

# NutePrune: Efficient Progressive Pruning with Numerous Teachers for Large Language Models

Anonymous ACL submission

## Abstract

The considerable size of Large Language Models (LLMs) presents notable deployment challenges, particularly on resource-constrained hardware. Structured pruning, offers an effective means to compress LLMs, thereby reducing storage costs and enhancing inference speed for more efficient utilization. In this work, we study data-efficient and resource-efficient structure pruning methods to obtain smaller yet still powerful models. Knowledge Distillation is well-suited for pruning, as the intact model can serve as an excellent teacher for pruned students. However, it becomes challenging in the context of LLMs due to memory constraints. To address this, we propose an efficient progressive Numerous-teacher pruning method (NutePrune). NutePrune mitigates excessive memory costs by loading only one intact model and integrating it with various masks and LoRA modules, enabling it to seamlessly switch between teacher and student roles. This approach allows us to leverage numerous teachers with varying capacities to progressively guide the pruned model, enhancing overall performance. Extensive experiments across various tasks demonstrate the effectiveness of NutePrune. In LLaMA-7B zero-shot experiments, NutePrune retains 97.17% of the performance of the original model at 20% sparsity and 95.07% at 25% sparsity.

## 1 Introduction

Large Language Models (LLMs) excel in language tasks (OpenAI, 2023; Touvron et al., 2023; Thoppilan et al., 2022; Scao et al., 2022), but their substantial size poses deployment and inference challenges (Frantar et al., 2022). Techniques like model pruning (Molchanov et al., 2016), knowledge distillation (Jiao et al., 2019), and quantization (Dettmers et al., 2023) have been proposed to address computational demands. The exploration of LLM pruning, especially structured pruning (Frantar and Alistarh,

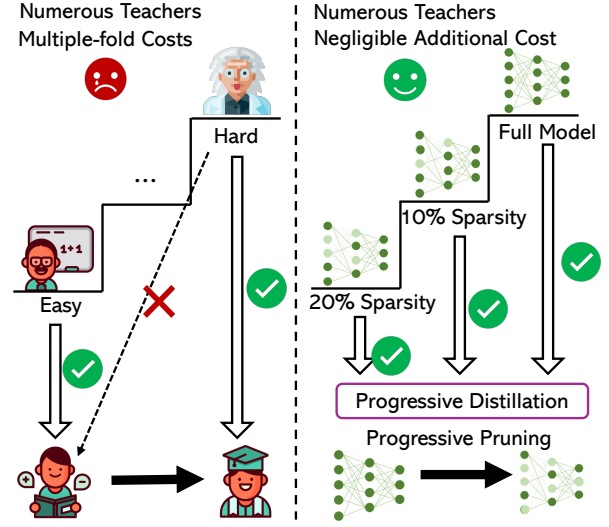


Figure 1: The advantage of our NutePrune. **Left:** Progressive distillation guides the student with teachers from easy to hard to avoid large capacity gap harming learning. But it suffers from multiple-fold costs of loading numerous teachers. **Right:** Our NutePrune leverages models with varying sparsity, enabling progressive distillation with negligible additional cost.

2023), holds great significance. Structured pruning reduces model size by removing coherent parameter groups, cutting inference costs on standard hardware. But it is more challenging than unstructured pruning in retaining the capabilities of LLMs (Hoeffler et al., 2021). Existing methods either adopt data-efficient approaches, causing a performance decline (Ma et al., 2023), or require extensive post-training to recover model performance (Xia et al., 2023). In this work, we investigate efficient methods to prune the model to higher sparsity without significant performance decline.

Knowledge distillation (KD) aims to train a more compact student model with supervision from a larger teacher model (Sanh et al., 2019; Gou et al., 2021). It's widely adopted and proven highly effective in the field of LLMs. Progressive learning,

utilizing intermediate teachers with a reduced gap in capabilities, has been demonstrated to improve performance in KD (Xiang et al., 2020). Previous work has shown that pruning with a distillation objective can improve performance (Xia et al., 2022). Distillation is particularly suitable for pruning since the full original model inherently serves as an excellent teacher for the pruned model (Sanh et al., 2020), which can offer a more detailed supervisory signal than conventional supervised training, enhancing the effectiveness of pruning with limited data (Lagunas et al., 2021).

However, applying this method in the realm of LLMs proves challenging. Given the vastness of an LLM, loading it onto GPUs consumes a substantial amount of memory. Introducing an additional teacher model requires twice the memory, making it impractical with limited memory resources. Furthermore, relying on a single teacher may not be the best practice (Liu et al., 2020; Wu et al., 2021). With the increasing gap of sparsity between teacher and student, the capacity gap is also widening, which toughens distillation. Employing multiple teachers with varying capacities can enhance the transfer of knowledge to students (Yuan et al., 2021). However, when it comes to the distillation of LLMs, memory consumption of multiple teachers becomes an even more pressing concern.

In this paper, we address the above challenges with an efficient progressive **Numerous-teacher** pruning method (NutePrune). Our motivation is demonstrated in Figure 1. NutePrune aims to diminish the capacity gap between the full teacher model and the highly sparse student, thereby alleviating the difficulty of distillation (Su et al., 2021; Mukherjee et al., 2023; Xiang et al., 2020). Instead of relying solely on a single full teacher, we instruct the student with many teachers with varying sparsity. To achieve this, we formulate pruning as an optimization problem where we learn masks to prune sub-modules while updating model parameters through LoRA (Hu et al., 2021). Specially, we load an intact model, serving dual roles as both a teacher and a student. In teacher mode, we incorporate the original model with collected frozen low-sparsity masks and corresponding LoRA modules. And in student mode, we incorporate it with learnable high-sparsity masks and LoRA modules. Since the masks and LoRA modules are highly parameter efficient, we collect and leverage numerous modules with different sparsity to incorporate numerous teachers and progressively prune the stu-

dent. Our contributions can be summarized as follows:

- We propose a novel distillation method that progressively guide the student using numerous teachers with varying sparsity to narrow the capacity gap. Through progressive KD, we achieve higher model sparsity without significant performance decline on limited data.
- Our NutePrune only loads one intact model and switch it between teacher and student modes by incorporating various masks and LoRA modules. This novel efficient distilling method for pruning enables using numerous teachers and introduces no extra memory cost, which is especially critical for LLMs.
- Extensive experiments across perplexity metric, commonsense reasoning, and MMLU demonstrate the effectiveness of our method.

## 2 Related Works

**Pruning for LLMs** For LLMs, SparseGPT (Frantar and Alistarh, 2023) and WANDA (Sun et al., 2023) employ unstructured pruning methods, while N:M sparsity (Zhou et al., 2021) is considered semi-structured. Despite the effectiveness of these methods, their intricate structures do not yield significant inference speedup on standard hardware (Frantar and Alistarh, 2023). Consequently, structured pruning has emerged as a recent consensus. CoFi (Xia et al., 2022) and nn pruning (Lagunas et al., 2021) are proposed for smaller language models like BERT (Devlin et al., 2018), often designed for specific tasks. CoFi loads both the teacher and student models, which is impractical for LLMs. Sheared-LLaMA (Xia et al., 2023) proposes pruning LLMs using a dynamic pre-training method, enhancing performance through extensive data and training resources.

However, concerns persist regarding limited memory and training resources for LLMs. In a pioneering effort, LLM-Pruner (Ma et al., 2023) prunes LLMs in one-shot and utilizes LoRA (Hu et al., 2021) for fine-tuning. LoRAPrune (Zhang et al., 2023) employs iterative pruning, replacing gradients on full weights with gradients on LoRA to calculate group importance. Compresso (Guo et al., 2023) leverages LoRA and elaborately designed prompts for training and inference. Meanwhile, LoRAShear (Chen et al., 2023) employs

LoRA and a dynamic fine-tuning scheme to recover knowledge.

**Knowledge Distillation (KD) for LLMs** KD (Hinton et al., 2015) has emerged as a vital technique to reduce inference costs while maintaining performance quality in the context of LLMs. Prior work of KD (Taori et al., 2023; Fu et al., 2023) mostly focus on black-box KD, using teacher’s generations to fine-tune the student. With the rise of open-source LLMs (Zhang et al., 2022; Touvron et al., 2023), interest in white-box KD is growing. White-box KD, leveraging teacher weights and logits, provides richer supervision signals, enhancing language abilities (Agarwal et al., 2023; Gu et al., 2023; Wen et al., 2023). Despite progress on small language models, significant performance gaps between large and small models persist (Achiam et al., 2023; Anil et al., 2023).

Progressive knowledge distillation (Xiang et al., 2020) has proven effective by using intermediate teachers to bridge the capacity gap with LLMs, especially in scenarios reliant on data generated by multiple teachers (Mukherjee et al., 2023). Orca (Mukherjee et al., 2023) first learns from easier examples from ChatGPT and then from harder ones from GPT-4, enhancing performance for smaller students in KD. However, applying white-box KD to LLMs poses challenges due to substantial memory requirements for loading both teacher and student models. This challenge becomes even more difficult when attempting to load multiple teachers.

### 3 Methodology

In this section, we first introduce how our NutePrune enables efficient knowledge distillation for structured pruning in 3.1. Then, to narrow capacity gap during distillation, we introduce the progressive knowledge distillation method that collects and incorporates numerous teachers in 3.2. The overview framework is illustrated in Figure 2.

#### 3.1 Efficient Distillation for Structured Pruning

We formulate structure pruning as a constrained optimization problem where we optimize masks and model to prune the structure to a target sparsity while maximizing performance. We use distillation loss instead of original language model loss and use  $L_0$  regularization to control the pruned sparsity. To mitigate substantial memory consumption during LLM training, we utilize LoRA for model updates,

making pruning the process of training these masks and LoRA parameters.

#### Learning masks to control the pruned structure

We allow for three types of structured pruning: attention heads, FFN intermediate dimensions, and hidden dimensions. It is achieved by learning corresponding masks  $\mathbf{z}_{head}, \mathbf{z}_{int}, \mathbf{z}_{hid} \in \{0, 1\}$ . Formally, the multi-head attention module  $\text{MHA}(x)$  and feed-forward networks  $\text{FFN}(x)$  of layer  $l$  are pruned as:

$$\text{MHA}^l(X) = \mathbf{z}_{hid} \cdot \sum_{h=1}^{N_{head}} \mathbf{z}_{head}^{l,h} \text{Att}^{l,h}(X). \quad (1)$$

$$\begin{aligned} \text{FFN}^l(X) &= \mathbf{z}_{hid} \\ &\cdot W_D^l \left( \mathbf{z}_{int}^l \cdot W_U^l(X) \cdot \text{Act}(W_G^l(X)) \right) \end{aligned} \quad (2)$$

where  $\text{Att}()$  is the attention module and  $\text{Act}()$  is the activation.  $W_D, W_U, W_G$  are down projection, up projection, and gating projection.

During mask training, we calculate the remaining size to obtain the expected sparsity  $\hat{s}$ :

$$\begin{aligned} \hat{s}(\mathbf{z}) &= \frac{1}{M} \cdot 4 \cdot d_h \cdot \sum_l^L \sum_h^{N_{head}} \sum_k^d \mathbf{z}_{head}^{l,h} \mathbf{z}_{hid}^k \\ &+ \frac{1}{M} \cdot 3 \cdot \sum_l^L \sum_i^{d_{int}} \sum_k^d \mathbf{z}_{int}^{l,i} \mathbf{z}_{hid}^k, \end{aligned} \quad (3)$$

where  $M$  denotes full model size.  $L$  is number of layers.  $d_h, N_{head}, d, d_{int}$  are head dimension, number of head, hidden dimension, and intermediate dimension, correspondingly.

All masking variables are learned as real numbers in  $[0, 1]$  during training. We follow (Louizos et al., 2017; Guo et al., 2023) and employ the augmented  $L_0$  regularization:

$$\begin{aligned} \mathbf{u} &\sim U(0, 1) \\ \mathbf{s} &= \text{sigmoid} \left( \frac{1}{\beta} \log \frac{\mathbf{u}}{1 - \mathbf{u}} + \log \alpha \right) \\ \tilde{\mathbf{s}} &= \mathbf{s} \times (r - l) + l \\ \mathbf{z} &= \min(1, \max(0, \tilde{\mathbf{s}})), \end{aligned} \quad (4)$$

where  $\mathbf{u}$  is uniformly sampled between 0 to 1.  $\alpha$  is the parameter to be learned and  $\beta$  is a hyperparameter.  $l, r$  is often  $-0.1$  and  $1.1$  to ensure most  $\mathbf{z}$  are either 0 or 1 after training.

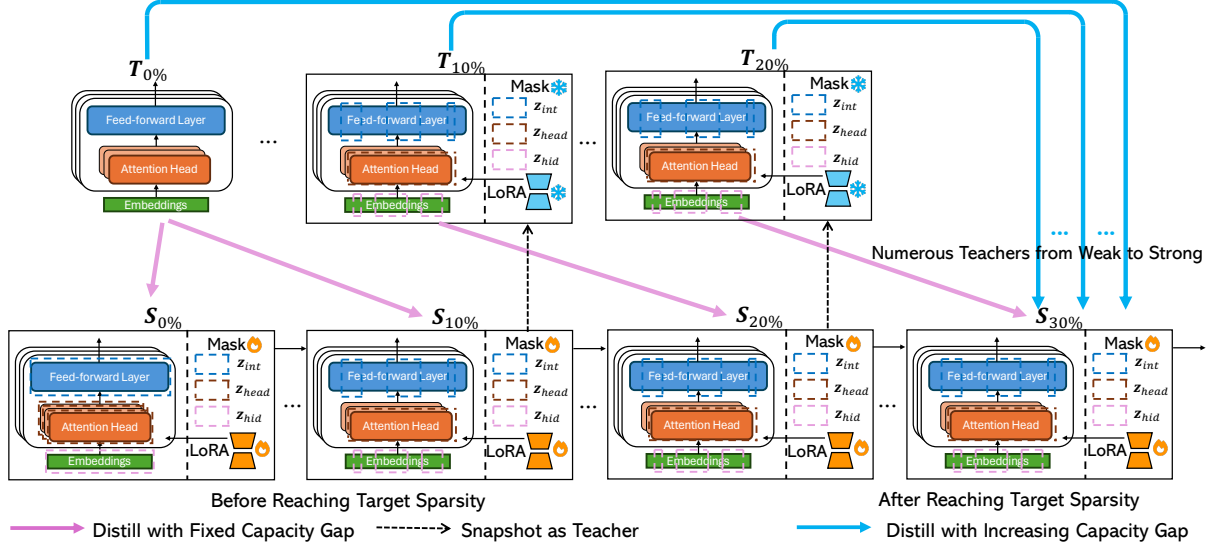


Figure 2: The overall framework of NutePrune. The pruned model is frozen and incorporated with learnable masks and LoRA. During pruning, the model is guided by numerous teachers. Before pruned to the target sparsity (e.g. 30%), it learns from teachers with a fixed capacity gap. Once the target sparsity is achieved, it continues to learn from all previous teachers from weak to strong. All these teachers are derived from snapshots of the student model itself. Since only the mask and LoRA modules are snapshotted, the additional memory cost is negligible.

To prevent models from drastically converging to different sizes, we follow (Wang et al., 2019) to use this Lagrangian term:

$$\mathcal{L}_0 = \lambda_1 \cdot (\hat{s} - t) + \lambda_2 \cdot (\hat{s} - t)^2, \quad (5)$$

where  $\lambda_1$  and  $\lambda_2$  are both learnable. This loss term  $\mathcal{L}_0$  will impose  $\hat{s}$  to gradually converge to target sparsity  $t$ .

**Updating parameters with LoRA** Considering massive memory usage during full fine-tuning for LLMs, we incorporate lightweight LoRA (Hu et al., 2021) modules into LLM weights to update parameters during pruning.

An incorporated module  $W'$  is consisted of the original weight  $W : \mathbb{R}^n \rightarrow \mathbb{R}^m$  and sequential LoRA weights parallel to  $W$ :

$$W'(X) = W(X) + W_B(W_A(X)), \quad (6)$$

where  $W_A : \mathbb{R}^n \rightarrow \mathbb{R}^r$ ,  $W_B : \mathbb{R}^r \rightarrow \mathbb{R}^m$  and  $r \ll m, n$ . During training,  $W$  is frozen and only  $W_A$  and  $W_B$  are learnable.

**Efficient distillation** Instead of simultaneously loading two massive models into memory, we propose to incorporate the frozen and intact model  $M$  with different lightweight masks and LoRA modules for the teacher and the student. Formally, let  $\mathbf{I} = \{\mathbf{z}, \mathbf{W}_A, \mathbf{W}_B\}$  denotes the set of all masks

and LoRA modules which is highly parameter efficient ( $|\mathbf{I}| \ll |\mathbf{M}|$ ). By incorporating  $\mathbf{I}$  into  $\mathbf{M}$ , we obtain  $\mathbf{M}_{\mathbf{I}}$ . The objective of knowledge distillation is the KL-divergence (Van Erven and Harremos, 2014) between teacher’s and student’s output probability distributions  $p$ :

$$\mathcal{L}_{KL} = D_{KL}(p(\mathbf{M}_{\mathbf{I}_S}, x), p(\mathbf{M}_{\mathbf{I}_T}, x)), \quad (7)$$

where  $x$  denotes training data.  $\mathbf{I}_S$  and  $\mathbf{I}_T$  denote the lightweight modules of student and teacher.

Additionally, intermediate layers of a teacher model can serve as effective targets for training a student model (Chen et al., 2021). This objective can be formulated as:

$$\mathcal{L}_{layer} = \sum_l^L \text{MSE}(\mathbf{h}_l(\mathbf{M}_{\mathbf{I}_S}, x), \mathbf{h}_l(\mathbf{M}_{\mathbf{I}_T}, x)), \quad (8)$$

where  $\mathbf{h}_l$  is the hidden embedding of the  $l$ -th layer. Therefore, the overall objective is:

$$\mathcal{L} = \mathcal{L}_{KL} + \alpha_1 \mathcal{L}_{layer} + \alpha_2 \mathcal{L}_0, \quad (9)$$

where  $\alpha_1, \alpha_2$  are hyperparameters to control the importance of different loss terms.

### 3.2 Progressive Knowledge Distillation with Numerous Teachers

All teachers are collected from the snapshot of students as the dotted line illustrated in Figure 2. To



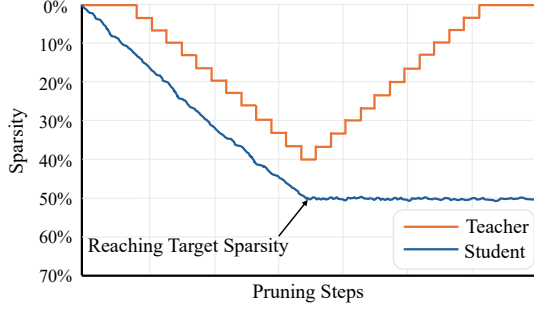


Figure 3: Illustration of the sparsity of teacher and student models during pruning. Take the example with the target sparsity  $t = 50\%$  and sparsity gap  $g = 10\%$ .

narrow the capacity gap between the intact teacher and high sparsity students, we leverage a novel progressive knowledge distillation (PKD) method for pruning. It consists of two stages when pruning a model from 0% sparsity as illustrated in Figure 3.

**Before reaching target sparsity** The sparsity of pruned model gradually increase from 0 to  $t$ . To narrow the sparsity gap, we set a fixed gap value  $g$  and make the pruned model  $S$  guided by teachers  $T$  whose sparsity  $\hat{s}(T)$  is approximately  $g$  less than  $\hat{s}(S)$ :  $\hat{s}(T) = \hat{s}(S) - g$ . These teachers are snapshots of previous students. The original intact model serves as the teacher for student  $\hat{s}(S) < g$ .

To avoid collecting too many teachers, we only collect teachers with an interval of  $i$ . Therefore, for any teacher with sparsity  $\hat{s}(T)$ , it is responsible for guiding a student set within a range of sparsity. We use  $\rightarrow$  to denote the relationship in which a teacher distills knowledge to students.

$$T \rightarrow \{S | \hat{s}(T) + g < \hat{s}(S) < \hat{s}(T) + g + i\}. \quad (10)$$

And the intact model  $\mathbf{M} = T_0$  is responsible for the early students whose sparsity is less than  $g + i$ :

$$T_0 \rightarrow \{S | \hat{s}(S) < g + i\}. \quad (11)$$

**After reaching target sparsity** When the pruned model reaches the target sparsity  $t$ , we proceed to the second stage of PKD. The model undergoes distillation by all preceding teachers, with a reduction of sparsity in the teachers. This gradual process guides the model’s learning trajectory from weaker to stronger knowledge and from easier to more challenging concepts. Throughout this stage, the sparsity of the pruned model  $\hat{s}$  remains close to the target sparsity  $t$ , while the masks  $\mathbf{z}$  and LoRA models  $\mathbf{W}_A, \mathbf{W}_B$  are continually optimized.

To receive sufficient instruction from the best model (the intact model  $\mathbf{M}$ ), the teacher model is maintained as  $\mathbf{M}$  during the final period.

### 3.3 Post Fine-tuning

After the pruning phase, to obtain better performance, we undergo a post fine-tuning stage following LLM-Pruner (Ma et al., 2023). We fix the masks and only fine-tune LoRA modules on the Stanford Alpaca (Taori et al., 2023) dataset.

## 4 Experiments

### 4.1 Experimental Setup

**Datasets** To assess the zero-shot ability of LLMs, we perform zero-shot classification tasks on seven commonsense reasoning benchmarks: BoolQ (Clark et al., 2019), PIQA (Bisk et al., 2020), Hel-laSwag (Zellers et al., 2019), WinoGrande (Sakaguchi et al., 2021), ARC-easy (Clark et al., 2018), ARC-challenge (Clark et al., 2018), and OpenBookQA (OBQA) (Mihaylov et al., 2018). We evaluate the general capcability of LLMs on the perplexity metric with WikiText (Merity et al., 2016) and PTB (Marcus et al., 1993) dataset. Additionally, We report the results on 5-shot MMLU to evaluate the in-context learning ability (Hendrycks et al., 2020), which consists of 57 tasks covering STEM, humanities, social science, etc.

**Models** We assess the performance of NutePrune on the LLaMA-1 family (Thoppilan et al., 2022), comprising LLaMA-7B, LLaMA-13B. Our experiments primarily center on pruning the LLaMA-7B model to provide a comprehensive comparison with previous studies. We primarily evaluate at sparsity levels of 20% and 50%.

**Baselines** Numerous works delve into pruning techniques for LLMs, including SparseGPT (Frantar and Alistarh, 2023) and WANDA (Sun et al., 2023). Considering the benefits of inference acceleration, we mainly focus on structure pruning methods. For open-source methods as far as we know, we implement LLM-Pruner (Ma et al., 2023) and Compresso (Guo et al., 2023) and conduct detailed comparisons with our NutePrune. For more recent works that are not publicly available, We assess NutePrune using the same settings as theirs for comparison with their reported results. This includes LoRAPrune (Zhang et al., 2023) and LoRAShear (Chen et al., 2023).

Ratio	Tune	Method	BoolQ	PIQA	HellaSwag	WinoGrande	ARC-e	ARC-c	OBQA	Avg.	★Avg.
0%		LLaMA-7B	73.18	78.35	72.99	67.01	67.45	41.38	42.40	63.25	66.39
20%		LLM-Pruner	59.39	75.57	65.34	61.33	59.18	37.12	39.80	56.82	59.01
		LoRAPrune	57.98	75.11	65.81	59.90	62.14	34.59	39.98	56.50	-
		†NutePrune	63.21	76.55	67.96	66.69	63.72	38.05	40.00	<b>59.46</b>	<b>63.03</b>
20%	✓	WANDA	65.75	74.70	64.52	59.35	60.65	36.26	39.40	57.23	-
		LLM-Pruner	69.54	76.44	68.11	65.11	63.43	37.88	40.00	60.07	61.94
		LoRAPrune	65.82	79.31	70.00	62.76	65.87	37.69	39.14	60.05	-
		LoRAShear	70.17	76.89	68.69	65.83	64.11	38.77	39.97	60.63	-
		Compresso	73.64	75.08	64.77	67.72	66.12	37.54	40.40	60.75	62.60
		‡NutePrune	72.69	76.71	68.99	65.51	65.49	38.48	40.20	61.15	63.57
		NutePrune	74.56	77.04	70.01	65.67	65.78	37.97	39.20	<b>61.46</b>	<b>64.39</b>
25%		†NutePrune	68.10	75.35	66.75	62.04	58.08	36.77	39.00	58.01	61.72
25%	✓	‡NutePrune	65.84	76.17	66.69	64.56	61.49	37.03	39.20	58.71	63.12
		NutePrune	68.99	77.20	67.90	65.04	63.76	37.80	40.20	60.13	63.78
50%		LLM-Pruner	52.32	59.63	35.64	53.20	33.50	27.22	33.40	42.13	40.94
		LoRAPrune	51.78	56.90	36.76	53.80	33.82	26.93	33.10	41.87	-
		†NutePrune	62.29	67.95	53.03	57.06	45.45	30.03	36.60	<b>50.35</b>	<b>53.14</b>
50%	✓	WANDA	50.90	57.38	38.12	55.98	42.68	34.20	38.78	45.43	-
		LLM-Pruner	61.47	68.82	47.56	55.09	46.46	28.24	35.20	48.98	48.97
		LoRAPrune	61.88	71.53	47.86	55.01	45.13	31.62	34.98	49.71	-
		LoRAShear	62.12	71.80	48.01	56.29	47.68	32.26	34.61	50.39	-
		Compresso	60.09	66.70	39.31	51.93	48.82	27.82	33.40	46.87	47.43
		‡NutePrune	62.20	69.91	53.87	57.77	46.59	31.74	35.8	51.13	53.94
		NutePrune	62.26	71.00	55.88	57.54	51.68	32.17	34.40	<b>52.13</b>	<b>54.91</b>

† only prunes the model by training masks without incorporating LoRA modules.

‡ prunes the model by co-training the masks and LoRA modules but without post fine-tuning on Alpaca.

★ includes our reproduced results with the newer version of *lm-evaluation-harness* that fixes a LLaMA tokenization issue. See Appendix A for detailed results.

Table 1: Zero-shot performance (%) of the compressed LLaMA models. **Bold** denotes the best average performance at the same setting.

**Implementation details** For pruning stage, we sample 20,000 sentences from the C4 (Raffel et al., 2020) dataset with a length of 512 tokens. We train with AdamW optimizer, a batch size of 16, and learning rates of 0.1 for masks and 0.001 for LoRA. We prune the model for 7 epochs. We use a linear sparsity schedule to gradually increase the target sparsity from 0 to the target ratio. We set sparsity warmup to 4 epochs for pruning of 20% sparsity and 1 epoch for 50%. The sparsity gap between the teacher and student  $g$  is 10% and the snapshot interval  $i$  of teachers is 1%. After pruning, we post fine-tune the pruned model on the Alpaca dataset (Taori et al., 2023) for 3 epochs. All experiments are conducted on one A100 GPU (80G).

## 4.2 Results

**Zero-shot performance** Table 1 demonstrates the zero-shot performance on commonsense reasoning tasks for compressed LLaMA models. We mainly present results with the previous version of *lm-evaluation-harness* for comprehensive comparison, which is widely employed in prior studies (Ma et al., 2023). Additionally, we present results with a newer version that addressed a tokenization issue

in LLaMA for more accurate evaluation, where we report our reproduced results of publicly available methods. The reported results include experiments for 20%, 25% and 50% sparsity levels, covering scenarios with and without parameter tuning.

The average performance of NutePrune consistently outperforms previous methods across all settings. For pruning without tuning, NutePrune outperforms LLM-Pruner by 2.64%/8.22% at 20%/50% sparsity, underscoring its ability to derive a more effective pruned structure compared to other methods. For pruning with LoRA constrained, NutePrune improves from 59.46%/50.35% to 61.46%/52.13% at 20%/50% sparsity, indicating co-training with LoRA could help recover model capability damaged by pruning. And with additional post fine-tuning on Alpaca, notably, it retains 97.17% of the performance of the original model at 20% sparsity and 95.07% at 25% sparsity. NutePrune exhibits more noticeable improvements at higher sparsity levels, proving the effectiveness of our PKD in mitigating the capacity gap between teacher and student. Specifically, NutePrune achieves comparable performance to other methods at 20% sparsity when operating at

25% sparsity. And even without parameter tuning, NutePrune achieves results comparable to other methods at 50% sparsity. These findings suggest that NutePrune is particularly effective in compressing models to higher sparsity levels.

**Perplexity** The perplexity (PPL) is a metric reflecting the generation ability of LLMs. The results on WikiText2 and PTB are presented in Table 2. Compresso prunes on an instruction dataset and depends on a meticulously crafted instructing prompt, enabling it to outperform LLM-Pruner on zero-shot tasks at low sparsity but yielding inferior performance on PPL, especially at high sparsity. In contrast, NutePrune utilizes a more general 20k C4 dataset, preserving overall language ability. This distinction allows NutePrune to avoid degradation of generalization and achieve superior performance in terms of PPL compared to both Compresso and LLM-Pruner.

Ratio	Tune	Method	WikiText2↓	PTB↓
0%		LLaMA-7B	5.68	8.81
20%		LLM-Pruner	9.96	15.61
		†NutePrune	<b>8.02</b>	<b>14.41</b>
20%	✓	LLM-Pruner	8.57	12.84
		Compresso	10.38	16.14
		NutePrune	<b>7.65</b>	<b>12.40</b>
25%		†NutePrune	9.04	14.35
25%	✓	NutePrune	7.85	13.10
50%		LLM-Pruner	98.1	224.54
		†NutePrune	<b>17.45</b>	<b>34.77</b>
50%	✓	LLM-Pruner	22.76	33.9
		Compresso	59.73	80.38
		NutePrune	<b>13.20</b>	<b>21.97</b>

Table 2: Perplexity metric on WikiText2 and PTB.

**Pruning of larger model** In addition to evaluating LLaMA-7B, we assess the larger LLaMA-13B with 20% sparsity. As demonstrated in Table 3, our approach yield an average zero-shot performance of 67.51%, which is only 0.12% lower than the full model and 1.75% higher than LLM-Pruner. It also outperforms LLM-Pruner in terms of PPL. Besides, in-context learning is a crucial capability for LLMs (Brown et al., 2020), particularly for larger language models such as LLaMA-13B. To assess it, we conduct evaluations using the MMLU with a 5-shot setting. Results demonstrates that NutePrune can also perform well on MMLU metrics with larger models.

**Inference latency** We test the inference latency by generating from 64 tokens to 256 tokens on

Method	Zero-Shot (%)	WikiText2↓	MMLU
0% Baseline	67.63	5.62	0.426
LLM-Pruner	65.76	6.95	0.351
NutePrune	<b>67.51</b>	<b>6.55</b>	<b>0.355</b>

Table 3: Performance of the compressed LLaMA-13B models with 20% sparsity.

vLLM (Kwon et al., 2023), which is a fast and widely deployed library for LLM inference and serving. The results are presented in Table 4. NutePrune achieves latency savings of 11% and 29% at sparsity levels of 20% and 50%, respectively. While LLM-Pruner can save slightly more latency due to its predefined neater structure, it comes at the cost of reduced flexibility in tailoring. However, as sparsity increases, the difference becomes negligible.

Method	20%	50%
0% Baseline	3.06	
LLM-Pruner	2.63(-14%)	2.17(-29%)
NutePrune	2.72(-11%)	2.18(-29%)

Table 4: Inference latency of pruned LLaMA-7B.

**Training cost** We report the memory and latency cost on different settings in Table 5. For extra GPU memory cost of PKD, NutePrune snapshot lightweight modules (masks and LoRA) of numerous teachers into CPU. Only one teacher module is loaded onto the GPU when needed, resulting in negligible memory cost compared with KD. In terms of extra time cost, compared with supervised training, KD requires one extra forward pass of teacher model, which is inevitable and cost 18.0% extra latency. When snapshotting a teacher or switching to a new teacher, due to the extremely low frequency of operations, the time can be ignored. Introducing  $\mathcal{L}_{layer}$  requires additional 32% memory which is also efficient compared to conventional KD.

Progressive	KD	$\mathcal{L}_{layer}$	Memory	Latency
			27.68	3.67
	✓		28.67	4.33
✓	✓		28.69	4.33
✓	✓	✓	38.00	5.52

Table 5: Training cost measured by average GPU memory (GB) and per step latency (s/iter).

### 4.3 Ablation Study

We validate the effectiveness of NutePrune and investigate which properties make for a good NutePrune. Results are average zero-shot performance with tuning but without post fine-tuning, unless otherwise stated.

**Effectiveness of PKD** To validate progressive knowledge distillation (PKD) in our NutePrune, we conduct ablation studies on various learning strategies. We eliminate the progressive schedule and adopt standard KD, where the intact model serves as the teacher throughout. Subsequently, we exclude the entire distillation procedure and employ the standard generative language model loss, specifically next-token prediction, to train masks and LoRA modules. The results presented in Table 6 demonstrate the critical role of KD in enhancing performance, with further improvements achieved through PKD. This phenomenon is particularly pronounced at higher sparsity.

Progressive	KD	20%	50%
✓	✓	<b>63.57</b>	<b>53.94</b>
	✓	63.19	52.73
		59.98	41.77

Table 6: NutePrune and variants at 20%/50% sparsity.

**Two stages of PKD** PKD includes one stage before reaching target sparsity and the other stage after that. Different progressive schedules are adopted. To assess the effectiveness of them, we conducted an ablation study at 50% sparsity under two training settings, as shown in Table 7: training masks only and co-training masks with LoRA. In a stage without a progressive schedule, the intact model serves as the teacher. For the masks-only scenario, adopting either stage 1 or 2 alone yields significant improvements over KD. And for co-training, significant improvement is observed when both stages are adopted simultaneously.

**Sparsity gap between teacher and student** During stage 1, the sparsity gap between teacher and student model is an important hyperparameter. We conduct experiments with various sparsity gaps and the results are presented in Table 8. A 10% gap is deemed appropriate to prevent a gap that is too small, as it may result in insufficient guidance, or a gap that is too large, which would also toughing distillation.

Stage 1	Stage 2	Avg.(%)	
		masks-only	co-train
✓	✓	<b>53.14</b>	<b>53.94</b>
✓		52.31	52.79
	✓	52.40	52.53
		51.83	52.73

Table 7: Performance of two stages of PKD.

Sparsity Gap	Avg.(%)
5%	53.62
10% (ours)	<b>53.94</b>
20%	53.04

Table 8: Performance of various sparsity gap.

**Snapshot interval of teachers** When taking snapshots of students as teachers, it is preferable to save as many teachers as possible to facilitate more comprehensive training. However, the additional CPU storage for these teachers incurs extra costs. As demonstrated in Table 9, selecting an interval of 1% leads to significant improvement over the 10% interval, and the associated extra storage is acceptable, which means we use 10 teachers for 20% sparsity and 40 teachers for 50% sparsity.

Snapshot Interval	CPU Storage	Avg.(%)
1% (ours)	728MB	<b>53.94</b>
10%	73MB	53.27

Table 9: Storage and performance of various intervals.

## 5 Conclusion

In this work, we propose NutePrune as a novel efficient progressive structured pruning method for LLMs. By minimizing the cost of KD, NutePrune aligns closely with the memory and time requirements of conventional training, allowing the utilization of numerous teachers for distillation. This approach mitigates the capacity gap between teacher and student, preserving the quality of the distillation process. We showcase the effectiveness of NutePrune, particularly at high sparsity levels, across various tasks and metrics for LLMs. The findings of this work contribute to the advancement of structured pruning techniques for LLMs, particularly in resource-constrained scenarios.



## 6 Limitations

We acknowledge the main limitation of this work is that we only evaluate our methods on LLaMA with 7B and 13B due to limited computation resources. And other model families, such as LLaMA-v2, are not included due to the absence of results for comparison from state-of-the-art methods. Recent work (Ma et al., 2023; Xia et al., 2023) proves that using extensive data for post-training could substantially enhance the performance, but it comes with a substantial increase in computational costs. We target on pruning on resource-constraint scenarios and leave pruning with extensive data for future work.

## References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Rishabh Agarwal, Nino Vieillard, Piotr Stanczyk, Sabela Ramos, Matthieu Geist, and Olivier Bachem. 2023. Gkd: Generalized knowledge distillation for auto-regressive sequence models. *arXiv preprint arXiv:2306.13649*.
- Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. 2023. Palm 2 technical report. *arXiv preprint arXiv:2305.10403*.
- Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. 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, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Defang Chen, Jian-Ping Mei, Yuan Zhang, Can Wang, Zhe Wang, Yan Feng, and Chun Chen. 2021. Cross-layer distillation with semantic calibration. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 7028–7036.
- Tianyi Chen, Tianyu Ding, Badal Yadav, Ilya Zharkov, and Luming Liang. 2023. Lorashear: Efficient large language model structured pruning and knowledge recovery. *arXiv preprint arXiv:2310.18356*.
- 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*.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms. *arXiv preprint arXiv:2305.14314*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- 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.
- Elias Frantar, Saleh Ashkboos, Torsten Hoefer, and Dan Alistarh. 2022. Gptq: Accurate post-training quantization for generative pre-trained transformers. *arXiv preprint arXiv:2210.17323*.
- Yao Fu, Hao Peng, Litu Ou, Ashish Sabharwal, and Tushar Khot. 2023. Specializing smaller language models towards multi-step reasoning. *arXiv preprint arXiv:2301.12726*.
- Jianping Gou, Baosheng Yu, Stephen J Maybank, and Dacheng Tao. 2021. Knowledge distillation: A survey. *International Journal of Computer Vision*, 129:1789–1819.
- Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang. 2023. Knowledge distillation of large language models. *arXiv preprint arXiv:2306.08543*.
- Song Guo, Jiahang Xu, Li Lyna Zhang, and Mao Yang. 2023. Compresso: Structured pruning with collaborative prompting learns compact large language models. *arXiv preprint arXiv:2310.05015*.
- 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*.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*.
- Torsten Hoefer, Dan Alistarh, Tal Ben-Nun, Nikoli Dryden, and Alexandra Peste. 2021. Sparsity in deep learning: Pruning and growth for efficient inference and training in neural networks. *The Journal of Machine Learning Research*, 22(1):10882–11005.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.

650	Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao	Victor Sanh, Lysandre Debut, Julien Chaumond, and	704
651	Chen, Linlin Li, Fang Wang, and Qun Liu. 2019.	Thomas Wolf. 2019. Distilbert, a distilled version	705
652	Tinybert: Distilling bert for natural language under-	of bert: smaller, faster, cheaper and lighter. <i>arXiv</i>	706
653	standing. <i>arXiv preprint arXiv:1909.10351</i> .	<i>preprint arXiv:1910.01108</i> .	707
654	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying	Victor Sanh, Thomas Wolf, and Alexander Rush. 2020.	708
655	Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E.	Movement pruning: Adaptive sparsity by fine-tuning.	709
656	Gonzalez, Hao Zhang, and Ion Stoica. 2023. Effi-	<i>Advances in Neural Information Processing Systems</i> ,	710
657	cient memory management for large language model	33:20378–20389.	711
658	serving with pagedattention. In <i>Proceedings of the</i>		
659	<i>ACM SIGOPS 29th Symposium on Operating Systems</i>	Teven Le Scao, Angela Fan, Christopher Akiki, El-	712
660	<i>Principles</i> .	lie Pavlick, Suzana Ilić, Daniel Hesslow, Roman	713
661	François Lagunas, Ella Charlaix, Victor Sanh, and	Castagné, Alexandra Sasha Luccioni, François Yvon,	714
662	Alexander M Rush. 2021. Block pruning for faster	Matthias Gallé, et al. 2022. Bloom: A 176b-	715
663	transformers. <i>arXiv preprint arXiv:2109.04838</i> .	parameter open-access multilingual language model.	716
664	Yuang Liu, Wei Zhang, and Jun Wang. 2020. Adap-	<i>arXiv preprint arXiv:2211.05100</i> .	717
665	tive multi-teacher multi-level knowledge distillation.	Weiyue Su, Xuyi Chen, Shikun Feng, Jiaxiang Liu,	718
666	<i>Neurocomputing</i> , 415:106–113.	Weixin Liu, Yu Sun, Hao Tian, Hua Wu, and Haifeng	719
667	Christos Louizos, Max Welling, and Diederik P Kingma.	Wang. 2021. Ernie-tiny: A progressive distillation	720
668	2017. Learning sparse neural networks through $l_0$	framework for pretrained transformer compression.	721
669	regularization. <i>arXiv preprint arXiv:1712.01312</i> .	<i>arXiv preprint arXiv:2106.02241</i> .	722
670	Xinyin Ma, Gongfan Fang, and Xinchao Wang. 2023.	Mingjie Sun, Zhuang Liu, Anna Bair, and J Zico	723
671	Llm-pruner: On the structural pruning of large lan-	Kolter. 2023. A simple and effective pruning ap-	724
672	guage models. <i>arXiv preprint arXiv:2305.11627</i> .	proach for large language models. <i>arXiv preprint</i>	725
673	Mitchell Marcus, Beatrice Santorini, and Mary Ann	<i>arXiv:2306.11695</i> .	726
674	Marcinkiewicz. 1993. Building a large annotated	Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann	727
675	corpus of english: The penn treebank.	Dubois, Xuechen Li, Carlos Guestrin, Percy Liang,	728
676	Stephen Merity, Caiming Xiong, James Bradbury, and	and Tatsunori B Hashimoto. 2023. Stanford alpaca:	729
677	Richard Socher. 2016. Pointer sentinel mixture mod-	An instruction-following llama model.	730
678	els. <i>arXiv preprint arXiv:1609.07843</i> .		
679	Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish	Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam	731
680	Sabharwal. 2018. Can a suit of armor conduct elec-	Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng,	732
681	tricity? a new dataset for open book question answer-	Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, et al.	733
682	ing. <i>arXiv preprint arXiv:1809.02789</i> .	2022. Lamda: Language models for dialog applica-	734
683	Pavlo Molchanov, Stephen Tyree, Tero Karras, Timo	tions. <i>arXiv preprint arXiv:2201.08239</i> .	735
684	Aila, and Jan Kautz. 2016. Pruning convolutional	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier	736
685	neural networks for resource efficient inference.	Martinet, Marie-Anne Lachaux, Timothée Lacroix,	737
686	<i>arXiv preprint arXiv:1611.06440</i> .	Baptiste Rozière, Naman Goyal, Eric Hambro,	738
687	Subhabrata Mukherjee, Arindam Mitra, Ganesh Jawa-	Faisal Azhar, et al. 2023. Llama: Open and effi-	739
688	har, Sahaj Agarwal, Hamid Palangi, and Ahmed	cient foundation language models. <i>arXiv preprint</i>	740
689	Awadallah. 2023. Orca: Progressive learning from	<i>arXiv:2302.13971</i> .	741
690	complex explanation traces of gpt-4. <i>arXiv preprint</i>	Tim Van Erven and Peter Harremos. 2014. Rényi diver-	742
691	<i>arXiv:2306.02707</i> .	gence and kullback-leibler divergence. <i>IEEE Trans-</i>	743
692	R OpenAI. 2023. Gpt-4 technical report. arxiv	<i>actions on Information Theory</i> , 60(7):3797–3820.	744
693	2303.08774. <i>View in Article</i> , 2:13.	Ziheng Wang, Jeremy Wohlwend, and Tao Lei. 2019.	745
694	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine	Structured pruning of large language models. <i>arXiv</i>	746
695	Lee, Sharan Narang, Michael Matena, Yanqi Zhou,	<i>preprint arXiv:1910.04732</i> .	747
696	Wei Li, and Peter J Liu. 2020. Exploring the limits	Yuqiao Wen, Zichao Li, Wenyu Du, and Lili Mou. 2023.	748
697	of transfer learning with a unified text-to-text trans-	f-divergence minimization for sequence-level knowl-	749
698	former. <i>The Journal of Machine Learning Research</i> ,	edge distillation. <i>arXiv preprint arXiv:2307.15190</i> .	750
699	21(1):5485–5551.	Chuhan Wu, Fangzhao Wu, and Yongfeng Huang. 2021.	751
700	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavat-	One teacher is enough? pre-trained language model	752
701	ula, and Yejin Choi. 2021. Winogrande: An adver-	distillation from multiple teachers. <i>arXiv preprint</i>	753
702	sarial winograd schema challenge at scale. <i>Commu-</i>	<i>arXiv:2106.01023</i> .	754
703	<i>nications of the ACM</i> , 64(9):99–106.	Mengzhou Xia, Tianyu Gao, Zhiyuan Zeng, and Danqi	755
		Chen. 2023. Sheared llama: Accelerating language	756
		model pre-training via structured pruning. <i>arXiv</i>	757
		<i>preprint arXiv:2310.06694</i> .	758

Mengzhou Xia, Zexuan Zhong, and Danqi Chen. 2022. Structured pruning learns compact and accurate models. *arXiv preprint arXiv:2204.00408*.

Liuyu Xiang, Guiguang Ding, and Jungong Han. 2020. Learning from multiple experts: Self-paced knowledge distillation for long-tailed classification. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part V 16*, pages 247–263. Springer.

Fei Yuan, Linjun Shou, Jian Pei, Wutao Lin, Ming Gong, Yan Fu, and Daxin Jiang. 2021. Reinforced multi-teacher selection for knowledge distillation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 14284–14291.

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, Chunhua Shen, Zhen Yang, Linlin Ou, Xinyi Yu, Bohan Zhuang, et al. 2023. Pruning meets low-rank parameter-efficient fine-tuning. *arXiv preprint arXiv:2305.18403*.

Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. 2022. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*.

Aojun Zhou, Yukun Ma, Junnan Zhu, Jianbo Liu, Zhi-jie Zhang, Kun Yuan, Wenxiu Sun, and Hongsheng Li. 2021. Learning n: m fine-grained structured sparse neural networks from scratch. *arXiv preprint arXiv:2102.04010*.

## A Zero-shot Performance with Newer Version

*lm-evaluation-harness* released a new version in June 2023 to assess the zero-shot performance of LLaMA<sup>1</sup>. This update addressed a tokenization bug specific to LLaMA, resulting in higher and more accurate performance results compared to the older version. Despite these improvements, current state-of-the-art reports continue to reference the older version. Consequently, we conducted experiments using both the new and old versions, and the detailed results for the new version are presented in Table 10.

## B Pruned Structure

To gain insights into the pruned model, we present a detailed overview of the pruned structure at sparsity levels of 20% and 50%. The original hidden

dimension is 4096, with a number of heads set at 32 and an intermediate dimension of 11008. Tables 11 and 12 reveal several observations. Notably, NutePrune tends to avoid pruning the hidden dimension, which aligns with the observation that pruning it may result in significant performance degradation (Ma et al., 2023). Regarding heads and intermediate dimensions, NutePrune tends to prune the the last few layers. This observation differs from LLM-Pruner, which asserts the importance of the last layers. Further analysis of this phenomenon is left for future work.

## C Generated Examples

We present generated examples from our pruned model using NutePrune at 20% sparsity. We provide examples of three types: without tuning (w/o tune), with tuning but without post-finetuning (w/ tune), and with tuning and post fine-tuning (w/ tune + post FT). The results are displayed in Table 13.

<sup>1</sup><https://github.com/EleutherAI/lm-evaluation-harness/pull/531>

Ratio	Tune	Method	BoolQ	PIQA	HellaSwag	WinoGrande	ARC-e	ARC-c	OBQA	Avg.
0%		LLaMA-7B	75.11	79.16	76.21	69.85	75.29	44.71	44.40	66.39
20%		LLM-Pruner	57.49	76.06	69.53	63.93	67.17	38.05	40.80	59.01
		‡NutePrune	70.21	76.93	71.66	68.27	71.09	40.44	42.60	<b>63.03</b>
20%	✓	LLM-Pruner	67.37	77.86	71.47	65.90	69.57	39.59	41.80	61.94
		Compresso	73.21	75.90	66.90	68.90	69.99	41.47	41.80	62.60
		‡NutePrune	73.79	77.37	72.27	67.48	72.77	38.91	42.40	63.57
		NutePrune	75.38	78.02	72.97	67.40	73.82	40.36	42.80	<b>64.39</b>
25%		‡NutePrune	71.53	76.50	70.60	65.98	69.11	39.93	38.40	61.72
25%	✓	‡NutePrune	72.91	77.42	70.34	68.11	70.92	41.55	40.60	63.12
		NutePrune	74.95	77.75	71.27	67.40	71.25	41.81	42.00	63.78
50%		LLM-Pruner	46.48	61.10	36.87	51.78	35.10	27.65	27.60	40.94
		‡NutePrune	65.38	69.04	55.08	61.33	55.72	30.80	34.60	<b>53.14</b>
50%	✓	LLM-Pruner	57.89	69.97	50.06	52.64	49.66	28.58	34.00	48.97
		Compresso	61.31	66.32	40.73	52.41	51.18	27.65	32.40	47.43
		‡NutePrune	67.25	70.67	56.64	59.83	57.07	31.74	34.40	53.94
		NutePrune	67.52	71.60	58.64	60.14	59.72	32.94	33.80	<b>54.91</b>

Table 10: Zero-shot performance of the compressed LLaMA models in the new version of lm-evaluation-harness. **Bold** denotes the best average performance at the same setting.

# Hidden Dim	4070							
Layer	1	2	3	4	5	6	7	8
# Head	23	22	30	22	29	27	30	28
# Intermediate Dim	5832	7820	9169	9187	8967	9163	9186	9112
Layer	9	10	11	12	13	14	15	16
# Head	30	30	31	30	32	27	30	30
# Intermediate Dim	9261	9165	9303	9695	10005	10258	10417	10564
Layer	17	18	19	20	21	22	23	24
# Head	30	29	26	25	23	21	16	21
# Intermediate Dim	10715	10759	10785	10790	10808	10778	10729	10707
Layer	25	26	27	28	29	30	31	32
# Head	14	15	6	7	11	8	7	9
# Intermediate Dim	10568	10366	10054	9403	8519	7297	6588	5164

Table 11: Detailed structure of compressed 20% LLaMA.

# Hidden Dim	4021							
Layer	1	2	3	4	5	6	7	8
# Head	24	21	27	18	26	21	20	20
# Intermediate Dim	3980	6216	7120	6590	5889	5731	5283	4944
Layer	9	10	11	12	13	14	15	16
# Head	25	21	27	26	23	21	25	21
# Intermediate Dim	4879	4563	4590	5040	5469	5832	6215	778
Layer	17	18	19	20	21	22	23	24
# Head	23	23	21	19	16	15	6	13
# Intermediate Dim	7706	8121	8005	8091	8232	7705	6879	6134
Layer	25	26	27	28	29	30	31	32
# Head	7	8	2	5	7	5	3	8
# Intermediate Dim	4912	3762	3107	2527	2357	2314	2931	2622

Table 12: Detailed structure of compressed 50% LLaMA.



Model	Generated Sentences
	Neural Sequential Model, especially transformers, are used to model the human brain. The model is based on the assumption that the brain is a system of interconnected modules, each of which is responsible for a specific function. The modules are connected by synapses, which are the synaptic connections between neurons. The synapsis are responsible for the transfer of information between the modules. The information is transferred in the form of a signal, which is a combination of the input signal and the output signal. The input signal is the signal that is received by the module, while the output is the output of the module. The output signal is then passed on to the next ...
w/o Tune	<p>AI can create a logo in seconds . It can also create a logo that is more creative than a human designer. AI can also create logos that are more creative and unique than a human can. AI is also able to create logos in a short time. AI has the ability to create logotypes. AI will be able to create a logo for a company in a short period of time. A logo is a symbol that is used to identify a company. A logo can be used to identify the company. AI also has the ability of creating logos that can be used in different ways. AI creates logos that have ...</p> <p>What's great about the holiday season ? The holiday spirit, of course! But it's also a time for giving. And giving is what we do best at the Chamber. We give you the opportunity to give back to your community. We've got a great list of charitable organizations that are looking for your help. And we've made it easy for you to give. Just click on the link below and you'll be taken to the website of the charity of your choice. You can make a donation online, or you can mail a check. Either way, you're giving back to the community. And that's what the ...</p>
	Neural Sequential Model, especially transformers, are used to model the human brain. The model is a neural network that is trained to recognize patterns in data. The network is trained by feeding it data and then testing it on new data. This is a form of machine learning. The neural network is a type of artificial neural network. The term neural network refers to a group of neurons that are connected together. The neurons are connected by synapses. The synapsis are the connections between the neurons. The connections are called weights. The weights are used by the network to determine how much influence each neuron has on the other neurons in the network ...
w/ Tune	<p>AI can create a logo in seconds , but it can't design a logo that will work for a company. AI is great for creating a logo, but not for designing a logo. A logo is a visual representation of a company's brand. It's a visual symbol that represents the company' logo. It is a symbol that is used to identify a company and its products. A good logo is one that is memorable and easy to recognize. A great logo is memorably and easy recognizable. A bad logo is not memorable or easy to recognize. A poorly designed logo is also not memorably or easy recognized. A ...</p> <p>What's great about the holiday season is that it's a time for giving. And what better way to give than to give the gift of a new home? Whether you're buying a home for the first time or you've been in the market for a while, there are a few things you should know before you start looking. 1. Know your budget. This is the most important thing to consider when you'll be looking for a home. You'll want to know how much you can afford to spend on a home, and you'd want to make sure you don't go over your budget when you start searching for homes ...</p>
	Neural Sequential Model, especially transformers, is a popular model for deep learning. It is a model that can be trained on large datasets. It is used for tasks such as image classification, natural language processing, and other tasks. It is an effective model for tasks such as image classification and natural language processing. It is also used for tasks like natural language understanding and natural film fiction. It is able to learn from large datasets and can be trained for various tasks. It can be used for tasks that require a large amount of data. It is capable of learning from large datasets. It is used to train neural networks for tasks such image classification ...
w/ Tune + Post FT	<p>AI can create a logo in seconds. AI Logo Designer. A Logo AI Can Create. A logo is a key element of any business. It helps to create a unique identity and helps to stand out from the crowd. A good logo can help to build trust and credibility. It can also help to increase the sales. Alogoai.com is a AI logo design service. We can design logos for any business, with the latest technology. We provide high quality logos at affordable prices. Logo design service provided by Alogo AI. Logos designed by AI are more accurate and effective. Logs designed by ...</p> <p>What's great about the holiday season is the opportunity to give back to the community. Whether it's donating to a local charity or giving to a needy family, it'll make you feel good to know that you're making a difference. Here are some ideas to get started. Volunteer at a local organization. If you have some free time, consider volunteering at a nonprofit organization. You can help with the holidays by helping with the decorations, helping with food preparation, or even helping with some of the administrative tasks. You'll feel great knowing that you are making a contribution to the organization. Donate ...</p>

Table 13: Generated Examples from the Compressed LLaMA-7B at 20% sparsity