

EfficientLLM: Unified Pruning-Aware Pretraining for Auto-Designed Compact Language Models

Xingrun Xing^{1,2,3}, Zheng Liu^{1*}, Shitao Xiao¹, Boyan Gao⁴, Yiming Liang^{1,2,3},
Haokun Lin^{2,3}, Xianlin Zeng⁵, Guoqi Li^{2,3}, Jiajun Zhang^{2,3*}

¹Beijing Academy of Artificial Intelligence,

²School of Artificial Intelligence, University of Chinese Academy of Sciences,

³Institute of Automation, Chinese Academy of Sciences,

⁴University of Oxford, ⁵China RongTong Asset Management Group

Correspondence: zhengliu1026@gmail.com, jjzhang@nlpr.ia.ac.cn

Abstract

Modern large language models (LLMs) driven by scaling laws achieve emergent intelligence in large model sizes. Recently, the increasing concerns about cloud costs, latency and privacy make it an urgent requirement to develop compact edge language models. Distinguished from direct pretraining that bounded by parameter scaling law, this work proposes the unified pruning-aware pretraining, focusing on pretraining compact models while preserving performance of much larger source models, termed EfficientLLM. It features following characteristics: 1) Pruning in Pretraining Corpus: we introduce minimal parameter groups to decouple LLMs and continuously optimize model architecture with classic pruning methods like LLM-Pruner and SparseGPT during pretraining. We reveal that it achieves top-quality compact language models to scale up LLM pruning to large scale pretraining. 2) Auto-Designed Architecture: the LLM architecture is auto-designed during saliency-driven pruning, unifying pretraining, architectural design, and parameter pruning into a single process. Based on these, EfficientLLM significantly outperforms directly pretrained baselines with $100M \sim 1B$ parameters, such as MobileLLM, SmoLLM, Qwen2.5-0.5B, OLMo-1B, Llama3.2-1B in common sense benchmarks, which bridges the performance gap between traditional LLM compression and direct pretraining. We open source on <https://github.com/Xingrun-Xing2/EfficientLLM>.

1 Introduction

Large Language Models (LLMs) have become a central component of modern AI systems (Achiam et al., 2023; Guo et al., 2025) and are increasingly transforming daily life, particularly in mobile edge applications. However, typical LLMs (Touvron

et al., 2023a), with 7 billion to 1 trillion parameters, require on-cloud deployment and continuous internet connectivity for interface. This places significant challenges in terms of latency, data-security and cloud-costs. Therefore, developing compact language models on resource-constrained devices becomes a recent tendency. For instance, MobileLLM (Liu et al., 2024b) focuses on sub-one billion model sizes, which would fit in the DRAM of smartphones without excessive consumption.

Direct pretraining is dominant in recent tiny language model pretraining. Some practices such as MobileLLM and PanGu- π -Pro (Tang et al., 2024) design deep-and-thin architectures for model efficiency. Other practices such as TinyLlama (Zhang et al., 2024a) and Qwen2.5-0.5B (Yang et al., 2024b) focus on scaling up pretraining data to 3T and 17T tokens. Based on best architectures and sufficient data, modern tiny models (Yang et al., 2024b; Groeneveld et al., 2024) are showing promise in reaching performance boundary. However, their overall performance appears to be somewhat locked by the parameter scaling law (Kaplan et al., 2020): given limited model size, simply scaling up pretraining data is inefficient. More importantly, the emergent intelligence (Brown et al., 2020) is only observed on larger model sizes, meaning tiny models may never achieve this by direct pretraining alone. What is the next to train more efficient compact models remains an open challenge.

In parallel, LLM compression (Ashkboos et al., 2024; Gu et al., 2024; Han et al., 2015) focus on retaining the performance of larger and stronger models while reducing computational cost. Despite its potential efficiency, existing methods (Sreenivas et al., 2024; Frantar and Alistarh, 2023; Xiao et al., 2023) compress LLM only using a small calibration dataset in post-training, which often results in significant performance degradation, making them unsuitable for top-quality compact language models. Recently, ShearedLlama (Xia et al., 2023) initial-

*Corresponding author.

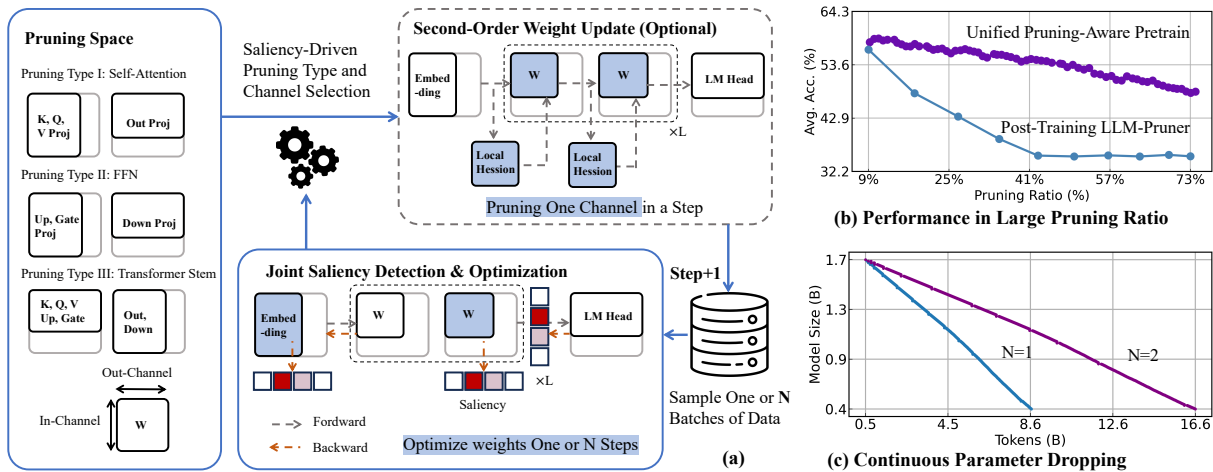


Figure 1: An overview of unified pruning-aware pretraining. (a) Towards **unified pruning, pretraining, and architecture design**, the training loop includes: joint saliency detection and weight optimizing, pruning type selection from pruning space, and weight updating. (b) Improve post-training pruning by pruning-aware pretraining. (c) Continuous model compression during pretraining.

izes from an optimized LLM, improving training efficiency. However, the constrained optimization (Platt and Barr, 1987) hinders scaling up pruning stage and the performance gap to direct pretraining still remains. This work extends the performance boundary of traditional pruning by scaling up to pretraining corpus, a crucial but underexplored approach in the LLM compression field.

This work proposes the unified pruning-aware pretraining to extend the efficiency boundary of compact language models. As shown in Fig. 1, we formulate pruning-aware pretraining as a unified framework for weight pruning, pretraining, and architecture design: 1) Compared with direct pretraining, pruning-aware pretraining leverages the performance of much larger optimized models, which direct pretraining smaller models never achieves. 2) Compared with post-training pruning, it scales up the pruning stage with pretraining data. As shown in Fig. 1 (b), pruning-aware pretraining scales up vanilla LLM-Pruner, achieving more than a 10% increase in accuracy. 3) Driven by saliency, the overall architecture can be auto-designed (Yu et al., 2020; Zoph et al., 2018) according to a predefined pruning space step by step. This work advances efficient compact language models as follows:

- We propose the unified pruning-aware pretraining, promoting LLM compression to the era of pretraining. General post-training methods like LLM-Pruner (Ma et al., 2023), SparseGPT (Frantar and Alistarh, 2023), and Wanda (Sun et al., 2023) could be scaling up in pretraining corpus.
- We explore the auto-designed architectures in

modern pretraining for the first time. Saliency-driven architectures are auto-searched via unified pruning and competitive with human practices.

- We pretrain a family of highly efficient language models in $100M \sim 1B$ sizes, named EfficientLLM, significantly exceeding direct pretraining with unified and scalable pruning.

2 Preliminary and Related Works

Compact Language Models. Towards accurate compact models, a lot of efforts explore the optimal training recipes: 1) data scale. OLMo-1B (Groeneveld et al., 2024), TinyLlama-1.1B (Zhang et al., 2024a), Qwen2.5-0.5B (Yang et al., 2024b) pretrain on 2T, 3T, and 17T tokens respectively, which is significantly larger than the optimal data sizes according to scaling law. 2) Architectures. MobileLLM (Liu et al., 2024b) shows that the deep-and-thin network and layer sharing achieve additional performance gains. However, direct pretraining is bounded by the scaling law, and can be data-inefficient. Llama3.2 (Dubey et al., 2024) and MiniTron (Sreenivas et al., 2024) introduce distillation and pruning for data-efficient training. Despite these advantages, the foundational questions in pretraining oriented LLM pruning remain unexplored: 1) Could LLM pruning process scale up? MiniTron only uses a small calibration dataset for pruning, while this work scales up pruning process in pretraining corpus. 2) Which is the efficient recipe to unify target pruning architecture design, pruning parameters, and pretraining? A unified training framework is proposed in this work.

LLM Pruning (Dong et al., 2024; Zhang et al., 2024b; Zhao et al., 2024; Bhaskar et al., 2024). We mainly focus on structural pruning to address hardware friendly compact language models. The most widely used LLM pruning is based on the Taylor expansion (LeCun et al., 1989; Hassibi et al., 1993; van der Ouderaa et al., 2023). By calibration, typical SparseGPT (Frantar and Alistarh, 2023) and Wanda (Sun et al., 2023) can only applied in semi-structured pruning; LLM-Pruner (Ma et al., 2023) only achieves 20% pruning ratio with reasonable accuracy. Even if pruning with finetuning, LoraPrune (Zhang et al., 2023) can only prune in 50% ratio. And there is an urgent requirement to scale up LLM pruning in the pretraining corpus. Another line of works learn to initialize from source model such as ShearedLlama (Xia et al., 2023) and NutePrune (Li et al., 2024) with less than 0.5B tokens. However, additional human-designed target architectures still required and the data scalability of pruning process remains to explore.

3 Unified Pruning-Aware Pretraining

According to scaling laws, both the scale of training data and the number of parameters are fundamental to the emergence of intelligence in modern LLMs. Direct pretraining of smaller models is inefficient in learning. Model compression methods, although based on pretrained large models, fail to meet the data scale requirements and suffer from significant performance drop. The principle of this work is to bridge the gap between direct pretraining and LLM compression by a unified training scheme.

Problem Formulation. Finding a sub-network from a pretrained LLM is non-trivial. Given an optimized LLM, post-training LLM pruning focuses on finding optimal channels in each layer towards a target architecture. However, for compact models, it is also crucial to design efficient target architectures: for example, MobileLLM shows the deeper architecture is much better than the wider for sub-billion LLMs. Therefore, beyond pruning, we also aim to identify the optimal target model architecture tailored to a specific source model. We formulate the architecture-agnostic pruning problem as:

$$\min_{a \in \mathcal{A}} \min_{c \in \mathcal{C}} \min_w \mathcal{L}_{pretrain}(a, c, w | \mathcal{M}), \quad (1)$$

where \mathcal{A} and \mathcal{C} are sub-architectures (Wu et al., 2019; Liu et al., 2018) and sub-channels sampled from the source model \mathcal{M} . We jointly optimize

pretraining loss through three factors: 1) the sub-architecture, a , 2) the sub-channels, c , and 3) the model weights, w . We outline the pruning-aware pretraining in Fig. 1 and detail each part as follows.

3.1 Defining Minimal Pruning Group

To design architectures automatically via pruning, we first define the minimal parameter groups as the minimal unit to prune in each step, which should be flexible enough to construct any shape transformers after pruning. Given a pretrained source model \mathcal{M} , the pruned model \mathcal{M}^* can be:

$$\mathcal{M}^* = \mathcal{M} - \sum_{t=1}^n g_t, \quad \text{s.t.} \quad \min_{g_t \in \mathcal{G}} \mathcal{L}_{pretrain}(\mathcal{M}) \quad (2)$$

where g_t is the mini-group of parameters pruned in step t , and \mathcal{G} is the pruning space formulated by defined mini-groups. According to Eq. 2, the pruning can be approximately decoupled by t steps and solved sequentially:

$$\begin{aligned} \mathcal{M}_t &= \mathcal{M}_{t-1} - g_t^* \\ \text{s.t.} \quad g_t^* &= \operatorname{argmin}_{g_t \in \mathcal{G}} \mathcal{L}_{pretrain}(g_t | \mathcal{M}_{t-1}). \end{aligned} \quad (3)$$

We first assume an optimal g_t^* in each pruning step to minimize pretraining loss, and solve how to acquire g_t^* in the next subsection. In each pruning step, an optimal mini-group of parameters are selected and dropped from the pretraining LLM \mathcal{M}_{t-1} , allowing the source model \mathcal{M} to adaptively reduce the number of parameters until a specific computation budget is met. For the fully structured pruning space, we impose two constraints in the design of mini-groups: 1) the hidden size, attention heads and intermediate size can be pruned flexibly; 2) the shape of different layers is the same.

Parameter Mini-Groups. For simplify, we indicate the *input layer group* as the *query, key, value projections* in attention blocks; or the *up, gate projections* in feed forward blocks. We indicate the *output layer group* as the *output projections* in attention blocks; or the *down projections* in feed forward blocks. As shown in Fig. 2, we define three basic pruning types and their mini-groups, $\mathcal{G}_{\text{attn}}$, \mathcal{G}_{ffn} , $\mathcal{G}_{\text{stem}}$:

1) Per-head pruning in self-attention blocks: when an attention head is pruned, all the corresponding output channels in the input layer group and input channels in the output layer group are pruned at the same time. We select the mini-group $\mathcal{G}_{\text{attn}}^{(\ell)}$ with the minimal saliency in the ℓ_{th} layer, and merge $\mathcal{G}_{\text{attn}}^{(\ell)}$

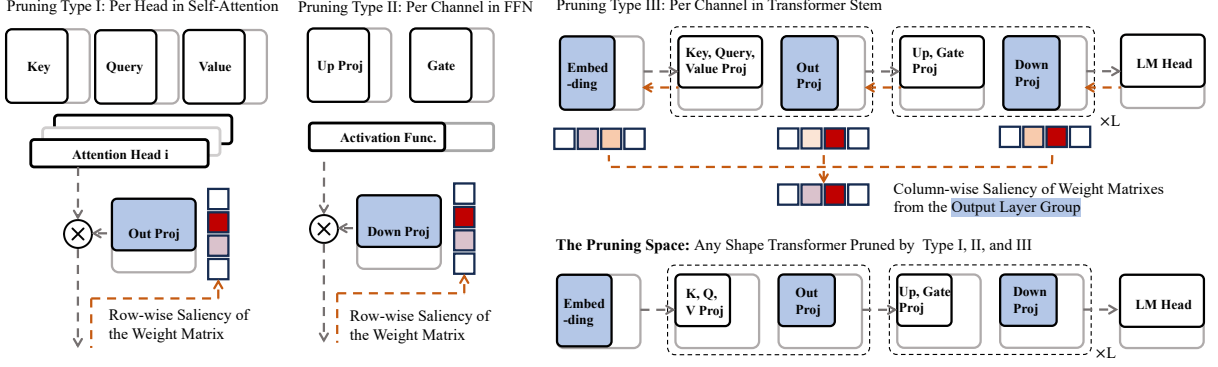


Figure 2: Three basic pruning types in the pruning space. All the weight metrics are shaped in $[D_{input}, D_{output}]$. In backpropagation (in orange), the saliency of the output layer group (in blue) is calculated according to Eq. 9.

in all layers as $\mathcal{G}_{\text{attn}}, \ell = 1, 2, \dots, n$:

$$\mathcal{G}_{\text{attn}} = \{W_{:,ij}^{(k,\ell)}, W_{:,ij}^{(q,\ell)}, W_{:,ij}^{(v,\ell)}, W_{ij,:}^{(o,\ell)}\}, \quad (4)$$

where $W_{:,ij}$ and $W_{ij,:}$ are column-wise and row-wise pruned; ij corresponds to channels of an attention head; and n is the overall blocks.

2) Per-channel pruning in feed-forward blocks: when a intermediate channel is pruned, the coupled channels include one output channel in the input layer group, and one input channel in the output layer group. We couple the mini-group $\mathcal{G}_{\text{ffn}}^{(\ell)}$ in the ℓ_{th} layer and merge $\mathcal{G}_{\text{ffn}}^{(\ell)}$ in all layers as \mathcal{G}_{ffn} :

$$\mathcal{G}_{\text{ffn}} = \{W_{:,i}^{(up,\ell)}, W_{:,i}^{(q,\ell)}, W_{i,:}^{(down,\ell)}\}. \quad (5)$$

3) Per-channel pruning in the transformer stem: When a channel of the transformer stem is pruned, one channel in the token embedding, one input channel in input layer group and one output channel in output layer group for every block, one input channel of the LM head projection is correspondingly pruned at the same time. We donate the stem mini-group as $\mathcal{G}_{\text{stem}}$:

$$\begin{aligned} \mathcal{G}_{\text{stem}} = & \{W_{i,:}^{(k,\ell)}, W_{i,:}^{(q,\ell)}, W_{i,:}^{(v,\ell)}, W_{:,i}^{(o,\ell)}\}, \dots \quad (6) \\ & \cup \{W_{i,:}^{(up,\ell)}, W_{i,:}^{(gate,\ell)}, W_{:,i}^{(down,\ell)}\}, \dots \\ & \cup \{\mathbf{w}_i^{(emb)}, W_{i,:}^{(head)}\}, \end{aligned}$$

where i should be the same in every blocks in Eq. 6; ij needn't the same across blocks in Eq. 4,5.

In each pruning step, the mini-groups are dynamically grouped by saliency, and we only choose among the 3 types to prune, and the final pruning space is 3^t with t -step pruning.

3.2 Optimizing Mini-Groups by Saliency

Based on the mini-groups, Eq.1 becomes a bi-level optimization of the mini-groups g and weights w :

$$\begin{aligned} \min_{g \in \mathcal{G}} \mathcal{L}_{\text{pretrain}}(g, w^* | \mathcal{M}), \quad (7) \\ \text{s.t. } w^* = \underset{w}{\operatorname{argmin}} \mathcal{L}_{\text{pretrain}}(w, g^* | \mathcal{M}), \end{aligned}$$

where the outer optimization could be solved by pruning a mini-group in each step as Eq.3, and the inner optimization could be directly solved by weight pretraining. The weight pretraining step and mini-group optimization (Eq.3) step alternate, and the model size drops continuously during pretraining as shown in Fig. 1 (c), until the pre-defined parameter budget is achieved. *Different from vanilla iterative pruning along a predefined trajectory, unified pruning selects its pruning trajectory according to the mini-group saliency at each step.*

Mini-Group Saliency. In each mini-group selection step, Taylor expansion evaluates the optimal mini-group g_t^* in Eq.3. For an optimized source model, loss of any weight w can be approximated by a second-order Taylor expansion around its optimal value w^* :

$$\mathcal{L}(w) \simeq \mathcal{L}(w^*) + \delta w^\top \nabla \mathcal{L}(w^*) + \frac{1}{2} \delta w^\top \mathbf{H}_{\mathcal{L}}(w^*) \delta w \quad (8)$$

where \mathcal{L} , $\nabla \mathcal{L}$, $\mathbf{H}_{\mathcal{L}}$ is the global loss, gradient, hessian matrix; and $\delta w = w - w^*$. We substitute Eq.8 into Eq.3:

$$\begin{aligned} g_t^* = & \underset{g_t \in \mathcal{G}}{\operatorname{argmin}} \mathcal{L}_{\text{pretrain}}(g_t | \mathcal{M}_{t-1}) \quad (9) \\ = & \underset{g_t \in \mathcal{G}}{\operatorname{argmin}} g_t^\top \nabla \mathcal{L}(\mathcal{M}_{t-1}) + \frac{1}{2} g_t^\top \mathbf{H}_{\mathcal{L}}(\mathcal{M}_{t-1}) g_t, \end{aligned}$$

where we omit the first term $\mathcal{L}(w^*) = \mathcal{L}(\mathcal{M}_{t-1})$ in Eq.8, because $\mathcal{L}(\mathcal{M}_{t-1})$ is the same in the step t

for the 3 mini-groups, $\mathcal{G} = \{\mathcal{G}_{\text{attn}}, \mathcal{G}_{\text{ffn}}, \mathcal{G}_{\text{stem}}\}$. And we could calculate mini-group saliency according to Eq. 9 (Ma et al., 2023).

Efficient Calculation. In practice, we only calculate the saliency of the output layer groups for efficiency. a neural network is a directed acyclic graph (DAG) (Liu et al., 2018). For each node in the graph, pruning all its inputs or all of its outputs is sufficient to prune the entire network. It saves 2 ~ 3 times computation with output layer group only calculation. Details are shown in Fig. 2: 1) *Pruning Type I*: we only calculate element-wise saliency matrix for the weights of the output projection, and then sum each column of the saliency matrix. We select $\mathcal{G}_{\text{attn}}^{(\ell)}$ based on the lowest row-wise saliency in the output projection weights. 2) *Pruning Type II*: we only calculate the element-wise saliency matrix for the down projection, and then sum each column. $\mathcal{G}_{\text{ffn}}^{(\ell)}$ with the lowest row-wise saliency are selected. 3) *Pruning Type III*: we already have all the element-wise saliency in the output layer group based on Type I and II. To evaluate saliency of the hidden state, we first sum each row of saliency matrices, and then, sum the saliency in all of the output layer group.

Hessian Approximations. Existing post-training methods such as LLM-Pruner (Ma et al., 2023), SparseGPT (Frantar and Alistarh, 2023), and Wanda (Sun et al., 2023) have proposed various Hessian approximations to approximate hessian matrices. By substituting Eq. 9, our framework can naturally extend these post-training pruning methods to the unified pretraining stage. One special case of LLM pruning is SparseGPT, which involves additional second-order weight updating. Without loss of generality, we present a column-wise decoupling of the inverse Hessian matrix computation at each step in Appendix A, and demonstrate that it has the same overall computation as the original second-order weight update.

4 Implementations

Models. To compare with the most general post-training pruning, EfficientLLM basically approximates Eq.9 as LLM-Pruner. In empirical results, we also generalize to SparseGPT and Wanda metrics. Various source models are evaluated: 1) In empirical results, we evaluate EfficientLLM with SmolLM-1.7B (Allal et al., 2024), Llama/Llama2-7B (Touvron et al., 2023a,b), and Qwen2.5-7B (Yang et al., 2024b) source mod-

els. 2) For the final EfficientLLM, we prolong and pretrain EfficientLLM-134M from the source model SmolLM-360M; EfficientLLM-457M and 1.1B from SmolLM-1.7B. After pruning-aware pre-training, the auto-designed architectures are shown in Appendix B.

Data Composition. EfficientLLM almost maintains a data distribution similar to the source model: 1) In empirical results, we sample from RedPajama-1T (Weber et al., 2024) as pretraining data with the Llama-7B source model. We directly use the FineWeb-Edu for Llama2/Qwen2.5-7B models. 2) In final EfficientLLM, our pretraining data composition is similar to SmolLM, including 220B tokens from FineWeb-Edu (Lozhkov et al., 2024), 28B tokens from Cosmopedia v2 (Ben Allal et al., 2024a), 4B tokens from Python-Edu (Ben Allal et al., 2024b), and 27.5B tokens randomly sampled from OpenWebMath (Paster et al., 2023).

Training. We apply the 2 stage pretraining as ShearedLlama (Xia et al., 2023) but explore in large scale: the prolonged unified pruning followed by the continued pretraining. For the final EfficientLLM-134M, 460M, and 1.1B, we pretrain 50.3B, 72.1B, and 36.7B tokens for unified pruning, where the token numbers depend on the number of iterations needed to reach the target parameter budget. The continued pretraining contains 500B, 500B, and 320B tokens respectively for higher performance. Note that the large-scale continued pretraining is not necessary, and 50B tokens also achieve competitive performance. All the training details are shown in Appendix C.

Evaluations. For pretrained base models, we follow Llama, MobileLLM, and ShearedLlama to evaluate Common Sense Reasoning tasks: ARC (Clark et al., 2018), BoolQ (Clark et al., 2019), HellaSwag (Zellers et al., 2019), OBQA (Mihaylov et al., 2018), PIQA (Bisk et al., 2020), and WinoGrande (Sakaguchi et al., 2021). The MMLU (Hendrycks et al., 2020) for Word Knowledge evaluation is also applied. For instruct finetuned model, we use the standard Alpaca-Eval (Li et al., 2023) with GPT-4o as the judge model.

5 Empirical Results

In this section, we evaluate unified pruning-aware pretraining with a relatively small amount of data. Based on the empirical results, we extend this approach to EfficientLLM in the next section. As shown in Table 1, 2, 5, we train around 450M tar-

Table 1: Comparison of different training schemes under **the same token budget**. 1) Direct pretraining: "hidden" or "source" indicates target architectures from the source model. 2) Comparison with ShearedLlama in >70% pruning ratio. 3) Comparison with LLM-Pruner with continued pretraining in smaller pruning ratios. We prune both LLaMA2-7B and Qwen2.5-7B for 4000 steps to 3.3B and 5.2B.

Model	Source	ARC-c	ARC-e	BoolQ	HellaSwag	OBQA	PIQA	WinoGrande	Avg.
Direct-hidden	–	26.88	55.85	52.39	37.10	30.20	66.65	50.51	45.65
Direct-source	–	28.67	59.68	56.30	38.41	32.00	65.67	49.96	47.24
Searched Arch. (Ours)	–	29.10	59.30	61.38	37.74	32.80	66.76	51.07	48.30
ShearedLlama-hidden	SmolLM-1.7B	28.41	57.41	60.98	39.79	30.80	66.81	54.06	48.32
ShearedLlama-source	SmolLM-1.7B	30.89	62.08	61.07	44.29	32.60	68.39	52.25	50.22
Unified FewShot Prun.	SmolLM-1.7B	30.63	62.67	61.22	44.31	34.00	68.66	53.43	50.25
Unified Prun. Pretrain.	SmolLM-1.7B	30.46	64.06	61.99	45.98	34.00	69.91	53.91	51.47
+Second-order Update	SmolLM-1.7B	30.97	63.22	60.86	46.51	35.00	69.64	55.09	51.61
LLM-Pruner	LLaMA2-7B	27.30	43.52	61.10	38.15	29.40	63.11	51.14	44.82
Unified Prun. Pretrain.	LLaMA2-7B	32.76	57.53	65.96	53.07	33.80	69.42	58.09	52.95
LLM-Pruner	Qwen2.5-7B	38.40	66.67	70.58	57.97	38.40	72.52	60.62	57.88
Unified Prun. Pretrain.	Qwen2.5-7B	40.10	70.50	76.79	61.26	40.80	73.34	61.72	60.64

Table 2: Generalization to various pruning metrics. LLM-Pruner (Ma et al., 2023), OBC (Frantar and Alistarh, 2022), Diag-Hessian (Sun et al., 2023; LeCun et al., 1989) metrics are embedded in the unified pruning. All results are without continued pretraining.

Model	ARC-c	ARC-e	BoolQ	HellaSwag	OBQA	PIQA	WinoGrande	Avg.
LLM-Pruner	22.18	26.64	49.63	25.62	27.80	51.20	50.59	36.24
Unified Prun. Pretrain.	29.18	57.79	59.57	41.93	34.40	66.97	52.88	48.96
OBC (SparseGPT)	25.51	25.88	37.86	26.81	29.20	51.41	50.51	35.31
Unified Prun. Pretrain.	29.44	60.02	62.02	42.27	33.40	67.41	53.35	49.70
Diag-Hess (Wanda, OBD)	25.85	25.88	50.80	26.06	30.40	51.69	48.22	36.99
Unified Prun. Pretrain.	29.95	60.14	60.83	41.68	32.80	66.38	52.88	49.24

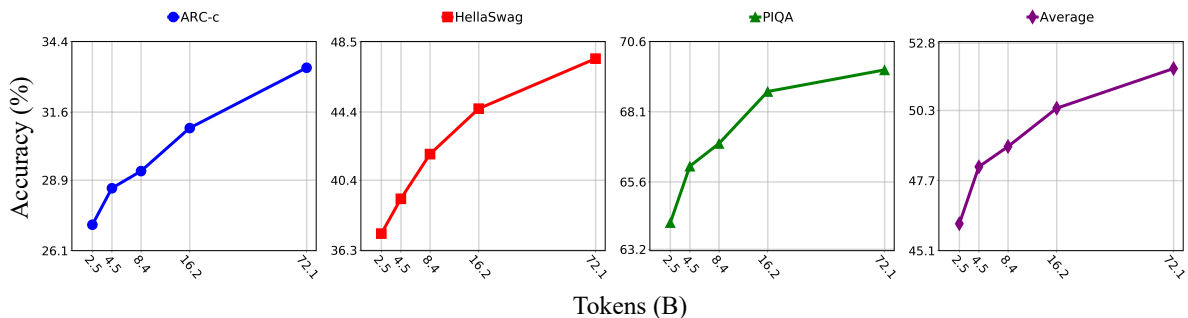


Figure 3: Scalability of unified pruning-aware pretraining without continued pretraining.

get models from the source model SmolLM-1.7B.

5.1 Comparison of Pruning Recipes

Training Schemes. We first compare the efficiency of different training schemes under the same 5B token budget: 1) Direct Pretraining: three architectures are pretrained from scratch using 5B tokens: (i) Direct-hidden: Prunes the hidden size of the source architecture. (ii) Direct-source: Uniformly scales down the source model. (iii) Searched architecture: Uses the architecture searched by Unified Pruning. 2) ShearedLlama: Requires manually specified pruning targets. We consider two variants: pruning hidden size or uniformly scaling the source model. 3) Unified Few-shot Pruning: Applies mini-group optimization on 128 samples using LLM-Pruner’s metric, followed by training on 5B tokens. 4) Unified Pruning-aware Pretrain-

ing: Follows the ShearedLlama pipeline, pruning with 1B tokens, then continuing training on the remaining 4B tokens.

We evaluate the necessity of unified architecture auto-design, pruning, and pretraining in Table 1: 1) **Architecture matters:** For directly pretrained compact models, the choice of architecture significantly impacts performance. The searched Arch. outperforms Direct-hidden by 2.7%. 2) **Pruning boosts compact models:** With the same architecture, Unified Pruning-aware Pretraining surpasses Search Arch. by 3.17%. 3) **Scaling up pruning can benefit more than finetuning:** Compared to Unified Few-shot Pruning, Pruning-aware pretraining allocates more tokens to the pruning stage and improves accuracy by 1.22%. Unlike ShearedLlama, Unified Pruning requires no manual architecture design. ShearedLlama-hidden adopts a subop-

Table 3: Comparisons of LLM pruning in Llama-7B. We scale up pruning-aware pretraining to 5B tokens for EfficientLLM. #Tuning donates whether to finetune after pruning. Most works report finetuned results.

#Ratio	Model	#Tuning	ARC-c	ARC-e	BoolQ	HellaSwag	OBQA	PIQA	WinoGrande	Avg.
50%	MaP	✓	30.63	49.32	39.69	42.49	31.40	66.81	50.67	44.43
	MvP	✓	26.79	44.07	59.94	40.98	31.80	63.06	55.64	46.04
	WANDA	✓	34.20	42.68	50.90	38.12	38.78	57.38	55.98	45.43
	LLM-Pruner	✓	28.24	46.46	61.47	47.56	35.20	68.82	55.09	48.98
	LoRAPrune	✓	31.62	45.13	61.88	47.86	34.98	71.53	55.01	49.72
	LoRAShear	✓	32.26	47.68	62.12	48.01	34.61	71.80	56.29	50.40
	Compresso	✓	27.82	48.82	60.09	39.31	33.40	66.70	51.93	46.87
	NutePrune	✗	31.74	46.59	62.20	53.87	35.80	69.91	57.77	51.13
	NutePrune	✓	32.17	51.68	62.26	55.88	34.40	71.00	57.54	52.13
	EfficientLLM	✗	30.80	52.15	62.29	54.70	35.20	71.33	56.75	51.89
	EfficientLLM	✓	34.04	64.81	64.83	60.12	34.60	73.88	61.48	56.25
70%	LLM-Pruner	✓	24.83	39.56	47.28	31.66	28.80	60.83	50.75	40.53
	NutePrune	✓	26.19	42.17	62.08	39.43	30.20	62.30	51.46	44.83
	EfficientLLM	✗	27.73	54.50	47.89	47.77	31.00	68.17	55.17	47.46
	EfficientLLM	✓	29.95	58.59	58.13	52.02	34.60	70.08	55.96	51.33

timal target, resulting in 3.15% accuracy drop.

Pruning Metrics. We embed well-studied post-training pruning metrics in our Unified Pruning-aware Pretraining, enhancing performance of existing methods. By replacing Eq. 9, we apply LLM-Pruner and diagonal Hessian-based metrics into Unified Pruning. We further extend this by adding second-order weight updates (Section 3.3) to support OBC-based methods. Reaching the pruning target takes about 4,000 steps with a batch size of 1M tokens. As shown in Table 2, LLM-Pruner, OBC, and diagonal Hessian metrics improve accuracy by 12.72%, 14.39%, and 12.25%, respectively. According to Appendix A.2, the second-order update performs better in small scale pretraining, and similar in large scale. For generality, we apply the LLM-pruner metric in main results.

Architecture Robustness. As shown in Appendix B.1, architectures are stably optimized according to dynamic saliency. In Appendix B.2, we evaluate different pruning trajectories to reach the searched architecture. Once we find the optimal architecture, a different trajectory can achieve the same or better results.

5.2 Comparison with LLM Pruning

We mainly focus on large pruning ratio because it is more practical to achieve highly efficiency based on heavy source LLMs. In Table 3, we scale up pruning-aware pretraining to only 5B tokens. We report both results with or without finetuning after pruning. Because previous works finetune in different settings, we finetune additional 1B tokens if with it. Notice that, even without finetuning, EfficientLLM exceeds all the according baselines. It is shown that existing LLM pruning is impractical in large pruning ratio. By simply scaling up LLM-Pruner metric in pruning-aware pre-

training, EfficientLLM significantly exceeds SoTA NutePrune 6.5% in 70% ratio without bells and whistles, while NutePrune integrates distillation and additional learnable masks. In 50% ratio, EfficientLLM exceeds LoRAPrune by 2.18% and 6.54% when with and without tuning. Experiments reveal that only scaling up the pruning stage to 5B tokens can achieve much higher performance than previous results, highlighting the importance of scalable pruning methods.

6 EfficientLLM

Language Modeling. For fair comparison, we collect main streams of compact language models in $100M \sim 1B$ sizes, evaluate in the same conditions (Appendix C.2), and make a benchmark in Table 4. 1) Early compact models like Pythia (Biderman et al., 2023) are direct pretrained in limited tokens and sub-optimal architectures, which largely hinder the performance. By leveraging unified pruning, EfficientLLM achieves both architecture and data efficiency. For instance, EfficientLLM-134M exceeds Pythia-410M by 4.13% average accuracy; EfficientLLM-1.1B with 50B tokens exceeds OLMo-1B, TinyLlama, Llama3.2-1B in accuracy. 2) Compared with the SoTA compact model MobileLLM (Liu et al., 2024b), EfficientLLM-134M exceeds MobileLLM-125M by 5.52% with the large scale model compression. 3) Modern SoTA industrial models scaling up pretraining tokens like Qwen (Yang et al., 2024a,b) and Llama3.2-1B, EfficientLLM-457M and 1.1B outperforms Qwen2.5-0.5B by 4.15% and Llama3.2-1B by 4.06% respectively with limited pretraining data. Because of pruning from source model, EfficientLLM-457M achieves higher accuracy with **62** \times and **16** \times fewer GPU hours than Qwen2.5-0.5B when using 50B and 500B continued pretrain-

Table 4: Zero-shot performance on World Knowledge and Common Sense Reasoning tasks. ‘‘Avg.’’ calculate among the 7 Common Sense Reasoning tasks. #Tokens count continued pretraining for EfficientLLM. All the results are evaluated on the same evaluation (Appendix C).

Model	#Tokens	#Params	MMLU	ARC-c	ARC-e	BoolQ	HellaSwag	OBQA	PIQA	WinoGrande	Avg.
Pythia-160M	300B	162M	26.43	22.27	37.84	43.33	29.97	26.40	58.87	49.96	38.38
Mamba-130M	1.2T	130M	27.65	24.49	47.56	54.68	35.11	29.00	64.69	53.35	44.13
MobileLLM-125M	1T	125M	27.58	24.32	46.38	60.34	38.15	28.40	65.13	52.41	45.02
SmolLM-135M	600B	135M	30.05	29.35	61.32	59.85	42.67	34.40	68.55	52.96	49.87
EfficientLLM	500B	134M	30.54	30.97	62.88	60.40	43.81	33.60	68.82	53.28	50.54
Pythia-410M	300B	405M	29.10	24.15	51.39	59.20	40.20	29.40	66.70	53.83	46.41
MobileLLM-350M	1T	345M	30.21	27.39	56.40	61.96	49.51	31.00	68.88	57.14	50.33
SmolLM-360M	600B	362M	33.89	36.26	70.16	55.23	53.51	37.60	71.38	57.22	54.48
Qwen2-0.5B	15T	494M	31.85	28.50	55.05	61.25	49.16	32.80	69.75	57.22	50.53
Qwen2.5-0.5B	17T	494M	33.37	32.17	64.44	61.99	52.09	35.20	70.29	56.20	53.20
EfficientLLM	50B	457M	33.09	35.92	70.50	59.85	53.16	35.00	72.69	56.27	54.77
EfficientLLM	500B	457M	34.54	38.40	72.10	62.42	56.84	40.40	73.83	57.46	57.35
Pythia-1B	300B	1.0B	30.14	26.96	56.86	60.04	47.15	31.20	70.29	52.88	49.34
TinyLlama-1.1B	3T	1.1B	32.30	30.29	60.40	56.85	59.13	35.80	73.07	59.04	53.51
ShearedLlama-1.3B	50B	1.3B	31.51	29.44	61.07	61.83	59.33	34.40	73.94	58.01	54.00
OLMo-1B	2T	1.2B	32.03	30.72	63.55	61.38	62.86	36.40	75.35	59.35	55.66
Llama3.2-1B	–	1.2B	36.31	31.48	65.28	63.88	63.69	37.40	74.59	60.54	56.69
EfficientLLM	50B	1.1B	36.71	40.36	73.61	62.39	60.24	40.20	75.19	61.25	59.03
EfficientLLM	320B	1.1B	37.71	42.24	73.48	67.09	64.09	41.80	75.41	61.17	60.75

Table 5: Ablation studies on the EfficientLLM-457M. ‘‘Direct Training’’ keeps the same architecture ratio as the source model. ‘‘Pruned Arch.’’ indicates the auto-designed architecture. Settings are the same as Sec. 5.1.

Model	ARC-c	ARC-e	BoolQ	HellaSwag	OBQA	PIQA	WinoGrande	Avg.
Direct Training	28.67	59.68	56.30	38.41	32.00	65.67	49.96	47.24
+Searched Arch.	29.10	59.30	61.38	37.74	32.80	66.76	51.07	48.30
+Unified FewShot Pruning	30.63	62.67	61.22	44.31	34.00	68.66	53.43	50.25
+Prolonged Pruning	34.13	66.16	60.49	49.87	35.80	70.67	54.30	53.06
+Prolonged Continued Pretrain.	35.92	70.50	59.85	53.16	35.00	72.69	56.27	54.77

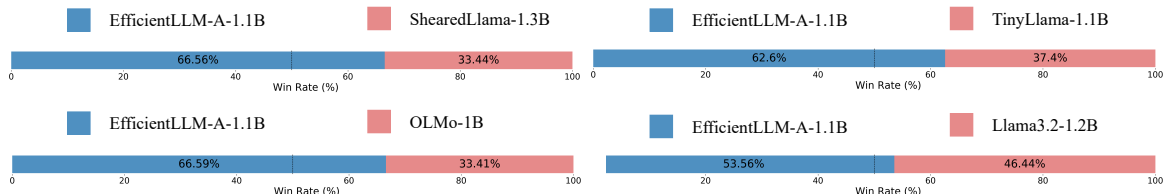


Figure 4: Win rate of EfficientLLM in the instruction tuning task.

ing tokens in Appendix D.1.

Ablation Studies. To evaluate each part of the unified pruning, we decouple into 3 basic designs to improve compact language model pretraining: the auto-designed architecture, pruning, and scalable pruning stage. As shown in Fig. 3, we scale up unified pruning according to Appendix C.3. As shown in Table 5, the pruned architecture, unified pruning, and scaling up pruning with 78B tokens continuously improve 5.82% accuracy. We further scale up the continued pretraining (or finetuning) to 50B tokens, as in ShearedLlama, and achieve 7.53% accuracy overall.

Instruction Tuning. We finetune EfficientLLM-1.1B and other top-quality open-source base models includes OLMo-1B, ShearedLlama-1.3B, TinyLlama-1.1B and Llama3.2-1B in the same condition. We finetune on the Alpaca dataset (Taori et al., 2023) with 52K instructions for 3 epochs. As shown in Fig. 4, EfficientLLM-1.1B significantly outperforms SoTA baselines, indicating the gener-

alization ability in the supervised finetuning (SFT). More case studies are shown in Appendix E.

Inference Speed & Quantization. Most edge devices are non-GPU environments. We deploy compact models using 1, 2, 4, and 8 Intel Xeon @ 2.90GHz CPUs respectively. As shown in Appendix D.2, when using 2 CPUs, EfficientLLM-457M speeds up $\times 8.7$ and $\times 3.2$ than MobileLLM-350M and Qwen2.5-0.5B respectively; EfficientLLM-1B speeds up $\times 7.3$ and $\times 1.2$ than OLMo-1B and Llama3.2-1B. In Appendix D.3, we further quantize EfficientLLM with 8 bit weights and 8 bit activations (8W8A) using the general OmniQuant (Shao et al., 2023). After 8W8A quantization, EfficientLLM-457M even improves 0.09% and EfficientLLM-1B only drops 0.27% average accuracy.

7 Conclusion

This work primarily advances the compact language model pretraining to exceed the traditional

direct pretraining. Distinguished from almost LLM compression in post-training, this work scales up existing pruning metric in the pretraining corpus, promoting LLM compression to the era of pretraining. Technically, minimal parameter groups are defined and optimized by saliency to address scalable target-agnostic pruning. The results reveal that even if vanilla LLM-Pruner can surpass SoTA pruning methods by scaling up and outperform direct pretraining compact models.

Limitations

One limitation of this work is the lack of evaluation on recent reasoning-oriented foundation models such as Qwen3 (Yang et al., 2025) and DeepSeek-V3 (Liu et al., 2024a). The pretraining corpora for reasoning models are largely inaccessible, which makes it difficult to obtain data with a similar distribution for fair comparison. Consequently, we exclude reasoning-focused evaluations in this study.

Acknowledgments

This work is supported by the Strategic Priority Research Program of Chinese Academy of Sciences under Grant XDA04080400 and Beijing Natural Science Foundation L259016.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Joshua Ainslie, James Lee-Thorp, Michiel de Jong, Yury Zemlyanskiy, Federico Lebrón, and Sumit Sanghai. 2023. Gqa: Training generalized multi-query transformer models from multi-head checkpoints. *arXiv preprint arXiv:2305.13245*.
- Loubna Ben Allal, Anton Lozhkov, Elie Bakouch, Leandro von Werra, and Thomas Wolf. 2024. Smolm - blazingly fast and remarkably powerful.
- Saleh Ashkboos, Maximilian L Croci, Marcelo Genari do Nascimento, Torsten Hoefler, and James Hensman. 2024. Slicept: Compress large language models by deleting rows and columns. *arXiv preprint arXiv:2401.15024*.
- Loubna Ben Allal, Anton Lozhkov, Guilherme Penedo, Thomas Wolf, and Leandro von Werra. 2024a. [Cosmopedia](#).
- Loubna Ben Allal, Anton Lozhkov, Guilherme Penedo, Thomas Wolf, and Leandro von Werra. 2024b. [Smolm-corporus](#).

Adithya Bhaskar, Alexander Wettig, Dan Friedman, and Danqi Chen. 2024. Finding transformer circuits with edge pruning. *arXiv preprint arXiv:2406.16778*.

Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O’Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, and 1 others. 2023. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning*, pages 2397–2430. PMLR.

Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, and 1 others. 2020. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 7432–7439.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, and 1 others. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. 2019. Boolq: Exploring the surprising difficulty of natural yes/no questions. *arXiv preprint arXiv:1905.10044*.

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. *arXiv preprint arXiv:1803.05457*.

Peijie Dong, Lujun Li, Zhenheng Tang, Xiang Liu, Xinglin Pan, Qiang Wang, and Xiaowen Chu. 2024. Pruner-zero: Evolving symbolic pruning metric from scratch for large language models. *arXiv preprint arXiv:2406.02924*.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, and 1 others. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Clémentine Fourrier, Nathan Habib, Thomas Wolf, and Lewis Tunstall. 2023. [Lighteval: A lightweight framework for llm evaluation](#).

Elias Frantar and Dan Alistarh. 2022. Optimal brain compression: A framework for accurate post-training quantization and pruning. *Advances in Neural Information Processing Systems*, 35:4475–4488.

Elias Frantar and Dan Alistarh. 2023. Sparsegpt: Massive language models can be accurately pruned in one-shot. In *International Conference on Machine Learning*, pages 10323–10337. PMLR.

- 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, and 5 others. 2024. [A framework for few-shot language model evaluation](#).
- Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya Harsh Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, and 1 others. 2024. Olmo: Accelerating the science of language models. *arXiv preprint arXiv:2402.00838*.
- Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang. 2024. Minillm: Knowledge distillation of large language models. In *The Twelfth International Conference on Learning Representations*.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Song Han, Huizi Mao, and William J Dally. 2015. Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding. *arXiv preprint arXiv:1510.00149*.
- Babak Hassibi, David G Stork, and Gregory J Wolff. 1993. Optimal brain surgeon and general network pruning. In *IEEE international conference on neural networks*, pages 293–299. IEEE.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*.
- Yann LeCun, John Denker, and Sara Solla. 1989. Optimal brain damage. *Advances in neural information processing systems*, 2.
- J Li, A Fang, G Smyrnis, M Ivgi, M Jordan, S Gadre, H Bansal, E Guha, S Keh, K Arora, and 1 others. Datacomp-lm: In search of the next generation of training sets for language models, 2024. [URL https://arxiv.org/abs/2406.11794](https://arxiv.org/abs/2406.11794).
- Shengrui Li, Junzhe Chen, Xueting Han, and Jing Bai. 2024. Nutep prune: Efficient progressive pruning with numerous teachers for large language models. *arXiv preprint arXiv:2402.09773*.
- Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. 2023. AlpacaEval: An automatic evaluator of instruction-following models.
- Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, and 1 others. 2024a. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*.
- Hanxiao Liu, Karen Simonyan, and Yiming Yang. 2018. Darts: Differentiable architecture search. *arXiv preprint arXiv:1806.09055*.
- Zechun Liu, Changsheng Zhao, Forrest Iandola, Chen Lai, Yuandong Tian, Igor Fedorov, Yinyang Xiong, Ernie Chang, Yangyang Shi, Raghuraman Krishnamoorthi, and 1 others. 2024b. MobileLLM: Optimizing sub-billion parameter language models for on-device use cases. *arXiv preprint arXiv:2402.14905*.
- Anton Lozhkov, Loubna Ben Allal, Leandro von Werra, and Thomas Wolf. 2024. [Fineweb-edu: the finest collection of educational content](#).
- Xinyin Ma, Gongfan Fang, and Xinchao Wang. 2023. Llm-pruner: On the structural pruning of large language models. *Advances in neural information processing systems*, 36:21702–21720.
- Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. 2018. Can a suit of armor conduct electricity? a new dataset for open book question answering. *arXiv preprint arXiv:1809.02789*.
- Keiran Paster, Marco Dos Santos, Zhangir Azerbayev, and Jimmy Ba. 2023. [Openwebmath: An open dataset of high-quality mathematical web text](#). *Preprint*, arXiv:2310.06786.
- John Platt and Alan Barr. 1987. Constrained differential optimization. In *Neural Information Processing Systems*.
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2021. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106.
- Wenqi Shao, Mengzhao Chen, Zhaoyang Zhang, Peng Xu, Lirui Zhao, Zhiqian Li, Kaipeng Zhang, Peng Gao, Yu Qiao, and Ping Luo. 2023. Omniquant: Omnidirectionally calibrated quantization for large language models. *arXiv preprint arXiv:2308.13137*.
- Sharath Turuvekere Sreenivas, Saurav Muralidharan, Raviraj Joshi, Marcin Chochowski, Mostofa Patwary, Mohammad Shoeybi, Bryan Catanzaro, Jan Kautz, and Pavlo Molchanov. 2024. Llm pruning and distillation in practice: The minitron approach. *arXiv preprint arXiv:2408.11796*.

- Mingjie Sun, Zhuang Liu, Anna Bair, and J Zico Kolter. 2023. A simple and effective pruning approach for large language models. *arXiv preprint arXiv:2306.11695*.
- Yehui Tang, Fangcheng Liu, Yunsheng Ni, Yuchuan Tian, Zheyuan Bai, Yi-Qi Hu, Sichao Liu, Shangling Jui, Kai Han, and Yunhe Wang. 2024. Rethinking optimization and architecture for tiny language models. *arXiv preprint arXiv:2402.02791*.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, and 1 others. 2023a. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, and 1 others. 2023b. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Tycho FA van der Ouderaa, Markus Nagel, Mart Van Baalen, Yuki M Asano, and Tijmen Blankevoort. 2023. The llm surgeon. *arXiv preprint arXiv:2312.17244*.
- Maurice Weber, Daniel Fu, Quentin Anthony, Yonatan Oren, Shane Adams, Anton Alexandrov, Xiaozhong Lyu, Huu Nguyen, Xiaozhe Yao, Virginia Adams, and 1 others. 2024. Redpajama: an open dataset for training large language models. *arXiv preprint arXiv:2411.12372*.
- Bichen Wu, Xiaoliang Dai, Peizhao Zhang, Yanghan Wang, Fei Sun, Yiming Wu, Yuandong Tian, Peter Vajda, Yangqing Jia, and Kurt Keutzer. 2019. Fbnet: Hardware-aware efficient convnet design via differentiable neural architecture search. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10734–10742.
- Mengzhou Xia, Tianyu Gao, Zhiyuan Zeng, and Danqi Chen. 2023. Sheared llama: Accelerating language model pre-training via structured pruning. *arXiv preprint arXiv:2310.06694*.
- Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. 2023. Smoothquant: Accurate and efficient post-training quantization for large language models. In *International Conference on Machine Learning*, pages 38087–38099. PMLR.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*.
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, and 1 others. 2024a. Qwen2 technical report. *Preprint, arXiv:2407.10671*.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, and 1 others. 2024b. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*.
- Jiahui Yu, Pengchong Jin, Hanxiao Liu, Gabriel Bender, Pieter-Jan Kindermans, Mingxing Tan, Thomas Huang, Xiaodan Song, Ruoming Pang, and Quoc Le. 2020. Bignas: Scaling up neural architecture search with big single-stage models. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part VII 16*, pages 702–717. Springer.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. Hellaswag: Can a machine really finish your sentence? *arXiv preprint arXiv:1905.07830*.
- Mingyang Zhang, Hao Chen, Chunhua Shen, Zhen Yang, Linlin Ou, Xinyi Yu, and Bohan Zhuang. 2023. Loraprune: Pruning meets low-rank parameter-efficient fine-tuning. *arXiv preprint arXiv:2305.18403*.
- Peiyuan Zhang, Guangtao Zeng, Tianduo Wang, and Wei Lu. 2024a. Tinyllama: An open-source small language model. *arXiv preprint arXiv:2401.02385*.
- Yingtao Zhang, Haoli Bai, Haokun Lin, Jialin Zhao, Lu Hou, and Carlo Vittorio Cannistraci. 2024b. Plug-and-play: An efficient post-training pruning method for large language models. In *The Twelfth International Conference on Learning Representations*.
- Bowen Zhao, Hannaneh Hajishirzi, and Qingqing Cao. 2024. Apt: Adaptive pruning and tuning pretrained language models for efficient training and inference. *arXiv preprint arXiv:2401.12200*.
- Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V Le. 2018. Learning transferable architectures for scalable image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8697–8710.

Appendix

A Second-Order Weight Updating

A.1 Efficient Calculation

Existing second order pruning applies the same Hessian matrix for the pruning weight detection and the remaining weight updating. However, calculating the global Hessian matrix is impossible in modern LLMs for its $\mathcal{O}(n^4)$ complexity. A common approach is to use the squared error at each layer as a proxy for the global loss: $\mathbf{H}_{\mathcal{L}} \simeq XX^T$, such as in SparseGPT (Frantar and Alistarh, 2023), OBC (Frantar and Alistarh, 2022). Although achieving the $\mathcal{O}(d_{row} \times d_{col}^2)$ complexity, Hessian matrices can not describe the global loss.

This work addresses this problem by decoupling the Hessian matrix in saliency detection and weight updating. To capture global saliency, we approximate with global diagonal Hessian matrices as LLM-Pruner for saliency detection; to reduce computational complexity, we apply the layerwise proxy Hessian, $\mathbf{H}_{\mathcal{L}} \simeq XX^T$, for weight updating. In each step, we prune a mini-group including only one column of weights in a layer, and the remaining weights are updated by $\delta w_p = -\frac{w_p}{[\mathbf{H}^{-1}]_{pp}} \cdot \mathbf{H}_{:,p}^{-1}$. To efficiently compute the p -th column of the inverse Hessian matrix $\mathbf{H}_{:,p}^{-1}$, it suffices to solve the linear equation $\mathbf{H}\mathbf{H}_{:,p}^{-1} = \mathbf{e}_p$ in a weight updating step.

A.2 Performance

As shown in Table 7, we evaluate second-order weight updating by pruning Llama2-7B to 1.3B with RedPajama dataset and SmoLLM-1.7B to 457M. The second-order updating improves in small scale pretraining, while their results become similar in large scale pretraining. We finally exclude the second-order updating in scaling up.

B Auto-Designed Architectures

B.1 Visualization

As shown in Fig.5, we visualize the pruning-aware pretraining. We prune SmoLLM-1.7B to EfficientLLM-457M. In Fig.5 (right), the self-attention parameter groups and FFN parameter groups are iteratively pruned in the initial stage. After 44.49B-token pretraining, the transformer stem parameter groups start to be pruned. This indicates that for the typical human-designed transformer shape, there are more redundant parameters in the attention head and the intermediate size of FFN compared with the transformer stem.

B.2 Architecture Robustness

We conduct further experiments to demonstrate how different pruning paths can lead to stable (and generally better) results. We began with 2000 steps of Unified Pruning to obtain an automatically designed model architecture, denoted as A^* . We then mark A^* as the target and restarted training from scratch, enforcing different pruning paths. Specifically, we randomly sampled three pruning paths:

- Prune in the order: Stem \rightarrow Attention \rightarrow FFN.
- Prune in the order: Attention \rightarrow Stem \rightarrow FFN.
- Prune in the order: FFN \rightarrow Stem \rightarrow Attention.

The experimental results show that, given a known target structure A^* , different pruning paths consistently lead to similar or even better performance. However, this does not imply that scaling up the pruning process is unimportant. On the contrary, it plays a key role in discovering a more accurate target architecture.

Table 6: Performance of different pruning trajectories (7 zero-shot average).

Pruning Trajectory	Avg. (7 zero-shot)
stem_attn_ffn	44.52%
attn_stem_ffn	44.11%
ffn_stem_attn	44.15%
Unified Pruning	43.92%

B.3 Architecture Comparisons

As shown in Table 8, we compare the auto-designed architectures by saliency via pruning and the best practices of human design, including MobileLLM and Qwen2/2.5-0.5B, OLMo-1B, ShearedLlama-1.3B. In EfficientLLM, the pruning ratio of hidden-size is smaller than attention heads and FFN intermediate channels driven by saliency.

As shown in Table 9, we compare the influence of different pruning metrics, including the classic LLM-Pruner (Ma et al., 2023), OBC (Frantar and Alistarh, 2022), and Diag-Hess (LeCun et al., 1989). The Diag-Hess only uses the second-order term in Eq.9, which applies the diagonal of the Hessian matrix for approximate calculation.

B.4 Cluster Attention

Pruning-aware pretraining could structurally prune the Group Query Attention (GQA) (Ainslie et al.,

Table 7: Comparisons between EfficientLLM and B for the second-order weight updating.

Model	Source	Tokens	ARC-C	ARC-E	BoolQ	HS	OBQA	PIQA	WG	Avg.
Unified Prun. Pretrain. +Second-order Update	Llama2-7B	6B	27.30	56.44	57.58	50.03	31.00	69.10	54.70	49.45
	Llama2-7B	6B	28.58	56.90	62.42	49.81	32.20	68.93	55.49	50.62
Unified Prun. Pretrain. +Second-order Update	SmolLM-1.7B	572B	38.40	72.10	62.42	56.84	40.40	73.83	57.46	57.35
	SmolLM-1.7B	572B	39.59	71.68	62.39	57.21	39.60	73.50	57.70	57.38

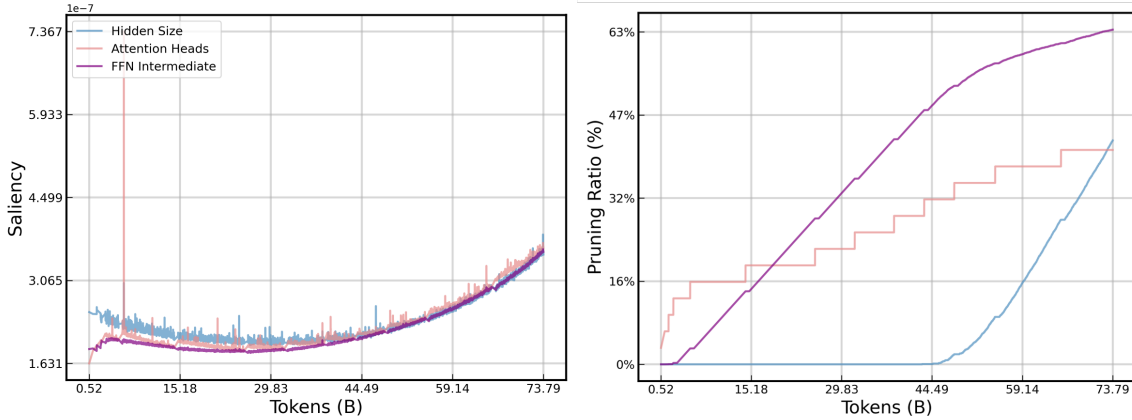


Figure 5: Visualization of pruning-aware pretraining. We plot the saliency of the three pruning types and their pruning ratio in training.

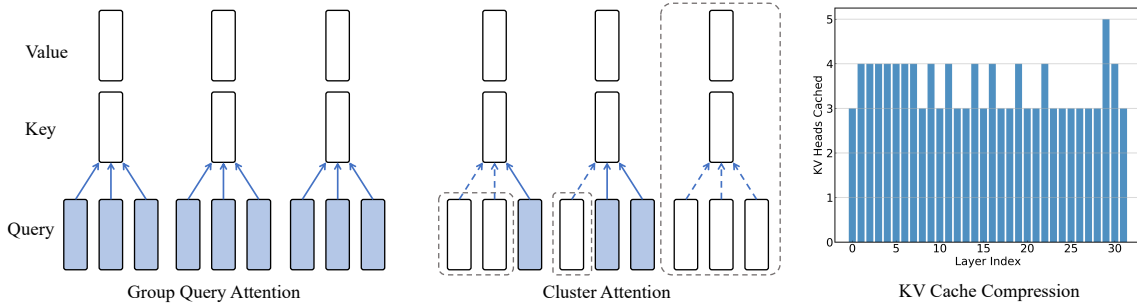


Figure 6: Group Query Attention (GQA) pruning. In the case of GQA, cluster attention can be obtained through pruning. After pruning, the number of query heads is the same in each layer, and the cluster attention compresses the KV Cache.

2023), which is commonly used for KV-cache compression in LLMs. When the source model adopts GQA, different pruning cases may arise:

- As illustrated in Fig. 9, in all cases considered below, the number of query attention heads remains the same across layers and follows standard self-attention. The difference lies in how keys and values are shared among queries.
- If all queries associated with a given key–value pair are pruned, the corresponding key and value are also pruned.
- If only a subset of queries associated with a key–value pair is pruned, the key and value are retained, which effectively results in a form of clustered attention.

We provide an example using EfficientLLM-134M in Fig. 9. Note that the source models for EfficientLLM-457M and EfficientLLM-1.1B do not employ GQA.

C Training and Evaluation Details

C.1 Training

Our training code and models will be fully open-sourced on GitHub and Huggingface. Detailed hyperparameters are shown in Table 10. Note that iterations determine the number of tokens in pruning-aware pretraining to achieve the target model size, which is not directly defined. It can be adjusted through batch size and the pruning frequency in each iteration.

Table 8: Architecture comparisons between EfficientLLM and human-designed models.

Model	Hidden Size	FFN Intermediate	Attention Heads	Head Dim	Layer
MobileLLM-125M	576	1536	9	64	30
EfficientLLM-134M	757	966	5	64	32
MobileLLM-350M	960	2560	15	64	32
Qwen2/2.5	896	4864	14	64	24
EfficientLLM-457M	1195	3006	19	64	24
MobileLLM-1B	1280	3584	20	64	54
ShearedLlama-1.3B	2048	5504	16	128	24
OLMo-1B	2048	8192	16	128	16
Llama3.2-1B	2048	8192	32	64	16
EfficientLLM-1.1B	2048	4870	24	64	24

Table 9: Architectures in different pruning metrics to scale up by pruning-aware pretraining. We compare the approximate 460M model size. “x1” indicates that the number of gradient descent steps and pruning steps in each iteration are 1:1.

Model	Hidden Size	FFN Intermediate	Attention Heads	Head Dim	Layer
LLM-Pruner x1 (Ma et al., 2023)	1169	3082	19	64	24
OBC x1 (Frantar and Alistarh, 2022)	1131	3258	19	64	24
Diag-Hess x1 (LeCun et al., 1989)	1963	1542	12	64	24

Table 10: Hyper-parameters in pruning-aware pretraining and continued pretraining stages.

Model	#Tokens	Learning Rate	WarmUp Steps	Batchsize	Text Length	#GPU
Pruning-134M	50.3B	2×10^{-3}	500	2 M	2048	32
Continued Pretrain-134M	500B	2×10^{-3}	10000	1 M	2048	32
Pruning-457M	72.1B	5×10^{-4}	500	1 M	2048	32
Continued Pretrain-457M	50B/500B	2×10^{-3}	10000	1 M	2048	40
Pruning-1.1B	36.7	5×10^{-4}	500	1 M	2048	32
Continued Pretrain-1.1B	50B/500B	5×10^{-4}	10000	1 M	2048	64

C.2 Evaluation

MMLU: According to Datacomp-lm (Li et al.) (Appendix G of Datacomp-lm) and SmoLLM (Allal et al., 2024), taking into account the log probabilities of complete answer sequences in MMLU is more related to weaker model performance, such as compact language models. Following SmoLLM (Allal et al., 2024), we apply the Lighteval-v0.7.0 (Fourrier et al., 2023) to evaluate MMLU zero-shot performance.

Common Sense Reasoning: Follow most of recent works (Xia et al., 2023; Ma et al., 2023; Li et al., 2024), we apply the widely used lm-evaluation-harness package (Gao et al., 2024) to evaluate zero-shot common sense reasoning tasks. To avoid different results introduced by different versions. We evaluate all the benchmarks with the

0.4.3 version. However, some previous works evaluate in older version 0.3.0, and we evaluate with the same version in Table 3, 50% pruning ratio. Finally, all the versions are the same.

C.3 Scalability

In Fig. 3, we evaluate the scalability of pruning-aware pretraining. According to Eq. 7, we set the ratio of pruning steps to gradient descent steps to 4:1, 2:1, 1:1, and 1:9 in a iteration, respectively. When the target model size is reached, the pruning-aware pretraining requires 2.5B, 4.5B, 8.4B, and 72.1B tokens of pretraining, respectively. Fig. 3 indicates that scaling up pruning-aware pretraining continuously improves pruning performance. By scaling up LLM pruning during pretraining, the upper boundary of LLM compression can be ex-

Table 11: Comparison of pre-training and pruning costs in GPU hours. PT indicates pretraining.

Model	PT Tokens	GPU Hours	Pruning GPU Hours	Continued PT GPU Hours	Total Hours	Acc. (%)
Qwen2.5-0.5B	17T	199467	-	-	199467	53.20
EfficientLLM-457M	-	-	72.1B	2166	3184	54.77
EfficientLLM-457M	-	-	72.1B	2166	500B	57.35

Table 12: Inference latency (in milliseconds) and accuracy on different models.

Model (ms)	1 CPU	2 CPUs	4 CPUs	8 CPUs	Acc. (%)
MobileLLM-350M	132.37	81.26	53.99	41.95	51.30
Qwen2.5-0.5B	54.51	30.08	19.87	13.54	53.20
EfficientLLM-457M	15.93	9.38	5.67	4.07	57.35
OLMo-1B	352.92	186.39	101.81	65.12	55.66
Llama3.2-1B	51.45	30.17	16.59	10.33	56.69
EfficientLLM-1B	49.77	25.52	16.00	10.32	60.75

Table 13: Performance comparison of models with different bit-widths.

Model	#Bits	ARC-C	ARC-E	BoolQ	HS	OBQA	PIQA	WG	Avg.
Qwen2.5-0.5B	bf16	32.17	64.44	61.99	52.09	35.20	70.29	56.20	53.20
EfficientLLM-457M	bf16	38.40	72.10	62.42	56.84	40.40	73.83	57.46	57.35
EfficientLLM-457M-Quant	8W8A	38.23	72.14	63.18	56.90	40.00	73.67	57.93	57.44
Llama3.2-1B	bf16	31.48	65.28	63.88	63.69	37.40	74.59	60.54	56.69
EfficientLLM-1B	bf16	42.24	73.48	67.09	64.09	41.80	75.41	61.17	60.75
EfficientLLM-1B-Quant	8W8A	41.64	73.32	66.51	63.98	42.40	75.24	60.30	60.48

tended. As shown in Table 14, EfficientLLM could recover almost the same performance as the source model in the 35% pruning ratio.

D Training and Inference Efficiency

D.1 Training Efficiency

As shown in Table 11, we evaluate pretraining speed under the same environment, as shown in the table. With pruning-aware pretraining, training EfficientLLM-457M requires x16 to x62 times fewer GPU hours compared to Qwen2.5-0.5B, while achieving higher accuracy.

D.2 Inference Efficiency

As shown in Table 12, we deploy EfficientLLM on non-GPU devices. We deploy on 1, 2, 4, and 8 Intel Xeon @ 2.90GHz CPUs respectively. Compared with SoTA compact models, EfficientLLM achieves both higher inference speed (ms/token) and average zero-shot accuracy.

D.3 Quantization

As shown in Table 13, we applied OmniQuant to perform 8-bit weight and activation quantization (8W8A) on EfficientLLM. The experimental results show that 8-bit quantization only causes minor disturbances to the results, further proving the effectiveness of the model acceleration.

Table 14: Comparison between SmolLM-1.7B and EfficientLLM-1.1B.

Metric	SmolLM	EfficientLLM
#Params	1.7B	1.1B
ARC-c	46.16	42.24
ARC-e	76.60	73.48
BoolQ	65.99	67.09
HellaSwag	65.74	64.09
OBQA	42.00	41.80
PIQA	75.95	75.41
WinoGrande	60.14	61.17
Avg.	61.80	60.75

E Case Study in Alpaca Finetuning

In the instruction finetuning task, we finetune EfficientLLM-1.1B, OLMo-1B, ShearedLlama-1.3B, and Llama3.2-1B in the Alpaca dataset. And then, we evaluate using the Alpaca-Eval. We sample some instructions in Alpaca-Eval and make comparisons as follows:

Input:

I want to eventually become a plant/manufacturing/production manager, and eventually a Chief Operating Officer. What sort of experience and post-graduate education path would you recommend? Please provide non-obvious and atypical career paths leading to such positions, on top of the classic ones.

Output (EfficientLLM-1.1B):

To become a plant/manufacturing/production manager and eventually a Chief Operating Officer, you will need to have a combination of experience and post-graduate education.

For experience, you can start by working in a production or manufacturing role. This will give you the opportunity to gain hands-on experience in the industry and learn about the day-to-day operations of a production or manufacturing facility. You can also gain experience by working in a related field, such as engineering, logistics, or supply chain management.

For post-graduate education, you can pursue a degree in a related field, such as industrial engineering, operations management, or supply chain management. You can also pursue a master's degree in a related field, such as operations management or supply chain management. Additionally, you can pursue professional certifications, such as the Certified Production and Inventory Management Professional (CPIMP) or the Certified Supply Chain Professional (CSCP).

In addition to experience and post-graduate education, you will also need to have excellent communication, organizational, and problem-solving skills. You will also need to be able to lead and motivate a team, as well as be able to think strategically and make decisions quickly.

Overall, to become a plant/manufacturing/production manager and eventually a Chief Operating Officer, you will need to have a combination of experience and post-graduate education, as well as excellent communication, organizational, and problem-solving skills. Additionally, you will need to be able to lead and motivate a team, as well as be able to think strategically and make decisions quickly.

Output (OLMo-1B):

To become a plant/manufacturing/production manager, you will need at least a Bachelor's degree in a relevant field, such as industrial engineering, manufacturing engineering, or production management. To become a Chief Operating Officer, you will need at least a Master's degree in a relevant field, such as business administration, management, or marketing. To become a plant/manufacturing/production manager, you will need at least a PhD in a relevant field, such as industrial engineering, manufacturing engineering, or production management. Additionally, to become a plant/manufacturing/production manager, you may also want to obtain a certification such as the Certified Production Manager (CPM) or Certified Quality Manager (CQM) to demonstrate your expertise in the field.

Output (Llama3.2-1B):

To become a plant/manufacturing/production manager and eventually a Chief Operating Officer, it is recommended to have a degree in engineering, business, or a related field, as well as relevant post-graduate education and experience. Examples of post-graduate education include an MBA, a Master's degree in engineering, or even a PhD in a related field. Relevant experience could include working in a manufacturing or production role, managing a team of engineers, or leading a plant or production department.

Output (ShearedLlama-1.3B):

To become a plant/manufacturing/production manager, you will need to have a bachelor's degree in a relevant field such as engineering, business, or manufacturing. After that, you can pursue a master's degree in a related field such as operations management, manufacturing engineering, or production and operations management. Additionally, you may want to pursue post-graduate education in areas such as project management, lean manufacturing, or supply chain management.

For a Chief Operating Officer position, you will need to have a bachelor's degree in a relevant field such as engineering, business, or manufacturing. After that, you can pursue a master's degree in a related field such as operations management, manufacturing engineering, or production and operations management. Additionally, you may want to pursue post-graduate education in areas such as project management, lean manufacturing, or supply chain management.

Overall, to become a plant/manufacturing/production manager or Chief Operating Officer, you will need to have a bachelor's degree in a relevant field such as engineering, business, or manufacturing, and then pursue a master's degree in a related field such as operations management, manufacturing engineering, or production and operations management. Additionally, you may want to pursue post-graduate education in areas such as project management, lean manufacturing, or supply chain management.

Input:

Why so many people like film photography when digital camera and mobile phone are more convenient and generate higher quality of photos?

Output (EfficientLLM-1.1B):

Many people still prefer film photography because of its unique aesthetic qualities. Film photography allows for more creative control, as the photographer can adjust the ISO, shutter speed, and aperture to achieve the desired effect. Additionally, film photography produces a unique look and feel that cannot be replicated with digital cameras and mobile phones. Film photography also allows for more experimentation, as the photographer can adjust the ISO, shutter speed, and aperture to achieve the desired effect. Finally, film photography can be a more personal and intimate experience, as it allows the photographer to capture moments in a way that digital cameras and mobile phones cannot.

Output (OLMo-1B):

Film photography has many advantages compared to digital camera and mobile phone, such as higher quality of photos, more creative control, and the ability to capture moments that would otherwise go unnoticed. Additionally, film photography requires more time and patience to develop the film, which adds an element of adventure and creativity to the process. Finally, many people find the physical process of developing film to be aesthetically pleasing, while digital camera and mobile phone processing can often produce photos with a less-than-desirable outcome.

Output (ShearedLlama-1.3B):

People like film photography because it allows for more creative control. With film photography, the photographer has the ability to adjust the exposure, focus, and other settings to create the desired effect. Digital cameras and mobile phones do not offer this level of creative control, making it difficult to achieve the desired effect. Additionally, film photography requires more thought and preparation than digital camera and mobile phone photography, making it more enjoyable and rewarding.

Output (Llama3.2-1B):

People like film photography because it produces higher quality of photos than digital cameras and mobile phones. Film photography is more expensive than digital cameras and mobile phones, but the quality of photos produced by film is much better than those produced by digital cameras and mobile phones. Additionally, film photography requires more skill and patience to take good photos, which many people enjoy.