

SHORTENED LLAMA: A SIMPLE DEPTH PRUNING FOR LARGE LANGUAGE MODELS

Bo-Kyeong Kim¹ Geonmin Kim¹ Tae-Ho Kim¹
 Thibault Castells¹ Shinkook Choi¹ Junho Shin¹ Hyoung-Kyu Song¹

¹Nota Inc.

{bokyeong.kim, geonmin.kim, thkim}@nota.ai

ABSTRACT

Structured pruning of modern large language models (LLMs) has emerged as a way of decreasing their high computational needs. Width pruning reduces the size of projection weight matrices (e.g., by removing attention heads) while maintaining the number of layers. Depth pruning, in contrast, removes entire layers or blocks, while keeping the size of the remaining weights unchanged. Most current research focuses on either width-only or a blend of width and depth pruning, with little comparative analysis between the two units (width vs. depth) concerning their impact on LLM inference efficiency. In this work, we show that a simple depth pruning approach can compete with recent width pruning methods in terms of zero-shot task performance. Our pruning method boosts inference speeds, especially under memory-constrained conditions that require limited batch sizes for running LLMs, where width pruning is ineffective. We hope this work can help deploy LLMs on local and edge devices. Code and model can be found at: <https://github.com/Nota-NetsPresso/shortened-llm>.

1 INTRODUCTION

The advancement of large language models (LLMs) (Touvron et al., 2023; OpenAI, 2023; Chowdhery et al., 2022; Zhang et al., 2022; Scao et al., 2022) has brought significant improvements in language-based tasks, enabling versatile applications such as powerful chatbots (Google, 2023; OpenAI, 2022). However, the deployment of LLMs is constrained by their intensive computational demands. To make LLMs more accessible and efficient for practical use, various optimization strategies have been actively studied over recent years (see Section 5 for details). This work focuses on *structured* pruning (Fang et al., 2023a; Li et al., 2017b), which removes groups of unnecessary weights and can facilitate hardware-agnostic acceleration.

In the context of compressing billion-parameter LLMs, LLM-Pruner (Ma et al., 2023) and FLAP (An et al., 2024) narrow the network width by pruning coupled structures (e.g., attention heads and their associated weight connections) while maintaining the number of layers. Sheared-LLaMA (Xia et al., 2024) reduces not only the network width but also its depth by entirely removing some layers. Despite the existence of pruning methods (Xia et al., 2022; Kurtic et al., 2023; Xia et al., 2024) that incorporate both width and depth aspects, there remains a gap in detailed analysis comparing these two factors (see Figure 1), specifically in relation to their impact on LLM inference efficiency.

Problem: Small-batch LLM Inference. In addition to substantial model sizes, LLM inference is distinguished by an autoregressive decoding mechanism, which predicts tokens one by one based on the input and the previously generated tokens. This token-by-token generation process often involves multiplying large matrices

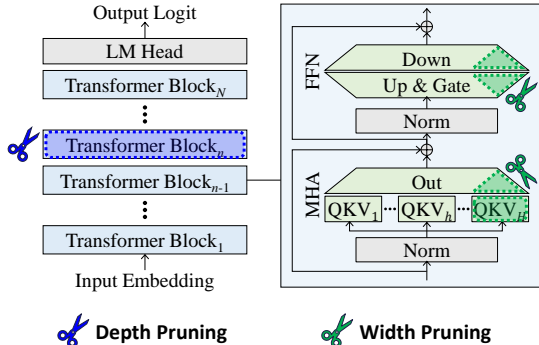


Figure 1: Comparison of pruning units. Width pruning reduces the size of projection weight matrices. Depth pruning removes Transformer blocks, or individual MHA and FFN modules.

(weights) with smaller matrices or vectors (activations). The primary bottleneck for inference efficiency is memory access operations rather than the speed of mathematical computations (referred to as ‘memory-bound’), leading to suboptimal use of GPU computing power (Kwon et al., 2023).

Though increasing batch sizes is a standard way to enhance GPU computation and throughput, it poses a risk of out-of-memory (OOM) errors, as depicted in Figure 2,¹ unless advanced system-level optimizations (Kwon et al., 2023; Jin et al., 2023) are applied. In this study, our focus is on accelerating the inference of LLMs under small-batch conditions caused by hardware restrictions. Such situations are relevant for deploying LLMs on memory-constrained local devices, which can enhance user experience and data privacy protection.

Approach Overview. Depth pruning is often regarded as being less effective in generation performance compared to width pruning, due to the elimination of bigger and coarse units. Contrary to this prevailing view, we show that a simple depth pruning coupled with a LoRA retraining (Hu et al., 2022) can rival recent width pruning studies for LLMs in zero-shot capabilities. Our contributions are summarized as follows:

- For small-batch scenarios, reducing weight shapes through width pruning does not improve LLM inference efficiency and can even degrade it when the resulting weight dimensions are unsuitable for GPU capabilities (see Figure 2). This aspect has been underexplored in prior works.
- We introduce a simple yet effective method for depth pruning of LLMs by exploring various design factors. Our compact LLMs, obtained by excluding several Transformer blocks, achieve actual speedups. They perform comparably to finely width-pruned models in zero-shot tasks.

2 METHOD: BLOCK PRUNING

An LLM is a stack of multiple Transformer blocks (Vaswani et al., 2017), each of which contains a pair of multi-head attention (MHA) and feed-forward network (FFN) modules (see Figure 1). We consider the Transformer block as the pruning unit to prioritize faster inference. Our approach is simple: after identifying unimportant blocks, we perform one-shot pruning and light retraining.

2.1 EVALUATION OF BLOCK-LEVEL IMPORTANCE

The linear weight matrix is denoted as $\mathbf{W}^{k,n} = [W_{i,j}^{k,n}]$ with a size of (d_{out}, d_{in}) , where k represents the operation type (e.g., a query projection in MHA or an up projection in FFN) in the n -th Transformer block. The weight scores are initially calculated at the output neuron level (Sun et al., 2024) and then summed² to assess the block-level importance. See Section B for comparative methods.

Taylor+. Assessing the error caused by the removal of a weight parameter helps in identifying its significance. For a given calibration dataset D , this can be expressed as the alteration in the training loss \mathcal{L} (LeCun et al., 1989; Molchanov et al., 2019): $|\mathcal{L}(W_{i,j}^{k,n}; D) - \mathcal{L}(W_{i,j}^{k,n} = 0; D)| \approx \left| \frac{\partial \mathcal{L}(D)}{\partial W_{i,j}^{k,n}} W_{i,j}^{k,n} \right|$, where we omit the second-order derivatives (Ma et al., 2023; Fang et al., 2023b). We

¹Using the Hugging Face’s Transformers library (Wolf et al., 2020), we ran the LLMs with 12 input tokens and varying output lengths for 20 batched runs after 10 warm-ups. Top: Peak GPU compute utilization (NVIDIA, 2018). Bottom: Mean latency over 20 runs, evaluated with 27% and 29% pruning for the 7B and 13B models, respectively. See Section A for additional results.

²In our exploration of various aggregation strategies (i.e., sum, mean, product, and max operations across module and block levels), summing the scores was effective at different pruning ratios.

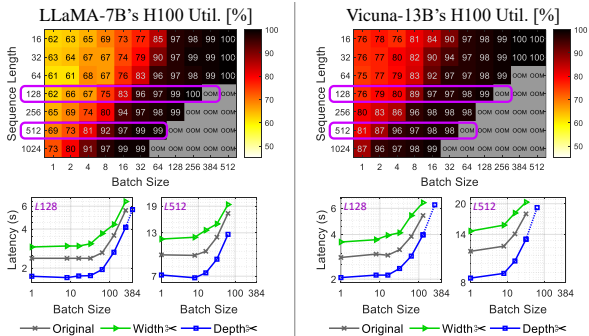
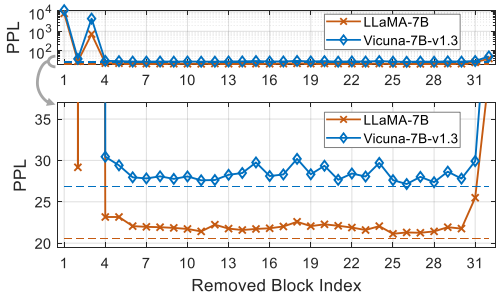


Figure 2: Top: GPU compute utilization of the original LLMs. Bottom: Latency results (L : target output length). Under similar parameters, our depth pruning yields higher speeds than width pruning (Ma et al., 2023).

define the block score as $I_{\text{Taylor}}^n = \sum_k \sum_i \sum_j \left| \frac{\partial \mathcal{L}(D)}{\partial W_{i,j}^{k,n}} W_{i,j}^{k,n} \right|$. The symbol ‘+’ denotes a heuristic that keeps the initial and final few blocks unpruned (Lee et al., 2021; Ma et al., 2023).

Perplexity (PPL). We physically remove each Transformer block and monitor its influence on PPL using the calibration set D : $I_{\text{PPL}}^n = \exp \left\{ -\frac{1}{SL} \sum_s \sum_l \log p_{\theta^n}(x_l^{(s)} | x_{<l}^{(s)}) \right\}$, where θ^n denotes the model without its n -th block, and $s = 1, \dots, S$ and $l = 1, \dots, L$ are the indices for sequences and tokens in D . The use of PPL can reflect the model’s behavior by being derived from the next-token prediction loss; it requires only the forward pass, avoiding the need to compute back-propagation gradients (Ma et al., 2023) and Hessian inverses (Kurtic et al., 2023), or to involve a mask learning stage (Xia et al., 2024). Figure 3 shows that several blocks are removable, showing only a slight effect on the PPL metric.



2.2 ONE-SHOT PRUNING AND RETRAINING

After sorting the block-level importance scores, we prune the less crucial blocks in a single step. Since every block has an identical configuration and it is easy to calculate the number of parameters for one block, we readily decide how many blocks should be removed to meet the target model size.

Figure 3: Estimated importance of each Transformer block on the calibration set. Blocks with lower PPL scores are pruned.

We efficiently retrain the pruned models with the low-rank adaptation (LoRA) method (Hu et al., 2022; Ma et al., 2023). The weight matrix of the adapted network is expressed as $W_0 + \Delta W = W_0 + BA$, where W_0 denotes the initial weight with a shape of $(d_{\text{out}}, d_{\text{in}})$. The update matrix ΔW is decomposed into two trainable parts, B and A with dimensions (d_{out}, r) and (r, d_{in}) , where r denotes a low rank. We show that LoRA has the potential to restore the performance of depth-pruned models. LoRA-based retraining can be efficiently completed on a single GPU in just a few hours. For example, retraining a model pruned by 20% from 7B parameters takes about 2 hours and 22GB VRAM, while a model reduced by 21% from 13B demands around 3 hours and 35GB VRAM.

3 EXPERIMENTAL SETUP

In this section, we outline the primary setup, while the details can be found in Section C.

Baseline and Data. LLM-Pruner (Ma et al., 2023), FLAP (An et al., 2024), and Wanda-sp (i.e., a structured variant (An et al., 2024) of Wanda (Sun et al., 2024)) serve as the baselines for width pruning. Following Ma et al. (2023), we randomly select 10 samples from BookCorpus (Zhu et al., 2015) to compute block-level significance during the pruning stage. We also use this calibration set for all the baseline methods to ensure a fair comparison. At the LoRA retraining stage, 50K samples of the refined Alpaca (Taori et al., 2023) are used.

Latency and Throughput. We follow Sheng et al. (2023) to measure the metrics. Given a batch size M and an output sequence length L (excluding the input length), the latency T represents the time required to handle the given prompts and produce ML output tokens. The throughput is computed as ML/T . We report the average results from 20 runs after the initial 10 warm-up batches.

4 RESULTS

This section presents key results. Additional results and analyses can be found in Appendix.

Main Comparison (Tables 1 and 6). The width pruning methods do not improve LLM inference efficiency. Under limited input (batch) scales, the processing speed largely hinges on the frequency of memory access operations. Addressing this issue by merely reducing matrix sizes is challenging, unless they are completely removed. The speed even worsens compared to the original model due to GPU-unfriendly operation dimensions (e.g., the hidden sizes of FFN are often not divisible by 8 (Table 5), which hinders the effective utilization of GPU Tensor Cores (Andersch et al., 2019)). In contrast, our depth pruning achieves speed gains by completely removing several Transformer blocks, resulting in fewer memory accesses and matrix-level operations between activations and

Table 1: Zero-shot results of the compressed (top) LLaMA-7B and (bottom) Vicuna-13B-v1.3. The results of Wanda-sp (Sun et al., 2024; An et al., 2024), FLAP (An et al., 2024), and LLM-Pruner (Ma et al., 2023) were obtained with their official codes. See Table 6 for detailed results.

Model		Zero-shot Performance			H100 80GB [‡]		RTX3090 24GB [‡]	
		PPL↓ WikiText2	PTB	Ave Acc↑ (%) [†]	Latency↓ (s)	Throughput↑ (tokens/s)	Latency↓ (s)	Throughput↑ (tokens/s)
LLaMA-7B (6.7B)		12.6	22.1	66.3	2.4	53.7	5.1	25.0
20% Pruned (5.5B)	Wanda-sp	21.4	47.2	51.8	3.1	41.7	7.6	16.7
	FLAP	17.0	30.1	59.5	3.2	40.5	7.7	16.5
	LLM-Pruner	17.6	30.4	61.8	3.0	43.2	6.0	21.4
	Ours: Taylor+	20.2	32.3	63.5	1.9	66.0	4.5	28.4
	Ours: PPL	17.7	30.7	61.9	1.9	66.0	4.5	28.4
35% Pruned (4.5B)	Wanda-sp	133.6	210.1	36.9	3.1	41.6	8.0	16.1
	FLAP	25.6	44.4	52.7	3.2	40.5	8.1	15.8
	LLM-Pruner	24.2	40.7	55.5	2.9	44.4	6.1	21.1
	Ours: Taylor+	33.2	58.5	55.4	1.6	80.1	3.4	37.8
	Ours: PPL	23.1	38.8	55.2	1.6	80.1	3.4	37.8

Model		Zero-shot Performance			H100 80GB [‡]		RTX3090 24GB [‡]	
		PPL↓ WikiText2	PTB	Ave Acc↑ (%) [†]	Latency↓ (s)	Throughput↑ (tokens/s)	Latency↓ (s)	Throughput↑ (tokens/s)
Vicuna-13B-v1.3 (13.0B)		14.7	51.6	68.3	2.8	45.5	OOM	OOM
21% Pruned (10.5B)	Wanda-sp	19.0	71.8	63.6	3.8	34.1	9.8	12.9
	FLAP	18.8	65.3	63.3	3.9	32.6	10.2	12.6
	LLM-Pruner	16.0	57.0	65.3	3.8	34.0	7.5	17.3
	Ours: Taylor+	18.1	61.6	66.7	2.3	55.7	5.4	23.9
	Ours: PPL	16.1	56.5	64.9	2.3	55.7	5.4	23.9
37% Pruned (8.3B)	Wanda-sp	36.6	123.5	52.7	3.8	33.8	10.5	12.6
	FLAP	28.7	96.2	58.3	3.9	32.9	9.7	13.2
	LLM-Pruner	22.2	74.0	59.7	3.6	35.6	7.1	18.0
	Ours: Taylor+	34.2	90.4	61.4	1.8	69.7	4.0	31.7
	Ours: PPL	22.1	73.6	59.1	1.8	69.7	4.0	31.7

[†]Average accuracy on seven commonsense reasoning tasks.

[‡]Measured with 12 input tokens, 128 output tokens, and a batch size of 1 on a single GPU.

Table 2: Impact of block pruning criteria. Results of 20%-pruned LLaMA-7B.

Pruning Criterion	PPL↓ on Wiki2	Ave Acc↑ (%)
Magnitude	7720.7	34.4
Magnitude+	19.4	56.1
Taylor	3631.7	35.5
Taylor+	20.2	63.5
PPL	17.7	61.9

Table 3: Impact of pruning units (individual MHA and FFN modules vs. Transformer blocks) on LLaMA-7B.

Pruning Unit	#Param	PPL↓ on Wiki2	Ave Acc↑ (%)
Module	5.3B	25.2	61.1
Block	5.3B	18.6	60.6
Module	4.6B	38.9	52.5
Block	4.5B	23.1	55.2

Table 4: Impact of calibration data volume. Results of 20%-pruned LLaMA-7B.

Metric	PPL↓ on Wiki2		Ave Acc↑ (%)	
	# Cal Samples	10	100	10
Wanda-sp	21.4	21.7	51.8	52.0
FLAP	17.0	17.5	59.5	59.9
LLM-Pruner	17.6	17.0	61.8	61.7
Ours: Taylor+	20.2	19.0	63.5	63.9
Ours: PPL	17.7	17.4	61.9	61.7

weights. Moreover, under the same retraining setup as Ma et al. (2023), our models achieve zero-shot scores on par with finely width-pruned models. See Section F for generation examples.

Importance Criteria for Block Pruning (Tables 2 and 7). The basic methods without the ‘+’ label remove the critical initial blocks, resulting in inferior performance. In addition, using weight magnitude alone for pruning choices could be unsuitable. The Taylor+ criterion enhances accuracy in commonsense reasoning tasks, while the PPL method leads to better generation quality without relying on heuristic selection of pruning candidates.

Structural Unit for Depth Pruning (Tables 3 and 8). While individual MHA and FFN modules, which are finer than Transformer blocks, can be pruned, our results show that entire block removal generally outperforms module-level pruning. This contradicts the common view that more granular pruning units are better. It may be suboptimal to treat the modules in isolation, considering their collaborative roles (i.e., MHA captures dependency relations (Vaswani et al., 2017), while skip connections and FFN prevent rank collapse in purely attention-driven networks (Dong et al., 2021)).

Calibration Data Volume (Tables 4 and 9). The calibration set is employed to assess the weight significance of width pruning baselines and the block-level importance of our method during the pruning phase. Upon varying the number of calibration samples in the BookCorpus dataset, the results for the examined methods remain fairly unchanged, suggesting the adequacy of 10 samples.

5 RELATED WORK

Numerous techniques have been developed towards efficient LLMs, including knowledge distillation (Fu et al., 2023; Hsieh et al., 2023), quantization (Frantar et al., 2023; Dettmers et al., 2022), and system-level inference acceleration (Dao, 2023; Kwon et al., 2023). See Zhu et al. (2023); Wan et al. (2023) for survey.

In this study, we focus on network pruning (LeCun et al., 1989), which has a long-standing reputation in the model compression field. Beyond its use in relatively small-scale convolutional networks (Li et al., 2017a; He et al., 2019) and Transformer models (Yu et al., 2022; Xia et al., 2022; Kurtic et al., 2023), pruning has recently begun to be applied to contemporary LLMs. Several studies (Frantar & Alistarh, 2023; Sun et al., 2024) employ unstructured and semi-structured (Aojun Zhou, 2021) pruning by zeroing individual neurons. SparseGPT (Frantar & Alistarh, 2023) addresses the layer-wise reconstruction problem for pruning by computing Hessian inverses. Wanda (Sun et al., 2024) introduces a pruning criterion that involves multiplying weight magnitudes by input feature norms. Despite the plausible performance of pruned models using zero masks, they necessitate specialized support for sparse matrix operations to ensure actual speedups.

In contrast, structured pruning removes organized patterns, such as layers (Fan et al., 2020; Jha et al., 2023), MHA’s attention heads (Voita et al., 2019; Michel et al., 2022), FFN’s hidden sizes (Nova et al., 2023; Santacroce et al., 2023), and some hybrid forms (Lagunas et al., 2021; Liu et al., 2021; Xia et al., 2022; Kwon et al., 2022; Kurtic et al., 2023), thereby improving inference efficiency in a hardware-agnostic way. To compress LLMs, FLAP (An et al., 2024) and LLM-Pruner (Ma et al., 2023) eliminate coupled structures in the aspect of network width while retaining the number of layers. Sheared-LLaMA (Xia et al., 2024) introduces a mask learning phase aimed at identifying prunable components in both the network’s width and depth. Our work explores the relatively untapped area of depth-only pruning for multi-billion parameter LLMs, which can markedly accelerate latency while attaining competitive results.

Strategies for skipping layers (Schuster et al., 2022; Corro et al., 2023; Raposo et al., 2024) effectively serve to decrease computational burdens. Moreover, depth pruning approaches (Song et al., 2024; Men et al., 2024; Tang et al., 2024) for LLMs have been proposed concurrently with our work, based on the detected redundancy in Transformer blocks.

6 CONCLUSION

By introducing a block pruning method, we conduct an in-depth comparative analysis on the impact of network width and depth on LLM compression. Our work involves the one-shot removal of Transformer blocks, determined by evaluating various design choices. Despite its simplicity, our method matches the zero-shot performance of recent width pruning works. Moreover, it offers significant inference speedups in resource-constrained scenarios that require running LLMs with limited batch sizes, where width pruning falls short. Future research will explore more powerful retraining techniques, including full parameter updates and knowledge distillation.

ACKNOWLEDGMENTS

We thank the Microsoft Startups Founders Hub program and the AI Industrial Convergence Cluster Development project funded by the Ministry of Science and ICT (MSIT, Korea) and Gwangju Metropolitan City for their generous support of GPU resources.

REFERENCES

- Yongqi An, Xu Zhao, Tao Yu, Ming Tang, and Jinqiao Wang. Fluctuation-based adaptive structured pruning for large language models. In *AAAI*, 2024. 1, 3, 4, 5, 10, 11, 12, 17
- Michael Andersch, Valerie Sarge, and Paulius Micikevicius. Tensor core dl performance guide. In *NVIDIA GTC*, 2019. 3, 12, 17
- Junnan Zhu Jianbo Liu Zhijie Zhang Kun Yuan Wenxiu Sun Hongsheng Li Aojun Zhou, Yukun Ma. Learning n:m fine-grained structured sparse neural networks from scratch. In *ICLR*, 2021. 5
- Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. Piqa: Reasoning about physical commonsense in natural language. In *AAAI*, 2020. 10

- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, et al. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, 2023. URL <https://lmsys.org/blog/2023-03-30-vicuna/>. 10
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022. 1
- Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. BoolQ: Exploring the surprising difficulty of natural yes/no questions. In *NAACL*, 2019. 10
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*, 2018. 10
- Luciano Del Corro, Allie Del Giorno, Sahaj Agarwal, Bin Yu, Ahmed Awadallah, and Subhabrata Mukherjee. Skipdecode: Autoregressive skip decoding with batching and caching for efficient llm inference. *arXiv preprint arXiv:2307.02628*, 2023. 5
- Tri Dao. Flashattention-2: Faster attention with better parallelism and work partitioning. *arXiv preprint arXiv:2307.08691*, 2023. 5
- Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. Llm.int8(): 8-bit matrix multiplication for transformers at scale. In *NeurIPS*, 2022. 5
- Yihe Dong, Jean-Baptiste Cordonnier, and Andreas Loukas. Attention is not all you need: Pure attention loses rank doubly exponentially with depth. In *ICML*, 2021. 4, 13
- EleutherAI. Language model evaluation harness (package version 3326c54). <https://github.com/EleutherAI/lm-evaluation-harness>, 2023. 10
- Angela Fan, Edouard Grave, and Armand Joulin. Reducing transformer depth on demand with structured dropout. In *ICLR*, 2020. 5
- Gongfan Fang, Xinyin Ma, Mingli Song, Michael Bi Mi, and Xinchao Wang. Depgraph: Towards any structural pruning. In *CVPR*, 2023a. 1
- Gongfan Fang, Xinyin Ma, and Xinchao Wang. Structural pruning for diffusion models. In *NeurIPS*, 2023b. 2
- Elias Frantar and Dan Alistarh. SparseGPT: Massive language models can be accurately pruned in one-shot. In *ICML*, 2023. 5
- Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. OPTQ: Accurate quantization for generative pre-trained transformers. In *ICLR*, 2023. 5
- Yao Fu, Hao Peng, Litu Ou, Ashish Sabharwal, and Tushar Khot. Specializing smaller language models towards multi-step reasoning. In *ICML*, 2023. 5
- Trevor Gale, Erich Elsen, and Sara Hooker. The state of sparsity in deep neural networks. In *ICML Workshop*, 2019. 9
- Google. An important next step on our ai journey. <https://blog.google/technology/ai/bard-google-ai-search-updates/>, 2023. 1
- Yang He, Ping Liu, Ziwei Wang, Zhilan Hu, and Yi Yang. Filter pruning via geometric median for deep convolutional neural networks acceleration. In *CVPR*, 2019. 5
- Cheng-Yu Hsieh, Chun-Liang Li, Chih-Kuan Yeh, et al. Distilling step-by-step! outperforming larger language models with less training data and smaller model sizes. In *Findings of ACL*, 2023. 5
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. In *ICLR*, 2022. 2, 3, 10
- Ananya Harsh Jha, Tom Sherborne, Evan Pete Walsh, Dirk Groeneveld, Emma Strubell, and Iz Beltagy. How to train your (compressed) large language model. *arXiv preprint arXiv:2305.14864*, 2023. 5
- Yunho Jin, Chun-Feng Wu, David Brooks, and Gu-Yeon Wei. S³: Increasing gpu utilization during generative inference for higher throughput. In *NeurIPS*, 2023. 2

- Eldar Kurtic, Elias Frantar, and Dan Alistarh. Ziplm: Inference-aware structured pruning of language models. In *NeurIPS*, 2023. 1, 3, 5, 10
- Woosuk Kwon, Sehoon Kim, Michael W. Mahoney, Joseph Hassoun, Kurt Keutzer, and Amir Gholami. A fast post-training pruning framework for transformers. In *NeurIPS*, 2022. 5
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *SOSP*, 2023. 2, 5
- François Lagunas, Ella Charlaix, Victor Sanh, and Alexander M. Rush. Block pruning for faster transformers. In *EMNLP*, 2021. 5
- Yann LeCun, John Denker, and Sara Solla. Optimal brain damage. In *NeurIPS*, 1989. 2, 5, 9
- Jaeho Lee, Sejun Park, Sangwoo Mo, Sungsoo Ahn, and Jinwoo Shin. Layer-adaptive sparsity for the magnitude-based pruning. In *ICLR*, 2021. 3, 9
- Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient convnets. In *ICLR*, 2017a. 5, 9
- Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient convnets. In *ICLR*, 2017b. 1
- Zejian Liu, Fanrong Li, Gang Li, and Jian Cheng. Ebert: Efficient bert inference with dynamic structured pruning. In *Findings of ACL*, 2021. 5
- Xinyin Ma, Gongfan Fang, and Xinchao Wang. Llm-pruner: On the structural pruning of large language models. In *NeurIPS*, 2023. 1, 2, 3, 4, 5, 9, 10, 11, 12, 15, 17
- Mitchell P. Marcus, Beatrice Santorini, and Mary Ann Marcinkiewicz. Building a large annotated corpus of English: The Penn Treebank. *Computational Linguistics*, 19(2):313–330, 1993. 10
- Xin Men, Mingyu Xu, Qingyu Zhang, Bingning Wang, Hongyu Lin, Yaojie Lu, Xianpei Han, and Weipeng Chen. Shortgpt: Layers in large language models are more redundant than you expect. *arXiv preprint arXiv:2403.03853*, 2024. 5
- Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models. In *ICLR*, 2017. 10
- Paul Michel, Omer Levy, and Graham Neubig. Are sixteen heads really better than one? In *NeurIPS*, 2022. 5
- Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct electricity? a new dataset for open book question answering. In *EMNLP*, 2018. 10
- Pavlo Molchanov, Arun Mallya, Stephen Tyree, Iuri Frosio, and Jan Kautz. Importance estimation for neural network pruning. In *CVPR*, 2019. 2, 9
- Azade Nova, Hanjun Dai, and Dale Schuurmans. Gradient-free structured pruning with unlabeled data. In *ICML*, 2023. 5
- NVIDIA. Useful nvidia-smi queries. <https://enterprise-support.nvidia.com/s/article/Useful-nvidia-smi-Queries-2>, 2018. 2, 9
- OpenAI. Introducing chatgpt. <https://openai.com/blog/chatgpt>, 2022. 1
- OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 1
- David Raposo, Sam Ritter, Blake Richards, Timothy Lillicrap, Peter Conway Humphreys, and Adam Santoro. Mixture-of-depths: Dynamically allocating compute in transformer-based language models. *arXiv preprint arXiv:2404.02258*, 2024. 5
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale. *arXiv preprint arXiv:1907.10641*, 2019. 10
- Michael Santacrose, Zixin Wen, Yelong Shen, and Yuanzhi Li. What matters in the structured pruning of generative language models? *arXiv preprint arXiv:2302.03773*, 2023. 5
- Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, et al. Bloom: A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*, 2022. 1

- Tal Schuster, Adam Fisch, Jai Gupta, Mostafa Dehghani, Dara Bahri, Vinh Tran, Yi Tay, and Donald Metzler. Confident adaptive language modeling. In *NeurIPS*, 2022. 5
- Ying Sheng, Lianmin Zheng, Binhang Yuan, Zhuohan Li, Max Ryabinin, et al. Flexgen: High-throughput generative inference of large language models with a single gpu. In *ICML*, 2023. 3
- Jiwon Song, Kyungseok Oh, Taesu Kim, Hyungjun Kim, Yulhwa Kim, and Jae-Joon Kim. Sleb: Streamlining llms through redundancy verification and elimination of transformer blocks. *arXiv preprint arXiv:2402.09025*, 2024. 5
- Mingjie Sun, Zhuang Liu, Anna Bair, and J. Zico Kolter. A simple and effective pruning approach for large language models. In *ICLR*, 2024. 2, 3, 4, 5, 9, 10, 11, 12
- Yehui Tang, Fangcheng Liu, Yunsheng Ni, Yuchuan Tian, Zheyuan Bai, Yi-Qi Hu, Sichao Liu, Shangling Jui, Kai Han, and Yunhe Wang. Rethinking optimization and architecture for tiny language models. *arXiv preprint arXiv:2402.02791*, 2024. 5
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, et al. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca, 2023. 3, 10
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023. 1, 10
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017. 2, 4, 13
- Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov. Analyzing multi-head self-attention: Specialized heads do the heavy lifting, the rest can be pruned. In *ACL*, 2019. 5
- Zhongwei Wan, Xin Wang, Che Liu, Samiul Alam, Yu Zheng, Zhongnan Qu, Shen Yan, Yi Zhu, Quanlu Zhang, Mosharaf Chowdhury, and Mi Zhang. Efficient large language models: A survey. *arXiv preprint arXiv:2312.03863*, 2023. 5
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, et al. Transformers: State-of-the-art natural language processing. In *EMNLP: System Demonstrations*, 2020. 2, 9, 10
- Mengzhou Xia, Zexuan Zhong, and Danqi Chen. Structured pruning learns compact and accurate models. In *ACL*, 2022. 1, 5
- Mengzhou Xia, Tianyu Gao, Zhiyuan Zeng, and Danqi Chen. Sheared llama: Accelerating language model pre-training via structured pruning. In *ICLR*, 2024. 1, 3, 5, 10
- Shixing Yu, Tianlong Chen, Jiayi Shen, Huan Yuan, Jianchao Tan, Sen Yang, Ji Liu, and Zhangyang Wang. Unified visual transformer compression. In *ICLR*, 2022. 5
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? In *ACL*, 2019. 10
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, et al. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*, 2022. 1
- Xunyu Zhu, Jian Li, Yong Liu, Can Ma, and Weiping Wang. A survey on model compression for large language models. *arXiv preprint arXiv:2308.07633*, 2023. 5
- Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, et al. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In *ICCV*, 2015. 3, 10

Appendix — Shortened LLaMA: A Simple Depth Pruning for LLMs

A GPU COMPUTE UTILIZATION AND LATENCY RESULTS

Using the Hugging Face’s Transformers library (Wolf et al., 2020) (without xFormers-optimized attention), we ran the LLMs with 12 input tokens and varying output lengths for 20 batched runs after 10 warm-ups. Figure 4 shows the results with pruning ratios of 27% for the 7B-parameter model and 29% for the 13B model. The top part is from peak GPU compute utilization (NVIDIA, 2018), while the bottom part displays the mean latency over 20 runs.

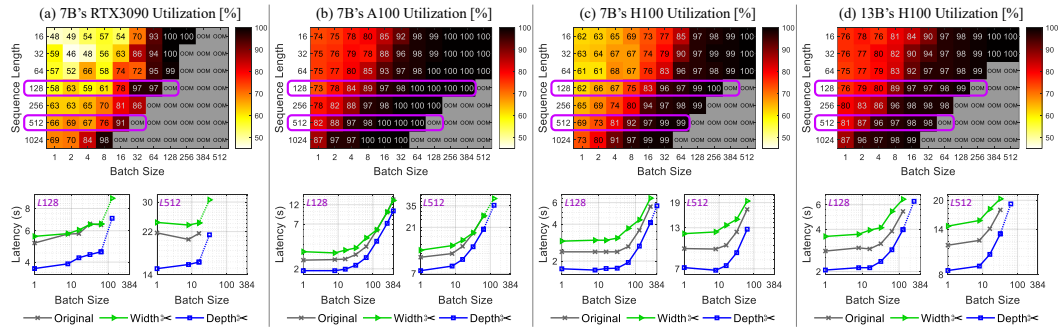


Figure 4: Top: GPU compute utilization of (a)–(c) running LLaMA-7B on different NVIDIA GPUs and that of (d) Vicuna-13B. LLM inference is typically constrained by memory access operations, resulting in lower GPU compute usage. Increasing batch sizes can enhance GPU utilization and throughput, but pushing this too far triggers out-of-memory (OOM) issues. Bottom: Latency results with varying batch sizes and target output lengths (labeled with L). Our depth pruning (blue lines) improves generation speeds over the original models (gray), while width pruning (Ma et al., 2023) is ineffective (green). The dotted lines show that pruned models can operate with larger batch sizes that cause OOM errors for the original model.

B EVALUATION OF BLOCK-LEVEL IMPORTANCE

We consider the following criteria to evaluate the significance of each Transformer block, ultimately selecting the Taylor+ and PPL metrics (see Tables 2 and 7). Specifically, the linear weight matrix is denoted as $\mathbf{W}^{k,n} = [W_{i,j}^{k,n}]$ with a size of (d_{out}, d_{in}) , where k represents the type of operation (e.g., a query projection in MHA or an up projection in FFN) within the n -th Transformer block. The weight importance scores are calculated at the output neuron level (Sun et al., 2024), followed by summing³ these scores to assess the block-level importance.

Magnitude. This metric (Li et al., 2017a) is a fundamental baseline in the pruning literature, assuming that weights with smaller norms are less informative. For the block-level analysis, we compute

$$I_{\text{Magnitude}}^n = \sum_k \sum_i \sum_j |W_{i,j}^{k,n}|.$$

Taylor. Assessing the error caused by the removal of a weight parameter helps in identifying its significance. For a given calibration dataset D , this can be expressed as the alteration in the training loss \mathcal{L} (LeCun et al., 1989; Molchanov et al., 2019): $|\mathcal{L}(W_{i,j}^{k,n}; D) - \mathcal{L}(W_{i,j}^{k,n} = 0; D)| \approx$

$$\left| \frac{\partial \mathcal{L}(D)}{\partial W_{i,j}^{k,n}} W_{i,j}^{k,n} \right|,$$

where we omit the second-order derivatives by following Ma et al. (2023). We define

$$\text{the block score as } I_{\text{Taylor}}^n = \sum_k \sum_i \sum_j \left| \frac{\partial \mathcal{L}(D)}{\partial W_{i,j}^{k,n}} W_{i,j}^{k,n} \right|.$$

Magnitude+ and Taylor+. Upon using the aforementioned metrics, the early blocks are labeled as unimportant, but their removal leads to severe performance drops. Similar to a popular heuristic (Gale et al., 2019; Lee et al., 2021), we preserve the first four and the last two blocks (Ma et al., 2023) by excluding them from the pruning candidates.

³In our exploration of various aggregation strategies (i.e., sum, mean, product, and max operations across module and block levels), summing the scores was effective at different pruning ratios.

Perplexity (PPL). Redundant blocks contribute less to the model’s outputs, and their removal leads to smaller degradation in PPL, a commonly used metric for language modeling tasks. In this context, we physically eliminate each block and monitor its influence on PPL using the calibration set D : $I_{\text{PPL}}^n = \exp \left\{ -\frac{1}{SL} \sum_s \sum_l \log p_{\theta^n}(x_l^{(s)} | x_{<l}^{(s)}) \right\}$, where θ^n denotes the model without its n -th block, and $s = 1, \dots, S$ and $l = 1, \dots, L$ are the indices for sequences and tokens in D . The use of PPL can reflect the model’s behavior by being derived from the next-token prediction loss; it requires only the forward pass, avoiding the need to compute back-propagation gradients (Ma et al., 2023) and Hessian inverses (Kurtic et al., 2023), or to involve a mask learning stage (Xia et al., 2024).

C EXPERIMENTAL SETUP

Model. Our testbed includes LLaMA-7B (Touvron et al., 2023) and Vicuna-{7B, 13B}-v1.3 (Chiang et al., 2023), which are famous open-source LLMs.

Baseline. We compare the two pruning units, network width vs. depth, using the same calibration dataset. The width pruning baseline methods are described below, and we utilize their official code for implementation. Table 5 shows the pruned architectures under similar numbers of parameters.⁴

- LLM-Pruner (Ma et al., 2023) employs a Taylor-based importance metric to remove attention heads from MHA and intermediate neurons from FFN. Local pruning is performed to select removable groups within the same module while maintaining uniform dimensions across the examined blocks. Adhering to their practice, the first and last few blocks remain unpruned. Their pruned models and ours are identically retrained with LoRA.
- FLAP (An et al., 2024) uses a fluctuation-based importance metric to explore the recoverability of feature maps after removing weight columns. Global pruning is applied, leading to different widths over distinct modules (see Table 5 for mean and standard deviation values). Instead of retraining, extra bias terms are added into pruned feature maps for performance restoration.
- Wanda-sp is presented in An et al. (2024) as a variant of Wanda (Sun et al., 2024) adjusted for structured pruning. The original metric was based on the product of weight magnitudes and input activation norms, which can be interpreted as addressing a local reconstruction objective. Wanda-sp extends this in a structured way while using common dimensions among different modules.

Data. Following Ma et al. (2023), we randomly select 10 samples from BookCorpus (Zhu et al., 2015) to compute block-level significance during the pruning stage. We also use this calibration dataset for the baseline methods to ensure a fair comparison. At the LoRA retraining stage, 50K samples of the refined Alpaca (Taori et al., 2023) are used.

Evaluation. Following Touvron et al. (2023), we measure zero-shot accuracy on commonsense reasoning datasets (i.e., BoolQ (Clark et al., 2019), PIQA (Bisk et al., 2020), HellaSwag (Zellers et al., 2019), WinoGrande (Sakaguchi et al., 2019), ARC-easy (Clark et al., 2018), ARC-challenge (Clark et al., 2018), and OpenbookQA (Mihaylov et al., 2018)) using the lm-evaluation-harness package (EleutherAI, 2023). We also report zero-shot PPL on WikiText2 (Merity et al., 2017) and PTB (Marcus et al., 1993).

Implementation Details. We utilize the Transformers library (Wolf et al., 2020) on a single NVIDIA A100 GPU (80GB VRAM). All the experiments involving 7B-parameter models can be conducted on a single NVIDIA RTX 3090 (24GB VRAM).

- At the pruning phase, we assess the significance of Transformer blocks using a small calibration set (containing 10 samples from BookCorpus (Zhu et al., 2015) with a sequence length of 128). For the PPL-based criterion, the calibration samples are fed into networks with a single block removed, and this step is iterated across all the blocks in the target model. For the Taylor+ method, we feed the calibration data into the original network to collect backward-gradient matrices. The pruning is completed efficiently within 1 to 2 hours for the 7B- and 13B-sized models.
- At the retraining stage, we apply a LoRA adapter (Hu et al., 2022) to every projection weight matrix by following Ma et al. (2023). We employ a LoRA rank of 8, a learning rate of 0.0001, and a batch size of 64 over 2 epochs. The retraining costs are notably low, with the entire process being executed on a single GPU. For example, retraining a 20%-pruned model from 7B parame-

⁴We used the parameter numbers from LLM-Pruner’s module-level pruning ratios of (25%, 35%, 45%) as the reference and adjusted the pruning ratios for our method and the other baselines.

Table 5: Pruned architectures on LLaMA-7B and Vicuna-{7B, 13B}-v1.3. While Wanda-sp (Sun et al., 2024; An et al., 2024), FLAP (An et al., 2024), and LLM-Pruner (Ma et al., 2023) reduce the network width, our method reduces the network depth. Using LLM-Pruner’s module-level pruning ratios of (25%, 35%, 45%) as benchmarks, we adjust others for comparable parameter numbers.

	Model	#Param	#Block [‡]	#Head [‡]	FFN-D [‡]
	Original 7B	6.7B	32	32	11008
20% [†]	Wanda-sp	5.5B	32	26	8807
	FLAP	5.4B	32	26.9 \pm 7.5	8577.4 \pm 2078.4
	LLM-Pruner	5.4B	32	24	8256
	Ours	5.5B	26	32	11008
27% [†]	Wanda-sp	4.9B	32	23	7816
	FLAP	4.9B	32	24.6 \pm 8.6	7497.1 \pm 2358.0
	LLM-Pruner	4.9B	32	21	7155
	Ours	4.9B	23	32	11008
35% [†]	Wanda-sp	4.5B	32	21	7156
	FLAP	4.5B	32	23.0 \pm 8.8	6781.1 \pm 2440.6
	LLM-Pruner	4.4B	32	18	6054
	Ours	4.5B	21	32	11008
	Original 13B	13.0B	40	40	13824
21% [†]	Wanda-sp	10.5B	40	32	11060
	FLAP	10.5B	40	33.7 \pm 8.9	10778.7 \pm 2316.0
	LLM-Pruner	10.3B	40	30	10368
	Ours	10.5B	32	40	13824
29% [†]	Wanda-sp	9.5B	40	29	9954
	FLAP	9.5B	40	31.1 \pm 10.6	9570.8 \pm 2601.0
	LLM-Pruner	9.2B	40	26	8985
	Ours	9.5B	29	40	13824
37% [†]	Wanda-sp	8.4B	40	26	8710
	FLAP	8.3B	40	27.5 \pm 11.3	8326.6 \pm 2874.9
	LLM-Pruner	8.2B	40	22	7603
	Ours	8.3B	25	40	13824

[†]Reduction ratio for the number of parameters.
[‡]#Block: #Transformer blocks; #Head: #attention heads of MHA; FFN-D: intermediate size of FFN.

ters takes about 2 hours and utilizes 22GB GPU memory, while a 21%-pruned model from 13B parameters requires approximately 3 hours and 35GB VRAM.

- At the inference stage, we maintain default configurations without employing xFormers-optimized attention or additional advanced features.

D ZERO-SHOT DOWNSTREAM TASK PERFORMANCE

Table 6: Zero-shot results of the compressed (top) LLaMA-7B, (middle) Vicuna-7B-v1.3, and (bottom) Vicuna-13B-v1.3. The width pruning methods of Wanda-sp (Sun et al., 2024; An et al., 2024), FLAP (An et al., 2024), and LLM-Pruner (Ma et al., 2023) often degrade inference efficiency due to the GPU-unfriendly weight sizes (Andersch et al., 2023). On the contrary, our depth pruning method attains actual speedups while performing comparably in zero-shot task scenarios.

Model		PPL↓		Average	BoolQ	Commonsense Reasoning Accuracy↑ (%)					Thr↑ (tokens/s) [‡]		
		Wiki2	PTB			PIQA	HellaSwag	WinoGrande	ARC-e	ARC-c	OBQA	H100	RTX3090
LLaMA-7B (6.7B)		12.6	22.1	66.3	75.0	78.7	76.2	69.9	75.3	44.7	44.4	53.7	25.0
20% Pruned (5.5B)	Wanda-sp	21.4	47.2	51.8	61.5	70.4	53.2	56.0	58.7	31.4	31.0	41.7	16.7
	FLAP	17.0	30.1	59.5	69.4	74.7	66.9	66.3	64.6	36.5	38.2	40.5	16.5
	LLM-Pruner	17.6	30.4	61.8	66.2	77.6	71.4	66.1	70.5	39.3	41.2	43.2	21.4
	Ours: Taylor+	20.2	32.3	63.5	75.7	75.7	71.5	69.1	69.9	41.6	40.8	66.0	28.4
	Ours: PPL	17.7	30.7	61.9	72.7	75.7	70.4	63.6	69.5	40.1	41.2	66.0	28.4
27% Pruned (4.9B)	Wanda-sp	50.4	106.9	42.1	62.0	60.4	33.2	52.8	37.6	23.0	25.4	41.7	16.0
	FLAP	21.3	37.1	55.8	68.2	70.6	61.0	64.1	58.8	31.4	36.8	40.2	16.5
	LLM-Pruner	20.5	36.1	58.7	62.8	75.5	67.2	64.9	63.5	36.8	40.2	44.0	22.9
	Ours: Taylor+	29.9	42.0	59.8	70.6	73.0	65.7	68.5	63.9	39.3	37.4	73.9	34.9
	Ours: PPL	20.7	36.0	57.6	66.6	73.1	63.7	60.4	64.3	36.0	39.2	73.9	34.9
35% Pruned (4.5B)	Wanda-sp	133.6	210.1	36.9	44.5	56.8	29.6	49.6	31.7	20.7	25.6	41.6	16.1
	FLAP	25.6	44.4	52.7	68.3	68.1	55.9	61.1	52.3	29.4	33.8	40.5	15.8
	LLM-Pruner	24.2	40.7	55.5	62.9	72.8	62.3	62.7	57.4	33.0	37.6	44.4	21.1
	Ours: Taylor+	33.2	58.5	55.4	62.5	69.2	60.7	66.8	57.4	34.5	36.8	80.1	37.8
	Ours: PPL	23.1	38.8	55.2	64.3	71.4	59.4	59.3	62.2	32.8	37.0	80.1	37.8
Vicuna-7B (6.7B)		17.1	63.2	65.9	78.1	77.3	73.9	69.5	74.3	44.3	43.8	53.7	25.0
20% Pruned (5.5B)	Wanda-sp	24.4	104.0	58.5	63.9	72.0	67.4	65.2	64.8	38.3	37.8	41.7	16.7
	FLAP	22.0	74.9	61.4	73.1	74.8	67.9	65.8	67.5	40.2	40.6	40.5	16.5
	LLM-Pruner	19.6	76.4	60.1	65.4	76.2	68.9	64.4	68.9	37.4	39.4	43.2	21.4
	Ours: Taylor+	21.0	72.3	62.5	78.7	74.8	69.4	68.5	68.2	38.7	39.6	66.0	28.4
	Ours: PPL	18.8	67.9	60.7	71.7	74.4	67.6	63.6	69.3	38.9	39.4	66.0	28.4
27% Pruned (4.9B)	Wanda-sp	36.5	177.6	50.9	49.0	67.1	57.2	59.2	57.6	33.7	32.4	41.7	16.0
	FLAP	27.9	88.3	57.1	72.0	71.5	62.0	61.2	61.2	35.4	36.6	40.2	16.5
	LLM-Pruner	22.7	87.9	57.1	60.8	74.3	65.9	60.9	64.4	34.6	38.8	44.0	22.9
	Ours: Taylor+	29.8	92.0	60.2	78.8	71.8	64.4	67.7	64.3	36.4	37.6	73.9	34.9
	Ours: PPL	23.0	78.2	56.1	66.4	72.9	60.6	59.2	63.1	33.8	37.0	73.9	34.9
35% Pruned (4.5B)	Wanda-sp	73.2	386.5	39.4	43.1	58.4	36.3	53.3	34.5	23.7	26.4	41.6	16.1
	FLAP	34.6	104.8	53.7	65.1	68.1	57.0	63.1	56.9	32.0	34.0	40.5	15.8
	LLM-Pruner	27.6	102.0	53.5	52.0	72.4	61.6	59.9	58.0	33.3	37.0	44.4	21.1
	Ours: Taylor+	35.0	110.3	55.0	64.0	69.6	59.3	66.5	57.5	33.3	35.2	80.1	37.8
	Ours: PPL	26.6	89.4	53.3	65.2	70.4	56.5	56.6	59.8	31.5	33.4	80.1	37.8
Vicuna-13B (13.0B)		14.7	51.6	68.3	82.8	78.3	77.0	71.2	75.4	47.7	45.4	45.5	OOM
21% Pruned (10.5B)	Wanda-sp	19.0	71.8	63.6	78.6	75.6	73.5	68.4	68.5	42.2	38.4	34.1	12.9
	FLAP	18.8	65.3	63.3	77.2	75.1	72.0	70.2	69.4	40.3	38.8	32.6	12.6
	LLM-Pruner	16.0	57.0	65.3	75.5	78.6	75.0	69.8	70.6	43.6	44.4	34.0	17.3
	Ours: Taylor+	18.1	61.6	66.7	83.0	76.8	75.1	72.8	72.5	44.5	42.4	55.7	23.9
	Ours: PPL	16.1	56.5	64.9	75.0	77.1	73.7	68.9	71.5	43.8	44.2	55.7	23.9
29% Pruned (9.5B)	Wanda-sp	23.4	84.9	60.0	71.5	74.2	68.7	65.1	64.3	36.8	39.4	33.7	13.5
	FLAP	22.8	78.8	61.6	75.9	73.7	67.9	66.4	67.3	38.0	42.0	33.0	12.1
	LLM-Pruner	19.0	66.4	62.7	68.3	77.1	72.0	69.7	68.6	40.0	43.4	35.8	15.0
	Ours: Taylor+	22.0	70.3	65.1	82.6	75.1	73.3	70.9	69.9	43.8	40.2	62.0	24.2
	Ours: PPL	18.1	62.2	62.0	67.5	75.6	70.6	65.5	70.9	43.3	40.2	62.0	24.2
37% Pruned (8.3B)	Wanda-sp	36.6	123.5	52.7	59.6	67.5	59.5	59.7	55.2	33.5	33.8	33.8	12.6
	FLAP	28.7	96.2	58.3	72.5	70.0	62.5	65.4	63.8	36.3	37.8	32.9	13.2
	LLM-Pruner	22.2	74.0	59.7	67.1	75.6	67.7	63.2	65.5	38.8	39.8	35.6	18.0
	Ours: Taylor+	34.2	90.4	61.4	78.5	71.3	69.2	69.9	64.2	40.5	36.6	69.7	31.7
	Ours: PPL	22.1	73.6	59.1	69.4	73.8	64.4	62.5	65.1	39.2	39.0	69.7	31.7

[‡]Throughput measured with 12 input tokens, 128 output tokens, and a batch size of 1 on a single GPU.

E ABLATION STUDY

E.1 IMPORTANCE CRITERIA FOR BLOCK PRUNING

Table 7 presents the results of block pruning using various significance criteria. The basic methods without the ‘+’ label fail to incorporate essential initial blocks, causing a decline in performance. The Magnitude+ method, which preserves these critical blocks, partially improves the scores; however, its effectiveness is still inferior compared to the other methods, indicating that relying solely on weight magnitude could be improper for pruning decisions. The Taylor+ criterion enhances accuracy in commonsense reasoning tasks, while the PPL method leads to better generation quality without relying on heuristic selection of pruning candidates.

Table 7: Comparison of pruning criteria on LLaMA-7B. The Taylor+ method excels in commonsense reasoning accuracy, while the PPL criterion leads to better generation performance.

Block Pruning		PPL↓		Ave Acc↑
Criterion		WikiText2	PTB	(%) [†]
20% Pruned (5.5B)	Magnitude	7720.7	10618.7	34.4
	Magnitude+	19.4	36.3	56.1
	Taylor	3631.7	4327.9	35.5
	Taylor+	20.2	32.3	63.5
	PPL	17.7	30.7	61.9
35% Pruned (4.5B)	Magnitude	8490.1	14472.1	34.9
	Magnitude+	36.9	61.1	49.3
	Taylor	7666.8	10913.1	35.3
	Taylor+	33.2	58.5	55.4
	PPL	23.1	38.8	55.2

[†] Average accuracy on seven commonsense reasoning tasks.

E.2 STRUCTURAL UNIT FOR DEPTH PRUNING

Pruning individual MHA and FFN modules, which are more fine-grained units than Transformer blocks, is also possible. To examine its effect, we measure the impact of removing each module on the PPL of the calibration set and selectively eliminate the unnecessary modules. The same LoRA retraining procedure is conducted.

Table 8 shows the results of depth pruning at different granularities. For the models with more than 5B parameters, removing individual MHA and FFN modules results in better downstream task accuracy but worse PPL compared to removing entire Transformer blocks. For smaller models than 5B, block-level pruning achieves superior results in terms of all the examined metrics. This differs from the common belief that removing finer units yields better performance.

Given the collaborative roles of the modules (i.e., MHA captures dependency relations (Vaswani et al., 2017), while skip connections and FFN prevent the rank collapse in purely attention-driven networks (Dong et al., 2021)), it may be suboptimal to treat them in isolation. Taking the 5.3B model in Table 8 as an example, module-level pruning results in consecutive FFNs in some positions, potentially impairing the model’s ability to handle word interactions. In contrast, with block removal, the loss of information could be compensated by neighboring blocks that serve similar functions.

Table 8: Comparison of depth pruning granularities on LLaMA-7B. Removing entire Transformer blocks instead of individual MHA and FFN modules generally yields better results.

Depth Pruning Unit	#Param	PPL↓		Ave Acc↑
		WikiText2	PTB	(%) [†]
Individual MHA & FFN	5.7B	20.8	34.8	63.1
Transformer Block	5.7B	16.9	29.3	62.8
Individual MHA & FFN	5.3B	25.2	41.3	61.1
Transformer Block	5.3B	18.6	33.1	60.6
Individual MHA & FFN	4.6B	38.9	58.7	52.5
Transformer Block	4.5B	23.1	38.8	55.2
Individual MHA & FFN	4.0B	63.2	88.9	48.3
Transformer Block	3.9B	31.1	47.3	50.6

[†] Average accuracy on seven commonsense reasoning tasks.

E.3 CALIBRATION DATA VOLUME

The calibration set is employed to assess the weight significance of width pruning baselines and the block-level importance of our method during the pruning phase.

Table 9 presents the results obtained by varying the number of calibration samples in the Book-Corpus dataset. The scores remain relatively stable for the examined methods, suggesting that 10 samples could be sufficient. However, our Taylor+ method encounters a drop in downstream task accuracy when 1K samples are used, leaving the exploration of calibration data characteristics for future research.

Table 9: Impact of calibration data volume. The results of 20%-pruned LLaMA-7B are reported.

Evaluation Metric	Method	# Calibration Samples			
		10	50	100	1000
PPL \downarrow on WikiText2	Wanda-sp	21.4	21.4	21.7	20.8
	FLAP	17.0	17.5	17.5	17.3
	LLM-Pruner	17.6	17.2	17.0	OOM \ddagger
	Ours: Taylor+	20.2	20.2	19.0	19.6
	Ours: PPL	17.7	17.2	17.4	17.4
Ave Acc \uparrow (%) \dagger	Wanda-sp	51.8	52.9	52.0	53.0
	FLAP	59.5	59.7	59.9	60.8
	LLM-Pruner	61.8	61.6	61.7	OOM \ddagger
	Ours: Taylor+	63.5	63.5	63.9	61.7
	Ours: PPL	61.9	61.5	61.7	61.7

\dagger Average accuracy on seven commonsense reasoning tasks.

\ddagger Out-of-memory error on an A100 (80GB) using the official code.

E.4 ONE-SHOT vs. ITERATIVE PRUNING

For one-shot pruning, multiple blocks are removed simultaneously from the original model, followed by just one phase of retraining. For iterative pruning, the removal of one block coupled with subsequent retraining is repeatedly performed. Here, we use the PPL-based importance criterion for selecting which blocks to remove.

Figure 5 compares the pruned networks before and after the retraining process. The iteratively pruned models yield better post-pruning results than one-shot pruned ones. However, a single retraining session after one-shot pruning leads to similar performance with iterative pruning. In light of the greatly reduced retraining budget, we opt for one-shot pruning.

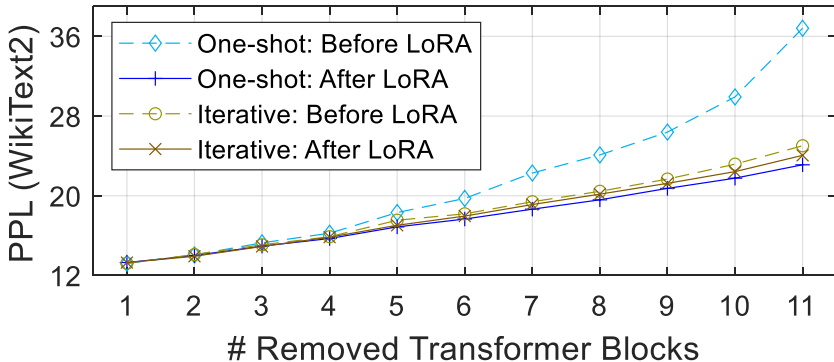


Figure 5: Comparison of one-shot and iterative block pruning on LLaMA-7B. The retraining phase of one-shot pruning effectively bridges the performance gap with iterative pruning.

F GENERATION EXAMPLES

Tables 10 and 11 present generation examples where some input prompts were sourced from Ma et al. (2023). In terms of linguistic flow and topical consistency, the pruned models yield sentences on par with those from the original model. However, as also noted in Ma et al. (2023), the output quality deteriorates when responding to factual questions or producing longer content. To overcome this, integrating more powerful retraining methods may be a valuable future direction.

Table 10: Generation examples from the original LLaMA-7B and 20%-compressed models.

Model	Generated Text
LLaMA-7B	What’s great about the holiday season, it’s a time when people focus less on work and more on families and friends. The holidays are a time when people can really focus on giving and helping others, which is something that the holidays are supposed to be about. In my blog post this week, I’ll be taking a look at a variety of causes that ...
Wanda-sp	What’s great about the holiday season, though, is its warm and inviting people, and what a chance to meet and share a cuppa with others who understand and have shared similar experiences, a feeling of empathy if you will, and a great opportunity to help others on a difficult journey. So I can say, all is well. But I’m not giving up the fight ...
FLAP	What’s great about the holiday season, and what makes it such an awesome time of the year, is that it’s the time of the year when people are so generous, and so goodhearted. When you get together with friends and family, and you go through the various things that you’re able to do, when you’re able to give people gifts, when people can see ...
LLM-Pruner	What’s great about the holiday season, is you can enjoy the company of friends and family. This is the perfect opportunity to gather together, and share all the special memories that you have together. These moments are the ones that make the holiday special. No matter what happens, or what you end up doing, it’s always good to get ...
Ours: Taylor+	What’s great about the holiday season, is that it’s a time of year when people can actually take a step back from the normal day-to-day routines and take some time out to relax and enjoy some well-deserved rest and rejuvenation. This might include indulging in a more comfortable or luxury holiday experience and this article aims to help ...
Ours: PPL	What’s great about the holiday season, in my mind, is that it brings together friends, family, and food. The holiday season is a great time to show some love and have some fun with the people that matter most. Whether you are celebrating Christmas, Hanukah, or New Years, the festivities will be full of joy and memories that will last ...
LLaMA-7B	AI can create a logo in seconds. Why are these new technologies threatening our jobs? AI can do most menial tasks in just a few seconds or minutes. These tasks can be very tedious for human beings to do. For instance, AI can identify patterns in large amounts of data so that complex statistical analysis can be carried out faster than a ...
Wanda-sp	AI can create a logo in seconds rather than years. To this point, this is only the beginning of what AI will do to change our human lives. As time goes on, it’s quite likely that we’ll see even greater use of robots. Even if this all sounds too dramatic, this may still happen. And AI is the best thing that can make it to happen. in 2019. All you ...
FLAP	AI can create a logo in seconds and it doesn’t cost you a penny. The new era of technology has opened up many exciting opportunities. Technology, like AI is being utilized for different purposes in different ways, and now, technology is becoming a solution to problems that would have been tedious or expensive only a few years ago. ...
LLM-Pruner	AI can create a logo in seconds, design a website from scratch and even write blog posts without any human intervention. This revolution is changing the way humans communicate and express ideas. In our view, this technology is a great asset but we’ll need to be smart in how we use it to prevent abuse or automation leading to job loss. So, ...
Ours: Taylor+	AI can create a logo in seconds or design a new website. There is a lot of software that exists in these two categories. It has become easier and faster for people to use. However, there is still the requirement for the designer or creative director to ensure that the designs are relevant and meaningful to the brand. The Future: AI for creative design ...
Ours: PPL	AI can create a logo in seconds, and it can design entire websites, all without human intervention. However, there is still a lot of work to be done to perfect this technology, including the development of natural language processing, machine learning and computer vision, as well as ensuring that AI designers are creating work that can be ...
LLaMA-7B	Neural network pruning is defined as the process of removing certain weights or neurons in a neural network to reduce the size of the neural network, and improve the training process. Neural networks pruning techniques have been applied successfully on classification, optimization, planning, and speech recognition tasks. A pruned ...
Wanda-sp	Neural network pruning is defined as any method that removes the redundant information by applying the neural network network. It is the term used to identify the systematic variables when an individual tries to find an algorithm. In addition, it is the method of classification used to define any set of algorithms that have been ...
FLAP	Neural network pruning is defined as a selection of neurons or weights in a neural network so that the information of the network is reduced to a specified level, and a process of re-training a neural network. The latter one is the most popular method, especially the neural network pruning. In recent years, this field has been studied in ...
LLM-Pruner	Neural network pruning is defined as reducing the complexity and size of an over-parameterized neural network. The pruning reduces the FLOPs and improves performance for the particular task. Pruning techniques have gained popularity in recent years due to their wide range of applications. This paper aims to explore the pruning ...
Ours: Taylor+	Neural network pruning is defined as reducing the weight connections and biases of a neural network, while retaining high performance. Pruning is an effective approach to reducing the number of parameters and improving the efficiency of neural networks. Recently, the concept of neural network shrinking has gained considerable ...
Ours: PPL	Neural network pruning is defined as the process of removing redundant neurons and connections between nodes in a neural network, while keeping the performance or accuracy to the maximum extent possible. This optimization is aimed at reducing the size and memory consumption of the network while maintaining its ability to perform task ...

Table 11: Generation examples from the original Vicuna-13B-v1.3 and 21%-compressed models.

Model	Generated Text
Vicuna-13B	What’s great about the holiday season, beyond the obvious time off and the chance to spend time with loved ones, is the opportunity to slow down and reflect on the year that’s passed. It’s a time to look back at all the accomplishments and growth you’ve experienced, and to set your sights on the year ahead. As a small business owner, the ...
Wanda-sp	What’s great about the holiday season, besides spending time with family and friends, is that it’s a time of giving. It’s a time to show appreciation for those in your life who make a difference and bring joy to your world. Whether it’s a gift of love, a gift of time, or a gift of money, it’s important to remember that giving is not just about ...
FLAP	What’s great about the holiday season, especially the days leading up to Christmas, is the opportunity to take a pause and recharge. It’s a time to reflect on the year and all the experiences we’ve had, good or bad. It’s a time to reflect on the people we love and the impact they’ve had on our lives. It’s a time to reflect on what we want ...
LLM-Pruner	What’s great about the holiday season, aside from the gifts and the holiday parties, is the opportunity to spend time with family and loved ones. One of the things that I love most about holiday seasons is the traditions that surround them, from decorating the tree to making special holiday meals, there’s a sense of nostalgia and warmth that ...
Ours: Taylor+	What’s great about the holiday season, is that it allows us to reflect on the year gone by and reflect on the moments that made it special. And while some people might choose to celebrate with the traditional traditions, some choose to do something different. One way to spice up your holiday season is by taking a trip to a far-off place. Why not ...
Ours: PPL	What’s great about the holiday season, in addition to the fun food, festivities, and spending time with loved ones, is that it creates the perfect opportunity to show people how much you appreciate them. From giving them a gift, to baking their favorite treat, or just taking a moment to acknowledge them in front of others – there’s something ...
Vicuna-13B	AI can create a logo in seconds with its complex algorithms and pre-set designs. But it can never truly understand the essence of a business, the emotions that go into creating a logo or the creativity and uniqueness that comes from human input. It’s time to put the human touch back into branding, where a designer and client can collaborate to ...
Wanda-sp	AI can create a logo in seconds using a logo design tool. However, a talented graphic designer can create a unique and meaningful logo that stands out from the crowd, captures the essence of the business, and appeals to the target audience. A good logo can also make a difference in how customers perceive the business and can even increase ...
FLAP	AI can create a logo in seconds. If you don’t like it, you can change it yourself or pay for a designer to do it for you. However, it’s important to note that a logo is more than just an image. It’s a visual representation of a brand or organization, and it communicates certain values and characteristics to the viewer. Therefore, it’s essential to ...
LLM-Pruner	AI can create a logo in seconds with just one click, but design is a human art, something that is beyond AI capabilities. The logo design is an important part of a company’s identity, and it’s a decision that affects the company’s visual branding and reputation for years to come. The logo is the visual representation of a company’s values ...
Ours: Taylor+	AI can create a logo in seconds, with or without a human’s input. Some people might like this, because they can create a logo for a small business on a budget. Other people, like myself, think that the human touch is important in art and design. A human designer can create a unique and meaningful logo by incorporating personal ...
Ours: PPL	AI can create a logo in seconds by combining a variety of different styles and shapes and patterns that make up the logos. However, a logo that is truly a reflection of your brand needs to be designed by a human, a graphic designer. They will bring a creative vision that will make your logo unique, appealing and memorable. They can add ...
Vicuna-13B	Neural network pruning is defined as the process of removing redundant or unnecessary connections in a neural network to reduce its size and computational complexity while preserving its performance. The primary motivation behind pruning is to improve the efficiency of deep learning models, especially in resource-constrained ...
Wanda-sp	Neural network pruning is defined as a method for removing connections in order to reduce the model size, thus allowing for computationally efficient models while retaining accuracy. This technique is useful in reducing the number of neurons in a neural network, which allows for fewer computations during training and reduces the ...
FLAP	Neural network pruning is defined as the process of reducing the size of a neural network by removing the least significant weights or neurons. The main motivation for pruning is to reduce the computational cost of neural networks while maintaining their accuracy. The key challenge of pruning is to identify the most effective ...
LLM-Pruner	Neural network pruning is defined as a technique that aims to reduce the size of a neural network by removing redundant and irrelevant connections between the neurons in the network. This approach is based on the observation that a large portion of the connections within the network is redundant and does not contribute to the overall ...
Ours: Taylor+	Neural network pruning is defined as the removal of redundant connections within a neural network to achieve a better model fit while retaining the network’s general accuracy. The goal of pruning is to reduce the computational cost and memory footprint of the network. One commonly used pruning method is called weight magnitude ...
Ours: PPL	Neural network pruning is defined as the task of removing unnecessary or redundant connections in a neural network while retaining its accuracy and performance. This is often done to reduce the memory usage and computational complexity of a neural network, which can be critical when running on devices with limited resources. In ...

G ADDITIONAL RESULTS OF INFERENCE EFFICIENCY

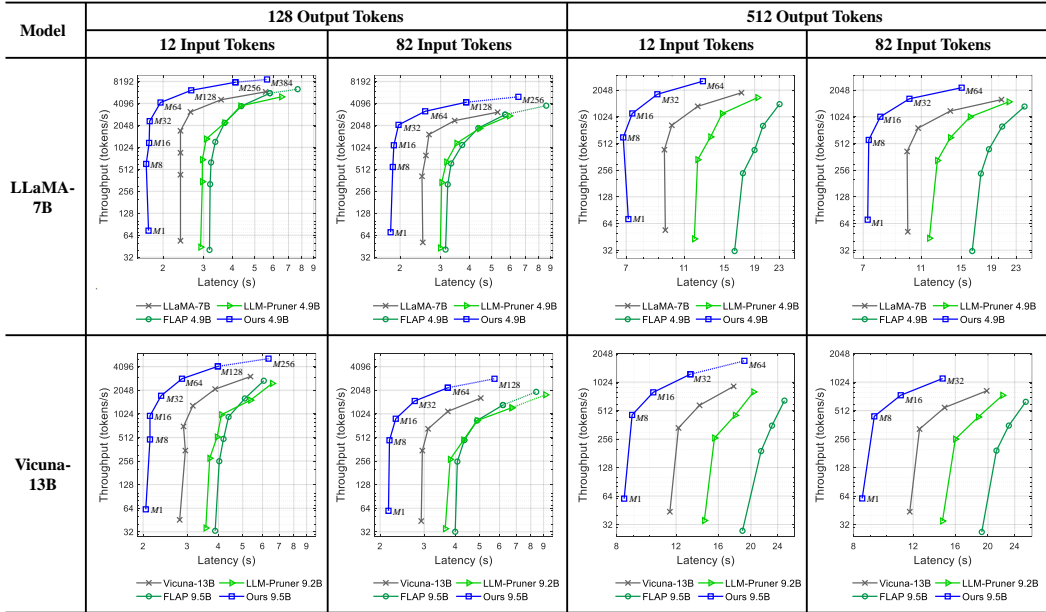


Figure 6: Inference efficiency of pruned models on an NVIDIA H100 GPU. Our depth pruning achieves a superior latency-throughput trade-off for various sequence lengths of input and output. In contrast, the width pruning of FLAP (An et al., 2024) and LLM-Pruner (Ma et al., 2023) degrades efficiency results due to GPU-unfriendly weight dimensions (Andersch et al., 2019) (e.g., the hidden sizes of FFN are often not divisible by 8). The markers labeled with M represent batch sizes. The dotted lines indicate that pruned models can operate with larger batch sizes, avoiding out-of-memory errors encountered by the original model.

H GPU MEMORY REQUIREMENTS

Table 12 shows the gains in GPU memory requirements from our depth-pruned models on NVIDIA H100 given 12 input tokens. The larger the batch size, the greater the improvement observed. Notably, our pruned models can handle an output length of 512 and a batch size of 64, unlike the original 13B-parameter model.

Table 12: GPU memory requirements for varying sequence lengths (L) and batch sizes (M). The results of the 7B and 13B models and our models with different pruning ratios are reported. Our approach effectively reduces the memory demands of the original models.

Model	$L128$			$L512$		
	$M1$	$M16$	$M64$	$M1$	$M16$	$M64$
7B	12.8GB	16.0GB	25.8GB	13.3GB	25.0GB	61.8GB
20%	10.5GB	13.1GB	21.1GB	10.9GB	20.4GB	50.4GB
27%	9.4GB	11.6GB	18.8GB	9.7GB	18.1GB	44.6GB
35%	8.6GB	10.7GB	17.2GB	9.0GB	16.6GB	40.8GB
13B	24.8GB	29.6GB	44.9GB	25.5GB	43.7GB	OOM
21%	19.9GB	23.8GB	36.0GB	20.5GB	35.0GB	OOM
29%	18.1GB	21.7GB	32.7GB	18.6GB	31.8GB	73.5GB
37%	15.7GB	18.8GB	28.3GB	16.1GB	27.5GB	63.5GB