-627b-token-cleaned-and-deduplicated-version-of-redpajama, 2023. URL https://huggingface.co/datasets/cerebras/SlimPajama-627B. 600
 - Mingjie Sun, Zhuang Liu, Anna Bair, and J Zico Kolter. A simple and effective pruning approach for large language models. In The Twelfth International Conference on Learning Representations.
 - Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. arXiv preprint arXiv:2403.05530, 2024.
 - Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijue Chen, Yanru Chen, Yuankun Chen, Yutian Chen, et al. Kimi k2: Open agentic intelligence. arXiv e-prints, pp. arXiv-2507, 2025.
 - Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. <u>Advances in neural information</u> processing systems, 30, 2017.
 - Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov. Analyzing multi-head self-attention: Specialized heads do the heavy lifting, the rest can be pruned. arXiv:1905.09418, 2019.
 - Yifei Yang, Zouying Cao, and Hai Zhao. Laco: Large language model pruning via layer collapse. arXiv preprint arXiv:2402.11187, 2024.
 - Biao Zhang and Rico Sennrich. Root mean square layer normalization. <u>Advances in neural</u> information processing systems, 32, 2019.
 - Qiang Zhang, Keyan Ding, Tianwen Lv, Xinda Wang, Qingyu Yin, Yiwen Zhang, Jing Yu, Yuhao Wang, Xiaotong Li, Zhuoyi Xiang, et al. Scientific large language models: A survey on biological & chemical domains. ACM Computing Surveys, 57(6):1–38, 2025.
 - Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. <u>arXiv</u> preprint arXiv:2303.18223, 1(2), 2023.

A EXPERIMENT SETTINGS

A.1 HYPER-PARAMETER CONFIGURATION

In all experiments we follow the two-stage pipeline described in Section 3.2. Below we detail the hyperparameter settings that were used to instantiate the compression plan, to construct the training groups, and to train the adapters. The values are the same for every model and compression rate unless explicitly noted. Also, the PnF are initialized as identity matrix, where only the diagonal elements are set to 1 otherwise 0.

A.2 COMPRESSION PLAN (PER-LAYER REDUCTION RATES)

For each target compression rate $c \in \{20\%, 30\%, 40\%\}$, we empirically driven target hidden-dimension targets for the multi-head self-attention (MHSA) and feed-forward network (FFN) sub-layers. The resulting dimensionalities are listed in Table 5. The notation indicates the target hidden size for each group in the order in which the groups are visited (from the output side toward the input side).

Table 5: Target hidden dimensions for MHSA and FFN at each compression rate. Each entry corresponds to a successive group of layers (see Figure 2).

Backbone	CR	MHSA	FFN
	20%	(72,72,-,-)	(3584, 3584, 4096, 4864)
OPT 2.7B	30%	(64, 72, 72, -, -, -, -)	(3328, 3328, 3840, 4608, 5632, 6144, 8192)
	40%	(64, 64, 72, 72, -, -, -, -)	(2560, 2560, 2816, 2816, 2816, 3840, 5888, 8096)
	20%	(-, -, -, -)	(5120, 5632, 7168, 7168)
OPT 6.7B	30%	(64, 80, 96, 112, -, -)	(4608, 5376, 6144, 8192, 10240, 13312)
	40%	(64, 64, 64, 64, 96, -)	(5632, 5376, 5120, 5120, 7168, 7168)
	20%	(-,-,-,-,-)	(2560, 2816, 3328, 4608, 9216)
Qwen3 4B	30%	(-,-,-,-,-)	(2560, 2560, 2560, 3072, 3584, 4864)
	40%	(-,-,-,-,-,-,-)	(2560, 2560, 2560, 2816, 2816, 3072, 3328, 5632)
	20%	(-,-,-,-)	(4096, 4352, 4864, 6144, 8704)
Qwen3 8B	30%	(-,-,-,-,-,-)	(4096, 4352, 4608, 4864, 4864, 5632, 7680)
	40%	(-,-,-,-,-,-,-,-)	(4096, 4352, 4608, 4608, 4352, 4608, 4608, 4608, 7936)

Interpretation of Table 5 Taking OPT 2.7B as an example, for a 20% reduction the first two groups (closest to the output) compress both the MHSA projection matrices to $r_{\rm head}=72$ and the FFN intermediate dimensions to $r_{\rm inter}=3584$. Subsequent groups use the next values in the list, while "-" denotes it retains the original dimension. At 30% and 40% the plan contains more groups, thereby spreading the reduction more gradually across the stack.

A.3 TRAINING SCHEDULE

The overall workflow of training is as follows. For each selected group G_k we:

- 1. Insert PnF adapters corresponding to index belonging to G_k
- 2. Train for E epochs while keeping all previously trained groups frozen
- 3. Proceed to G_{k+1} until G_{n_q}

Through out the entire experiments, the number of epochs is fixed to E := 1, giving a total of n_g iteration.

A.4 FOLDING STEP

After the final group has been trained, each adapter pair is merged into its corresponding projection matrix W by the closed-form multiplication. No additional fine-tuning is performed after folding,

which guarantees that the resulting model has exactly the same architecture and runtime characteristics as the original uncompressed model.

A.5 BASELINE RECOVERY FINE-TUNING SETTINGS

For the recovery-fine-tuning (RFT) stage we adopt LoRA, since LoRA fine-tuning is widely used in recent work. To ensure a fair comparison, we fix the low-rank dimension to r=16 for every LoRA experiment. Unless a particular method explicitly restricts its scope, LoRA is applied to all transformer layers—both the multi-head self-attention (MHSA) and feed-forward network (FFN) sub-layers.

B ADDITIONAL RESULTS

B.1 COMPARISON WITH BASELINE METHODS

In this section we compare our proposed approach with several baselines across a broader set of conditions. We evaluate four backbone models—Qwen-3-4B-Base, Qwen-3-8B-Base, OPT-2.7B, and OPT-6.7B—and we assess performance in both zero-shot and five-shot settings. Across all experiments, our method consistently yields the highest average score (Avg), closely matching the performance of the uncompressed baseline for each backbone.

The same trend observed in the zero-shot experiments holds in the five-shot setting. Our compression method consistently outperforms the baselines across all compression rates, and the performance gap widens on knowledge-intensive benchmarks. Thus, the superior performance of our approach is preserved when a few exemplars are provided.

Table 6: Performance of the different compression methods on Qwen3-4B-Base on five-shots setting.

Method	CR	Sample	PIQA	HS	WG	CSQA	ARC-e	ARC-c	OBQA	boolq	SIQA	mmlu	ld	Avg	RP
Baseline	0% (4.02B)	-	0.7889	0.7532	0.7206	0.8198	0.8674	0.6425	0.4500	0.8654	0.5502	0.7319	0.6501	0.7127	1.0000
	20% (3.53B)		0.6980	0.5612	0.6425	0.3030	0.6902	0.4130	0.3480	0.7746	0.4641	0.3250	0.4487	0.5153	0.7230
Slice GPT	30% (3.06B)	600k	0.6409	0.4661	0.6085	0.2293	0.5370	0.2952	0.3120	0.6911	0.4181	0.2651	0.3656	0.4390	0.6160
	40% (2.65B)		0.5832	0.3857	0.5596	0.1925	0.4158	0.2440	0.2780	0.511	0.3909	0.2672	0.2928	0.3746	0.5256
	20% (3.22B)		0.7236	0.5840	0.6425	0.7273	0.7016	0.4249	0.3680	0.7679	0.4698	0.6192	0.4496	0.5889	0.8264
LaCo	30% (2.81B)	600k	0.6398	0.475	0.5841	0.3194	0.5556	0.3362	0.2820	0.7028	0.4252	0.2863	0.3043	0.4464	0.6264
	40% (2.41B)		0.6300	0.4136	0.5509	0.2080	0.4996	0.2944	0.2880	0.6242	0.4083	0.2810	0.2550	0.4048	0.5680
	20% (3.22B)		0.7448	0.5572	0.5241	0.2015	0.7428	0.4292	0.3880	0.5474	0.4544	0.2895	0.3974	0.4797	0.6730
LLM-Streamline	30% (2.81B)	600k	0.6724	0.4333	0.5059	0.1891	0.5883	0.3054	0.3180	0.6012	0.4027	0.2538	0.3049	0.4159	0.5836
	40% (2.41B)		0.5865	0.3468	0.5643	0.1957	0.3742	0.2611	0.2800	0.3841	0.3602	0.2295	0.0060	0.3262	0.4577
	20% (3.22B)		0.7008	0.5520	0.6014	0.5766	0.7189	0.4573	0.3280	0.6914	0.4631	0.5167	0.4644	0.5519	0.7743
Short GPT	30% (2.81B)	600k	0.6088	0.3142	0.5138	0.1974	0.4196	0.2747	0.2480	0.3847	0.3561	0.2446	0.0134	0.3250	0.4561
	40% (2.41B)		0.5294	0.2564	0.4972	0.2080	0.2950	0.2568	0.2460	0.3869	0.3439	0.2370	0.0000	0.2961	0.4154
	20% (3.22B)		0.7559	0.6714	0.6772	0.8003	0.7739	0.4955	0.4100	0.8355	0.5417	0.6771	0.6055	0.6585	0.9240
Ours	30% (2.82B)	600k	0.7233	0.5847	0.6343	0.6798	0.7070	0.4008	0.4000	0.7602	0.4955	0.5412	0.5411	0.5880	0.8250
	40% (2.41B)		0.6912	0.5134	0.5783	0.5030	0.6186	0.3487	0.3540	0.7283	0.4517	0.3956	0.4757	0.5144	0.7218

Table 7: Performance of the different compression methods on Opt 6.7B in zero-shot setting.

Method	CR	Sample	PIQA	HS	WG	CSQA	ARC-e	ARC-c	OBQA	boolq	SIQA	mmlu	ld	Avg	RP
Baseline	0% (6.66B)	-	0.7644	0.6719	0.6543	0.2031	0.6002	0.3473	0.3760	0.6612	0.4278	0.2505	0.6769	0.5121	1.0000
	20% (5.49B)		0.7165	0.5657	0.6204	0.1916	0.5055	0.2961	0.3560	0.6235	0.4206	0.2500	0.5632	0.4645	0.9070
Slice GPT	30% (4.77B)	600k	0.7013	0.5220	0.6093	0.1957	0.4735	0.2875	0.3320	0.6064	0.3976	0.2421	0.4890	0.4415	0.8621
	40% (4.07B)		0.6589	0.4709	0.5604	0.1982	0.4495	0.2671	0.3280	0.5835	0.3899	0.2290	0.4017	0.4125	0.8054
	20% (5.25B)		0.6866	0.5310	0.6014	0.2064	0.4899	0.2995	0.3280	0.6214	0.4165	0.2503	0.5088	0.4491	0.8769
LaCo	30% (4.64B)	600k	0.6213	0.3890	0.5446	0.1974	0.3965	0.2560	0.2980	0.6214	0.3735	0.2463	0.1764	0.3746	0.7315
	40% (4.04B)		0.5930	0.3391	0.5170	0.1957	0.3481	0.2363	0.2740	0.6211	0.3613	0.2371	0.0638	0.3442	0.6722
	20% (5.25B)		0.7361	0.6037	0.6172	0.1761	0.5745	0.3191	0.3320	0.6324	0.4165	0.2470	0.5492	0.4731	0.9238
LLM-Streamline	30% (4.64B)	600k	0.6953	0.4204	0.5588	0.1974	0.5198	0.2850	0.3260	0.6330	0.3904	0.2381	0.2791	0.4130	0.8065
	40% (4.04B)		0.6284	0.3430	0.5288	0.1966	0.4491	0.2304	0.2960	0.6217	0.3464	0.2311	0.1186	0.3627	0.7083
	20% (5.25B)		0.5044	0.2597	0.5051	0.1957	0.2668	0.2594	0.2720	0.3783	0.3515	0.2295	0.0000	0.2929	0.5720
Short GPT	30% (4.64B)	600k	0.5065	0.2578	0.4917	0.1957	0.2597	0.2568	0.2860	0.3783	0.3418	0.2295	0.0000	0.2913	0.5687
	40% (4.04B)		0.5065	0.2579	0.4878	0.1957	0.2601	0.2491	0.2980	0.3783	0.3454	0.2295	0.0000	0.2917	0.5695
	20% (5.32B)		0.7403	0.6126	0.6461	0.2146	0.5886	0.3278	0.3600	0.6666	0.4207	0.2567	0.6135	0.4952	0.9671
Ours	30% (4.66B)	600k	0.7126	0.5321	0.6127	0.1998	0.5495	0.3069	0.3340	0.6496	0.4140	0.2512	0.5269	0.4627	0.9035
.	40% (3.99B)		0.6417	0.4926	0.5920	0.1966	0.4877	0.2874	0.3260	0.6382	0.3949	0.2464	0.4728	0.4342	0.8479

Table 8: Performance of the different compression methods on Opt 2.7B in zero-shot setting.

Method	CR	Sample	PIQA	HS	WG	CSQA	ARC-e	ARC-c	OBQA	boolq	SIQA	mmlu	ld	Avg	RP
Baseline	0% (2.65B)	-	0.7481	0.6063	0.6101	0.1990	0.5438	0.3131	0.3520	0.6027	0.4212	0.2567	0.6361	0.4808	1.0000
	20% (2.23B)		0.6654	0.4682	0.5904	0.2031	0.4322	0.2637	0.3300	0.5257	0.3838	0.2415	0.4108	0.4104	0.8537
Slice GPT	30% (1.94B)	600k	0.6300	0.4228	0.5635	0.1966	0.4175	0.2585	0.3060	0.5168	0.3705	0.2316	0.3551	0.3881	0.8072
	40% (1.66B)		0.5865	0.3674	0.5343	0.1957	0.3742	0.2509	0.2820	0.3982	0.3602	0.2301	0.2880	0.3516	0.7313
	20% (2.10B)		0.6697	0.4629	0.5612	0.1957	0.4356	0.2782	0.3080	0.6223	0.3899	0.2436	0.4768	0.4222	0.8781
LaCo	30% (1.86B)	600k	0.6197	0.3677	0.5627	0.2113	0.3699	0.2415	0.2880	0.5832	0.3853	0.2330	0.1469	0.3645	0.7581
	40% (1.63B)		0.5762	0.3006	0.5193	0.1957	0.3308	0.2261	0.2920	0.5920	0.3561	0.2312	0.0279	0.3316	0.6897
	20% (2.10B)		0.7100	0.5471	0.6038	0.1974	0.5097	0.2867	0.3240	0.6058	0.4053	0.2537	0.5692	0.4557	0.9478
LLM-Streamline	30% (1.86B)	600k	0.6763	0.4016	0.5438	0.1966	0.4609	0.2585	0.3160	0.6012	0.3756	0.2344	0.2876	0.3957	0.8230
	40% (1.63B)		0.6023	0.3122	0.5114	0.1949	0.3788	0.2150	0.2760	0.6119	0.3454	0.2298	0.0778	0.3414	0.7101
	20% (2.10B)		0.6692	0.4476	0.5745	0.1941	0.4457	0.2696	0.3080	0.5929	0.3904	0.2315	0.3155	0.4035	0.8393
Short GPT	30% (1.86B)	600k	0.5354	0.2715	0.5083	0.1982	0.3081	0.2381	0.2600	0.3789	0.3459	0.2301	0.0029	0.2979	0.6197
	40% (1.63B)		0.5152	0.2677	0.5067	0.1974	0.2908	0.2500	0.2600	0.3810	0.3423	0.2315	0.0035	0.2951	0.6138
	20% (2.11B)		0.7235	0.5012	0.6088	0.2023	0.5139	0.2922	0.3460	0.6287	0.4243	0.2500	0.5666	0.4598	0.9563
Ours	30% (1.85B)	600k	0.6908	0,4615	0.5741	0.1981	0.4724	0.2782	0.3180	0.6157	0.4132	0.2462	0.5407	0.4347	0.9042
	40% (1.58B)		0.6642	0.4205	0.5449	0.1957	0.4486	0.2759	0.2940	0.5861	0.4020	0.2388	0.4584	0.4117	0.8564

Table 9: Performance of the different compression methods on Opt 6.7B in five-shot setting.

Method	CR	Sample	PIQA	HS	WG	CSQA	ARC-e	ARC-c	OBQA	boolq	SIQA	mmlu	ld	Avg	RP
Baseline	0% (6.66B)		0.7704	0.6797	0.6598	0.1867	0.6982	0.3703	0.3920	0.7012	0.4785	0.2634	0.6451	0.5314	1.0000
	20% (5.49B)		0.7187	0.5652	0.6211	0.1981	0.5984	0.3293	0.3600	0.5492	0.4206	0.2622	0.4189	0.4583	0.8625
Slice GPT	30% (4.77B)	600k	0.6921	0.5221	0.6314	0.1826	0.5699	0.3063	0.3280	0.5318	0.4124	0.2553	0.3623	0.4358	0.8202
	40% (4.07B)		0.6561	0.4669	0.5912	0.1859	0.5173	0.2790	0.3220	0.5028	0.3935	0.2666	0.2925	0.4067	0.7654
	20% (5.25B)		0.6915	0.5318	0.6069	0.2146	0.5244	0.3038	0.3280	0.6217	0.4355	0.2595	0.4935	0.4556	0.8573
LaCo	30% (4.64B)	600k	0.6170	0.3914	0.5375	0.1998	0.4411	0.2730	0.2840	0.6220	0.3817	0.2549	0.1300	0.3757	0.7069
	40% (4.04B)		0.5919	0.3399	0.5312	0.1949	0.3733	0.2406	0.2660	0.6211	0.3541	0.2542	0.0324	0.3454	0.6500
	20% (5.25B)		0.7426	0.6207	0.5943	0.2006	0.6485	0.3455	0.3700	0.6519	0.4600	0.2522	0.5356	0.4929	0.9275
LLM-Streamline	30% (4.64B)	600k	0.6219	0.3529	0.5099	0.1810	0.4428	0.2338	0.2640	0.5927	0.3572	0.2496	0.0714	0.3525	0.6633
	40% (4.04B)		0.5811	0.2982	0.4964	0.1998	0.3577	0.2167	0.2560	0.5838	0.3326	0.2433	0.0213	0.3261	0.6136
	20% (5.25B)		0.5060	0.2606	0.5233	0.1957	0.2622	0.2594	0.2680	0.3783	0.3490	0.2295	0.0000	0.2938	0.5529
Short GPT	30% (4.64B)	600k	0.4984	0.2562	0.4957	0.1957	0.2563	0.2474	0.2800	0.3783	0.3423	0.2295	0.0000	0.2891	0.5440
	40% (4.04B)		0.5054	0.2552	0.4972	0.1957	0.2546	0.2534	0.2820	0.3783	0.3464	0.2295	0.0000	0.2907	0.5470
	20% (5.32B)		0.7647	0.6255	0.6319	0.2080	0.6477	0.3423	0.3720	0.6729	0.4683	0.2610	0.6032	0.5089	0.9576
Ours	30% (4.66B)	600k	0.7323	0.5273	0.6221	0.1909	0.5905	0.3167	0.3520	0.6461	0.4468	0.2547	0.5081	0.4716	0.8874
	40% (3.99B)		0.6896	0.4673	0.6038	0.1959	0.5343	0.2819	0.3320	0.6086	0.4292	0.2501	0.3951	0.4353	0.8191

Table 10: Performance of the different compression methods on Opt 2.7B in five-shot setting.

Method	CR	Sample	PIQA	HS	WG	CSQA	ARC-e	ARC-c	OBQA	boolq	SIQA	mmlu	ld	Avg	RP
Baseline	0% (2.65B)	-	0.7481	0.6068	0.6204	0.1884	0.6469	0.3311	0.3580	0.6272	0.4550	0.2579	0.6010	0.4946	1.0000
	20% (2.23B)		0.6757	0.4632	0.5770	0.1933	0.5080	0.2918	0.3100	0.4205	0.4099	0.2457	0.3037	0.3999	0.8085
Slice GPT	30% (1.94B)	600k	0.6322	0.4179	0.5746	0.2015	0.4609	0.2551	0.3000	0.4477	0.3991	0.2538	0.2663	0.3826	0.7736
	40% (1.66B)		0.5936	0.3612	0.5383	0.2080	0.3880	0.2449	0.2800	0.4349	0.3756	0.2480	0.1974	0.3518	0.7113
	20% (2.10B)		0.6746	0.4600	0.5825	0.1925	0.4886	0.2824	0.2900	0.6217	0.4252	0.2628	0.4221	0.4275	0.8643
LaCo	30% (1.86B)	600k	0.6186	0.3690	0.5588	0.1900	0.3986	0.2491	0.2600	0.6211	0.3705	0.2465	0.1025	0.3622	0.7324
	40% (1.63B)		0.5745	0.2973	0.5130	0.2023	0.3350	0.2287	0.2600	0.6208	0.3561	0.2366	0.0155	0.3309	0.6690
	20% (2.10B)		0.7198	0.5554	0.6006	0.1990	0.5871	0.3012	0.3260	0.6000	0.4385	0.2512	0.4925	0.4610	0.9321
LLM-Streamline	30% (1.86B)	600k	0.6436	0.4228	0.5138	0.1818	0.4524	0.2627	0.2720	0.5422	0.3689	0.2570	0.2327	0.3773	0.7628
	40% (1.63B)		0.5539	0.2823	0.5075	0.1990	0.3338	0.2099	0.2500	0.5673	0.3336	0.2505	0.0165	0.3186	0.6441
•	20% (2.10B)		0.6442	0.4013	0.5604	0.1916	0.4566	0.2637	0.2980	0.5621	0.3935	0.2505	0.1970	0.3835	0.7754
Short GPT	30% (1.86B)	600k	0.5152	0.2553	0.5257	0.1966	0.2727	0.2457	0.2800	0.3783	0.3413	0.2295	0.0000	0.2946	0.5956
	40% (1.63B)		0.5011	0.2572	0.5193	0.2007	0.2685	0.2654	0.2820	0.3783	0.3413	0.2342	0.0000	0.2953	0.5970
-	20% (2.11B)		0.7107	0.5642	0.6099	0.1901	0.5835	0.3101	0.3320	0.6300	0.4302	0.2534	0.5110	0.4659	0.9420
Ours	30% (1.85B)	600k	0.6794	0.4540	0.5741	0.2015	0.5243	0.2894	0.3300	0.5701	0.4291	0.2588	0.4518	0.4330	0.8754
	40% (1.58B)		0.6518	0.4096	0.5551	0.1966	0.4827	0.2777	0.3080	0.5498	0.4230	0.2503	0.3678	0.4066	0.8220

Table 11: Performance of the different compression methods on Qwen3-8B-Base with five shots setting. The pretrained backbone and its compressed variants are evaluated on the same set of benchmarks and compression rates as in Table 1. For each compression rate, the best result is shown in **boldface** and the second-best in <u>underlined</u> text.

Method	CR	PIQA	HS	WG	CSQA	ARC-e	ARC-c	OBQA	boolq	SIQA	mmlu	ld	Avg	RP
Baseline	0% (8.19B)	0.815	0.795	0.770	0.856	0.880	0.681	0.490	0.882	0.572	0.770	0.671	0.744	1.000
	20% (6.52B)	0.714	0.632	0.686	0.329	0.747	0.462	0.396	0.781	0.496	0.356	0.527	0.557	0.749
Slice GPT	30% (5.71B)	0.676	0.553	0.642	0.275	0.621	0.361	0.370	0.696	0.443	0.275	0.456	0.488	0.656
	40% (4.91B)	0.627	0.451	0.594	0.201	0.494	0.279	0.318	0.614	0.415	0.255	0.363	0.419	0.564
	20% (6.65B)	0.736	0.651	0.671	0.709	0.748	0.493	0.406	0.534	0.503	0.604	0.546	0.600	0.807
LaCo	30% (5.88B)	0.694	0.535	0.600	0.506	0.629	0.358	0.318	0.673	0.456	0.408	0.471	0.514	0.690
	40% (5.10B)	0.617	0.403	0.572	0.215	0.487	0.297	0.276	0.623	0.402	0.251	0.256	0.400	0.538
	20% (6.65B)	0.774	0.613	0.561	0.238	0.769	0.446	0.402	0.548	0.477	0.268	0.462	0.505	0.680
LLM-Streamline	30% (5.88B)	0.724	0.500	0.553	0.194	0.673	0.338	0.346	0.450	0.418	0.243	0.310	0.432	0.580
	40% (5.10B)	0.608	0.364	0.568	0.196	0.392	0.266	0.310	0.451	0.382	0.230	0.010	0.343	0.462
	20% (6.65B)	0.574	0.301	0.494	0.197	0.353	0.260	0.252	0.592	0.346	0.247	0.003	0.329	0.443
Short GPT	30% (5.88B)	0.561	0.278	0.494	0.195	0.327	0.227	0.252	0.493	0.347	0.256	0.002	0.312	0.420
	40% (5.10B)	0.540	0.258	0.512	0.198	0.307	0.230	0.256	0.417	0.348	0.229	0.000	0.300	0.403
	20% (6.55B)	0.788	0.718	0.730	0.796	0.819	0.540	0.442	0.853	0.546	0.660	0.653	0.686	0.922
Ours	30% (5.74B)	0.753	0.654	0.679	0.658	0.737	0.442	0.392	0.786	0.511	0.501	0.574	0.608	0.817
	40% (4.91B)	0.724	0.587	0.637	0.526	0.680	0.402	0.352	0.775	0.463	0.404	0.487	0.549	0.738

C STATEMENT OF LARGE-LANGUAGE-MODEL (LLM) USAGE

The authors acknowledge that a large-language-model (LLM) was employed as a general-purpose assistance tool during the preparation of this manuscript. Specifically, the following tasks were supported by the LLM under the direct supervision of the authors:

- Formatting and LaTeX assistance The LLM supplied LaTeX snippets for tables, equations, and figure captions (e.g., Table 5 and the hyper-parameter description). The authors integrated these snippets into the manuscript and performed all final compilation and formatting checks.
- Language polishing The LLM was used to improve readability, correct grammar, and adjust stylistic tone across the entire manuscript. The final wording reflects the authors' own decisions after thorough review.

All content generated by the LLM was fully supervised, fact-checked, and substantially revised by the human authors before inclusion in the final version. No portion of the manuscript was submitted to the LLM for autonomous generation without subsequent author verification.

The authors affirm that the intellectual contributions, experimental design, data analysis, and conclusions are entirely their own work, and that the LLM served only as an auxiliary writing and editing aid