

# A Simple and Effective $L_2$ Norm-Based Strategy for KV Cache Compression

Alessio Devoto<sup>‡\*</sup> Yu Zhao<sup>†\*</sup> Simone Scardapane<sup>‡</sup> Pasquale Minervini<sup>†§</sup>  
<sup>‡</sup>Sapienza University of Rome <sup>†</sup>The University of Edinburgh <sup>§</sup>Miniml.AI  
{alessio.devoto, simone.scardapane}@uniroma1.it  
{yu.zhao, p.minervini}@ed.ac.uk

## Abstract

The deployment of large language models (LLMs) is often hindered by the extensive memory requirements of the Key-Value (KV) cache, especially as context lengths increase. Existing approaches to reduce the KV cache size involve either fine-tuning the model to learn a compression strategy or leveraging attention scores to reduce the sequence length. We analyse the attention distributions in decoder-only Transformers-based models and observe that attention allocation patterns stay consistent across most layers. Surprisingly, we find a clear correlation between the  $L_2$  norm and the attention scores over cached KV pairs, where a low  $L_2$  norm of a key embedding usually leads to a high attention score during decoding. This finding indicates that *the influence of a KV pair is potentially determined by the key embedding itself before being queried*. Based on this observation, we compress the KV cache based on the  $L_2$  norm of key embeddings. Our experimental results show that this simple strategy can reduce the KV cache size by 50% on language modelling and needle-in-a-haystack tasks and 90% on passkey retrieval tasks without losing accuracy. Moreover, without relying on the attention scores, this approach remains compatible with FlashAttention, enabling broader applicability.

## 1 Introduction

Handling long contexts is desirable for large language models (LLMs), as it allows them to perform tasks that require understanding long-term dependencies (Liu et al., 2024; Fu et al., 2024; Chen et al., 2023; Staniszewski et al., 2023; Zhao et al., 2024; Tworowski et al., 2024). A key component for modelling long context is the KV cache, which stores the keys and values of past tokens in memory to avoid recomputing them during generation.

\*Equal contribution.

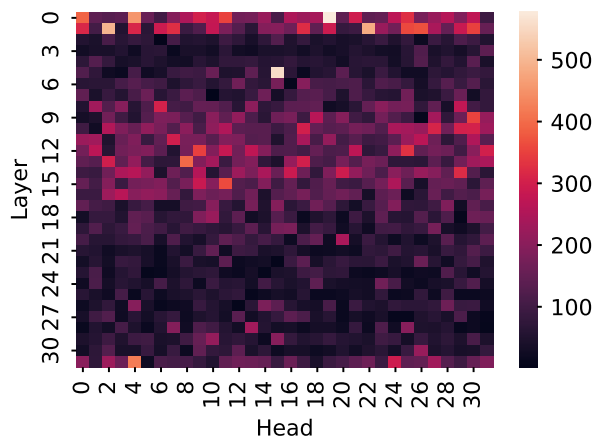


Figure 1: ALR, as defined in Eq. (3), for each head and layer in Llama2-7b. A lower value means a higher correlation between  $L_2$  norm and attention score.

However, processing long-context inputs often results in a high decoding latency since it requires repeatedly reading a potentially large KV cache from high-bandwidth memory (HBM) to the streaming multiprocessor (SM) during decoding (Fu, 2024). Consequently, the practical deployment of LLMs is frequently hindered by hardware limitations. To address the issue of KV cache growth, various KV cache compression methods have been proposed. These methods can be broadly categorised into trainable approaches, which involve modifications to the model architecture (Ainslie et al., 2023), or fine-tuning regime to inherently manage KV cache size (Nawrot et al., 2024), and non-trainable approaches, which apply post-hoc compression techniques to reduce the cache footprint without altering the underlying model (Li et al., 2024; Zhang et al., 2024b). While these methods have shown promise, they often involve complex algorithms or significant computational overhead, limiting their practicality; for example, post-hoc compression algorithms usually evict KV pairs based on attention scores, which is not compatible with FlashAttention (Dao et al., 2022) and thus prevents their

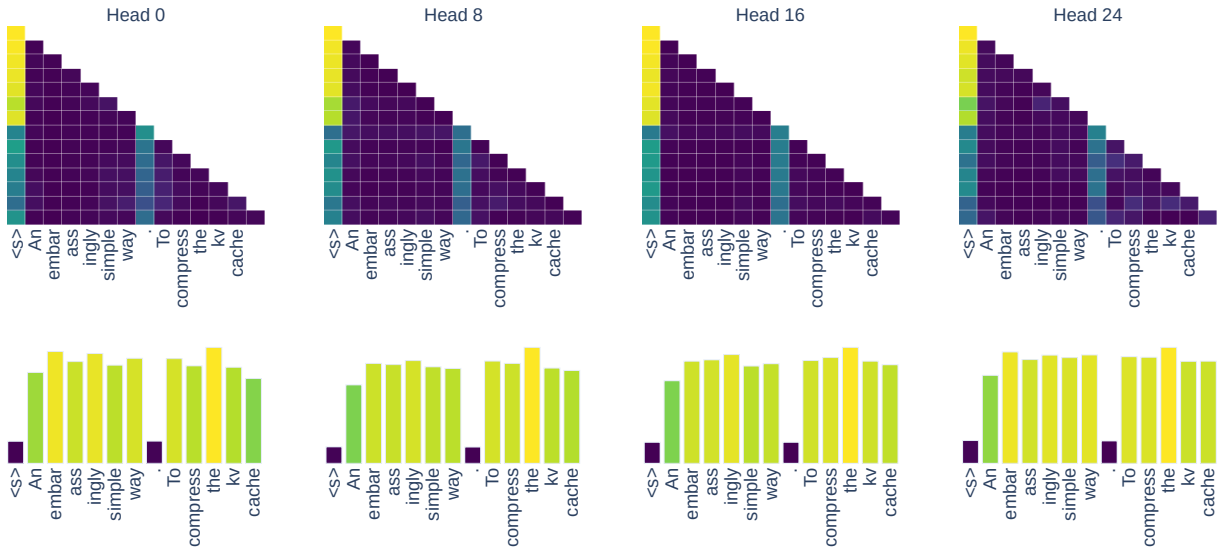


Figure 2: Five heads at layer 9 of Llama2-7b. Attention score (top) and  $L_2$  norm (bottom) are highly correlated. We observe similar patterns across most layers and for a wide range of inputs. More examples provided in Appendix D

applications in modern LLMs inference systems.

We show that the  $L_2$  norm of cached keys has a high correlation with attention scores. More specifically, we observe that a low  $L_2$  norm of a key embedding usually leads to a high attention score during decoding. Based on this observation, we propose a simple and highly effective strategy for KV cache compression: *keeping in memory only the keys with lowest  $L_2$  norm, and the corresponding values*. Unlike many existing methods, our heuristic can be applied off-the-shelf to any transformer-based decoder-only LLM without the need for additional training or significant modifications. More importantly, our method estimates the influence of cached key-value pairs without the need to compute the attention scores. Therefore, unlike other compression methods (Holmes et al., 2024; Li et al., 2024), it can be easily integrated with the popular FlashAttention (Dao et al., 2022).

Our experimental results demonstrate that this heuristic allows maintaining model performance in language modelling tasks and in tasks that require the model to store and retrieve the most critical information, such as passkey retrieval (Mohtashami and Jaggi, 2023) and needle-in-a-haystack tasks (Kamradt, 2023).

## 2 Patterns in the Attention Matrix

We first examine the attention scores on the language modelling task for a range of popular LLMs. By analysing the key embeddings and the attention

distribution, we observe that key embeddings with low  $L_2$  norm are often associated with higher attention scores. In Fig. 2, we provide an example using Llama-2-7b (Touvron et al., 2023), where the first row presents the attention distribution over the KV pairs, and the second row presents the  $L_2$  norm of each key embedding. We observe that the tokens with high attention scores, such as " $\langle s \rangle$ " and ".", have significantly lower  $L_2$  norm values than others. While Xiao et al. (2024) already observed peaked attention distributions for specific tokens, and Darcet et al. (2024) pointed out the influence of high  $L_2$  norm tokens on attention maps, we are the first, to the best of our knowledge, to point out the correlation between the  $L_2$  norm of the key embeddings and attention score. Based on our observation, we consider the following research question: can we compress the KV cache based on the  $L_2$  norm of the key embeddings?

An intuitive way to estimate the influence of compressing the KV cache is by examining the attention scores that are dropped due to the compression. In the following, we formally define this influence.

Given a prompt consisting of  $n$  tokens  $(x_1, x_2, \dots, x_n)$ , the LLM first encodes them into a KV cache—this step is referred to as the *pre-filling phase*. Then, the model autoregressively generates the next token  $x_{n+1}$ . When performing KV cache compression, some key-value pairs may be dropped and thus cannot be attended to. We define the attention loss caused by the compression as the sum

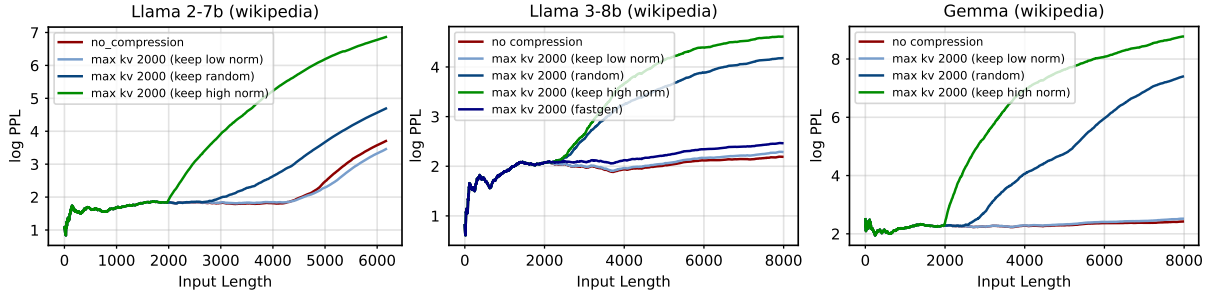


Figure 3: Perplexity for Llama 2-7b, Llama 3-8b and Gemma on language modelling task on wikipedia dataset. Additional results on coding dataset are available in Appendix B

of the attention scores associated with the dropped KV pairs:

$$\mathcal{L}_{l,h}^m = \sum_{p \in D_{l,h}} a_{l,h,p}, \quad (1)$$

where  $a_{l,h,p}$  is the attention score of the  $p$ -th token in the layer  $l$ , head  $h$ . In Eq. (1),  $D_{l,h}$  denotes the positions of  $m$  pairs of dropped KV,  $|D_{l,h}| = m$ , which depends on the compression method. An ideal compression algorithm aims to drop the KV pairs with the lowest attention scores, which will have less impact on the output. However, such attention scores are unavailable for a compression algorithm since it needs  $x_{n+1}$  to query the full KV cache in advance. Instead, we drop KV pairs with the highest  $L_2$  norm in key embeddings and use attention loss caused by ideal compression as the reference:

$$\mathcal{Y}_{l,h}^m = \mathcal{L}_{l,h}^m - \mathcal{L}_{l,h}^{m,ref}, \quad (2)$$

where  $\mathcal{L}_{l,h}^{m,ref}$  is the reference attention loss, and  $\mathcal{Y}_{l,h}^m$  is a non-negative value. A lower  $\mathcal{Y}_{l,h}^m$  indicates a lower difference and thus a higher correlation between the attention score and the  $L_2$  norm. To measure the overall difference between ideal attention score-based compression and  $L_2$  norm-based compression, we sum up the  $\mathcal{Y}_{l,h}^m$  over different numbers of compressed KV pairs:

$$\mathcal{Y}_{l,h} = \sum_{m=1}^n \mathcal{Y}_{l,h}^m. \quad (3)$$

We name the  $\mathcal{Y}_{l,h}$  as ALR, which indicates the Attention Loss for a compression method using ideal attention loss as Reference.

In Fig. 1, we plot the  $\mathcal{Y}$  across layers and heads. We observe that heads in the first two layers and some middle layers around the 12th layer have relatively high  $\mathcal{Y}$  values. The heads in other layers have lower  $\mathcal{Y}$  values, indicating a high correlation

between  $L_2$  norm and attention score. By leveraging this correlation, we can compress the KV cache based on the  $L_2$  norm of key embeddings. Optionally, we can skip the compression at the layers with low correlation. We show ablation experiments skipping layers in Appendix B.

### 3 Experiments

We evaluate our method on language modelling and two long-context modelling tasks, i.e., needle-in-a-haystack and passkey retrieval. Based on the observation supported by Fig. 1, the heads in the first two layers usually have a low correlation between  $L_2$  norm and attention score, so we do not perform compression on these layers as default. We conduct experiments to investigate the impact of compression on different layers in Appendix A.

**Language Modelling** For language modelling, we let the KV cache grow until a specific pre-defined length and subsequently start to discard the tokens with the highest  $L_2$  norm. We show in Fig. 3 that evicting even up to the 50% of KV Cache does not impact perplexity. Perplexity increases, as expected, once we exceed the pre-training context length. We show more results, including next token accuracy in Appendix B.

To further verify that keys with low  $L_2$  norm capture significant information, we test other eviction strategies, i.e. keeping tokens with highest  $L_2$  norm and keeping random tokens. It is clear from Fig. 3 that discarding tokens with low  $L_2$  impairs performance, even more so than random discarding, thus highlighting the importance of these low  $L_2$  norm keys.

**Pressure Test on Long-Context Tasks** The needle-in-a-haystack task (Kamradt, 2023) and passkey retrieval task (Mohtashami and Jaggi, 2023) are two synthetic tasks that are widely used

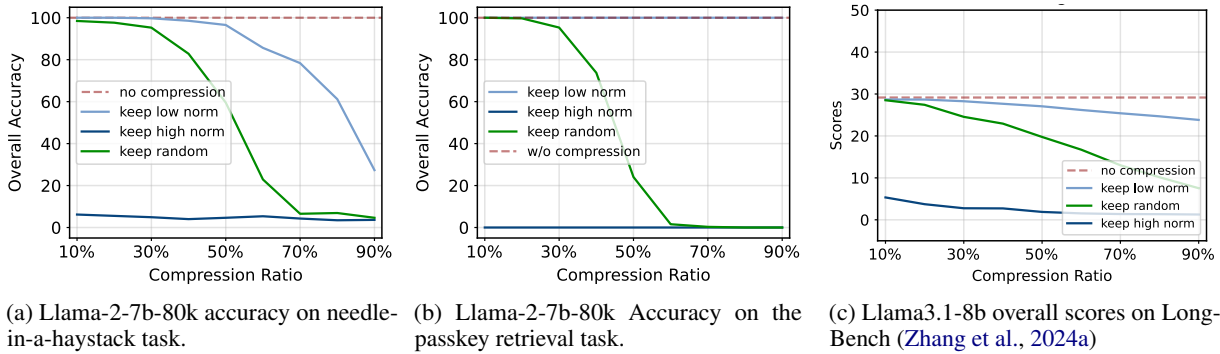


Figure 4: Score on long-context tasks for Llama-2-7b-80k and Llama 3.1

to pressure test the long-context modelling capability of LLMs. In both tasks, the model needs to identify and retrieve the important information from a long context to generate correct answers. Thus, these tasks test the compression method’s ability to keep important KV pairs and drop redundant ones.

In Figure 4a and Figure 4b, we present the experimental results of Llama-2-7b-80k (Fu et al., 2024). We analyse additional models in Appendix C. The model retains its performance on the needle-in-a-haystack task with 30% KV cache compression (Figure 4a) and maintains 99% accuracy with 50% compression. It also achieves 100% accuracy on the passkey retrieval task, even with 90% KV cache compression (Figure 4b).

We compare other eviction strategies, like keeping KV pairs with low  $L_2$  norm, with high  $L_2$  norm, and random. In Figure 4a and Figure 4b, we observe that the model cannot answer correctly when keeping only high  $L_2$  norm KV pairs, obtaining near zero accuracy. When we randomly compress the KV cache, the performance decreases significantly faster than keeping low  $L_2$  norm KV pairs. The above analysis indicates that KV pairs with low  $L_2$  norm are critical to generating the correct answer and thus contain important information.

**Experiments on LongBench** Additionally, we evaluate on LongBench (Zhang et al., 2024a). We test on several subsets, including NarrativeQA (Kociský et al., 2018), Qasper (Dasigi et al., 2021), HotpotQA (Yang et al., 2018), 2WikiMQA (Ho et al., 2020), and QMSum (Zhong et al., 2021). We report the results for the recently released long context Llama3.1 in Figure 4c. In addition, we show the complete per-subset results in Appendix C. The experimental results show that compressing the KV cache with low  $L_2$  norm only introduces a small ac-

curacy decrease even when compressing 50% KV cache, while compressing KV cache with high  $L_2$  norm results in almost zero accuracy.

**Comparison with FastGen** Like the majority of methods in the literature, FastGen (Holmes et al., 2024) utilises attention scores, which makes it incompatible with the popular FlashAttention (Dao et al., 2022), thereby limiting its efficiency and usability. For a fair comparison, we implement FastGen without using the attention scores, i.e., we only consider local, punctuation and special tokens. We perform experiments on language modelling with the Llama3 model (Dubey et al., 2024). Our method still outperforms FastGen with up to 50% KV cache eviction. We show the results in Figure 5.

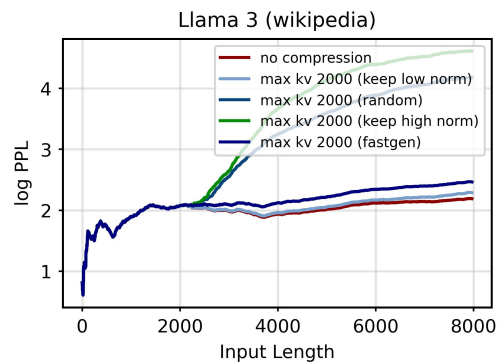


Figure 5: Perplexity of Llama3-8b on the wikipedia dataset when compared to FastGen (only local, special and punctuation tokens).

**Relationship between embedding and  $L_2$  norm**

After identifying a correlation between the  $L_2$  norm of token key embeddings and attention scores, we performed a further exploration by analyzing the key prjections in the KV cache. We found that tokens with lower  $L_2$  norm show sparse activations, with few dimensions having high values while most remain near zero, indicating limited use of the vec-

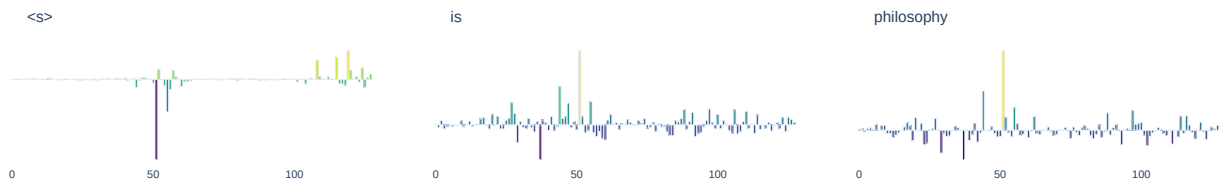
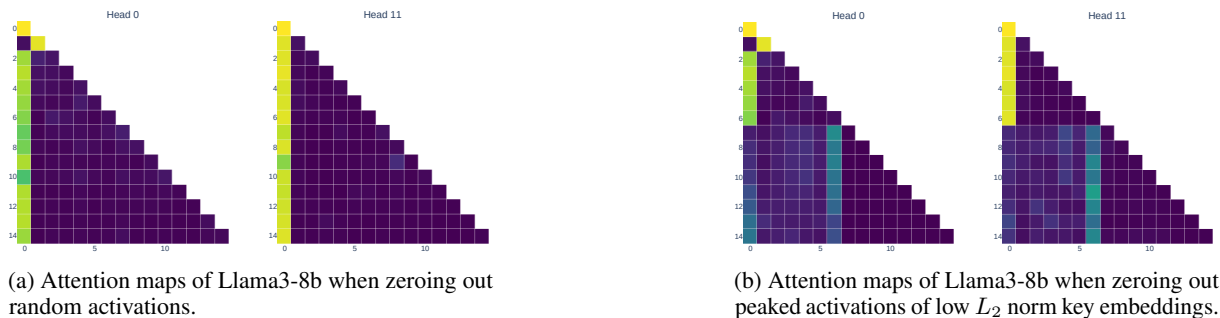


Figure 6: Key projections of the bos token  $\langle s \rangle$  vs other tokens. Each value represents the activation in a specific dimension for the embedding of the key projection. We found similar patterns across almost all heads and layers and in multiple texts. Only a few peaked activations ( $\sim 50$ ,  $\sim 56$  and  $\sim 120$ ) control the attention mechanism (see Figure 7). More plots like this in Appendix E



(a) Attention maps of Llama3-8b when zeroing out random activations.

(b) Attention maps of Llama3-8b when zeroing out peaked activations of low  $L_2$  norm key embeddings.

Figure 7: How the attention maps change if we set to zero a random activation (a) vs the specific peaked activations in the keys (b). In this example we set the values at iteration 5 during generation.

tor space (Figure 6). This sparsity aligns with the concept of "sink" tokens (Xiao et al., 2024), where many queries align with certain tokens, increasing their attention scores. We hypothesise that the lower  $L_2$  norm reflects a partial use of the available embedding space, leading to increased attention for these tokens. To test this, we zeroed out the dimensions responsible for the peaked activations in low  $L_2$  norm key embeddings and observed significant changes in attention maps during generation (Figure 7). However, randomly altering dimensions did not produce the same effect. This finding suggests that the  $L_2$  norm may serve as a proxy for the extent to which an embedding utilises the available vector space and, consequently, the degree to which it influences attention. Lower  $L_2$  norm appears to correspond to embeddings that drive disproportionately high attention values due to their alignment with a common "sink" direction.

## 4 Related Works

Recent long-context LLMs like Gemini-Pro-1.5 (Reid et al., 2024), Claude-3 (Anthropic, 2024), and GPT4 (Achiam et al., 2023) can process hundreds of thousands of tokens, but face high decoding latency. To address this, works like PageAttention (Kwon et al., 2023), Infinite-LLM (Lin et al., 2024), and vAttention (Prabhu et al., 2024) propose efficient memory management strategies. Others

focus on KV cache compression: DMC (Nawrot et al., 2024) uses dynamic trainable token merging, while H2O (Zhang et al., 2024b), FastGen (Ge et al., 2023), and SnapKV (Li et al., 2024) employ various attention-based training-free compression strategies. Unlike these methods, we uniquely use key embedding  $L_2$  norm for compression. While (Darcet et al., 2024) had previously found that high  $L_2$  norm hidden states aggregate important information, we are the first, to the best of our knowledge, to discover and leverage the correlation between low  $L_2$  norm key embeddings and high attention scores for efficient KV cache compression.

## 5 Conclusion

We introduced a simple yet highly effective strategy for KV cache compression in LLMs based on the  $L_2$  norm of key embeddings. We show that there is a significant correlation between the  $L_2$  norm of a key embedding and its attention score. Leveraging this observation, we compress the KV cache by retaining only those keys with the lowest  $L_2$  norm. Our experimental results on various tasks show that our compression strategy maintains the predictive accuracy of the model while significantly reducing the memory footprint. Our approach is straightforward and can be applied directly to any transformer-based, decoder-only LLM.



## 6 Limitations

While our research offers valuable insights, we tested only on relatively small models (Llama family and Gemma up to 8 billion parameters). In future work, we will assess our method on larger-scale models to ensure our findings generalize. Additionally, while we show that the  $L_2$  norm played a significant role in our experiments, we do not have a comprehensive theoretical explanation for why this is the case. Understanding the underlying reasons behind the importance of the  $L_2$  norm would require further theoretical exploration and empirical validation. Finally, we observed (Figure 1) that compressing based on  $L_2$  norm can be less effective depending on the layer and head considered, and we intend to investigate per-head compression ratios to leverage this observation.

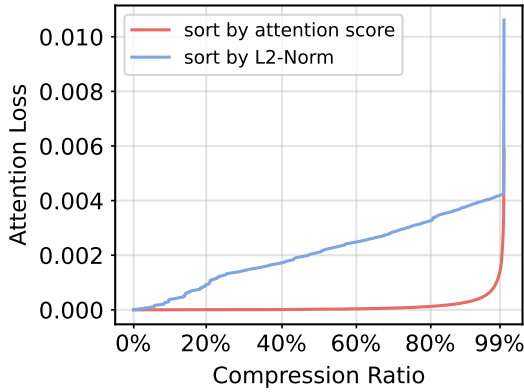
## 7 Acknowledgments

This work was supported by Sapienza Grant RM1221816BD028D6 (DeSMOS). Yu Zhao was supported in part by the UKRI Centre for Doctoral Training in Natural Language Processing, funded by UK Research and Innovation (grant EP/S022481/1) and the University of Edinburgh, School of Informatics.

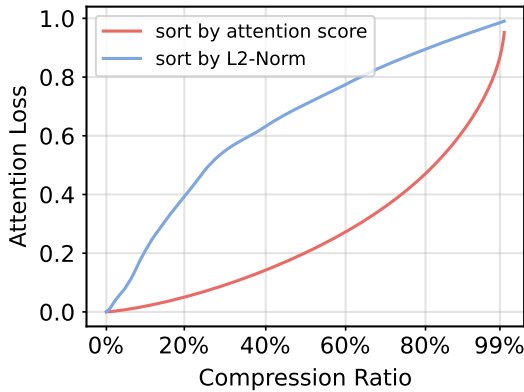
## References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Alvenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Joshua Ainslie, James Lee-Thorp, Michiel de Jong, Yury Zemlyanskiy, Federico Lebron, and Sumit Sanghai. 2023. [GQA: Training generalized multi-query transformer models from multi-head checkpoints](#). In *The 2023 Conference on Empirical Methods in Natural Language Processing*.
- AI Anthropic. 2024. The claude 3 model family: Opus, sonnet, haiku. *Claude-3 Model Card*.
- Yukang Chen, Shengju Qian, Haotian Tang, Xin Lai, Zhijian Liu, Song Han, and Jiaya Jia. 2023. Longlora: Efficient fine-tuning of long-context large language models. *arXiv preprint arXiv:2309.12307*.
- Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. 2022. FlashAttention: Fast and memory-efficient exact attention with IO-awareness. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- Timothée Darcet, Maxime Oquab, Julien Mairal, and Piotr Bojanowski. 2024. [Vision transformers need registers](#). In *The Twelfth International Conference on Learning Representations*.
- Pradeep Dasigi, Kyle Lo, Iz Beltagy, Arman Cohan, Noah A. Smith, and Matt Gardner. 2021. [A dataset of information-seeking questions and answers anchored in research papers](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021*, pages 4599–4610. Association for Computational Linguistics.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Yao Fu. 2024. Challenges in deploying long-context transformers: A theoretical peak performance analysis. *arXiv preprint arXiv:2405.08944*.
- Yao Fu, Rameswar Panda, Xinyao Niu, Xiang Yue, Hananeh Hajishirzi, Yoon Kim, and Hao Peng. 2024. Data engineering for scaling language models to 128k context. *arXiv preprint arXiv:2402.10171*.
- Suyu Ge, Yunan Zhang, Liyuan Liu, Minjia Zhang, Jiawei Han, and Jianfeng Gao. 2023. Model tells you what to discard: Adaptive kv cache compression for llms. *arXiv preprint arXiv:2310.01801*.
- Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. 2020. [Constructing A multi-hop QA dataset for comprehensive evaluation of reasoning steps](#). In *Proceedings of the 28th International Conference on Computational Linguistics, COLING 2020, Barcelona, Spain (Online), December 8-13, 2020*, pages 6609–6625. International Committee on Computational Linguistics.
- Connor Holmes, Masahiro Tanaka, Michael Wyatt, Ammar Ahmad Awan, Jeff Rasley, Samyam Rajbhandari, Reza Yazdani Aminabadi, Heyang Qin, Arash Bakhtiari, Lev Kurilenko, et al. 2024. DeepSpeed-fastgen: High-throughput text generation for llms via mii and deepspeed-inference. *arXiv preprint arXiv:2401.08671*.
- Greg Kamradt. 2023. Needle in a haystack - pressure testing llms. [https://github.com/gkamradt/LLMTest\\_NeedleInAHaystack](https://github.com/gkamradt/LLMTest_NeedleInAHaystack).
- Tomás Kociský, Jonathan Schwarz, Phil Blunsom, Chris Dyer, Karl Moritz Hermann, Gábor Melis, and Edward Grefenstette. 2018. [The narrativeqa reading comprehension challenge](#). *Trans. Assoc. Comput. Linguistics*, 6:317–328.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient

- memory management for large language model serving with pagedattention. In *Proceedings of the 29th Symposium on Operating Systems Principles*, pages 611–626.
- Yuhong Li, Yingbing Huang, Bowen Yang, Bharat Venkitesh, Acyr Locatelli, Hanchen Ye, Tianle Cai, Patrick Lewis, and Deming Chen. 2024. Snapkv: Llm knows what you are looking for before generation. *arXiv preprint arXiv:2404.14469*.
- Bin Lin, Tao Peng, Chen Zhang, Minmin Sun, Lanbo Li, Hanyu Zhao, Wencong Xiao, Qi Xu, Xiafei Qiu, Shen Li, et al. 2024. Infinite-llm: Efficient llm service for long context with distattention and distributed kvcache. *arXiv preprint arXiv:2401.02669*.
- Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2024. Lost in the middle: How language models use long contexts. *Transactions of the Association for Computational Linguistics*, 11:157–173.
- Amirkeivan Mohtashami and Martin Jaggi. 2023. Landmark attention: Random-access infinite context length for transformers. *arXiv preprint arXiv:2305.16300*.
- Piotr Nawrot, Adrian Łańcucki, Marcin Chochowski, David Tarjan, and Edoardo Ponti. 2024. Dynamic memory compression: Retrofitting LLMs for accelerated inference. In *Forty-first International Conference on Machine Learning*.
- Ramya Prabhu, Ajay Nayak, Jayashree Mohan, Ramachandran Ramjee, and Ashish Panwar. 2024. vattention: Dynamic memory management for serving llms without pagedattention. *arXiv preprint arXiv:2405.04437*.
- Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*.
- Konrad Staniszewski, Szymon Tworkowski, Yu Zhao, Sebastian Jaszczur, Henryk Michalewski, Lukasz Kuciński, and Piotr Miłoś. 2023. Structured packing in llm training improves long context utilization. *ArXiv*, abs/2312.17296.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Szymon Tworkowski, Konrad Staniszewski, Mikołaj Patek, Yuhuai Wu, Henryk Michalewski, and Piotr Miłoś. 2024. Focused transformer: Contrastive training for context scaling. *Advances in Neural Information Processing Systems*, 36.
- Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. 2024. Efficient streaming language models with attention sinks. In *The Twelfth International Conference on Learning Representations*.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W. Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*, pages 2369–2380. Association for Computational Linguistics.
- Xinrong Zhang, Yingfa Chen, Shengding Hu, Zihang Xu, Junhao Chen, Moo Hao, Xu Han, Zhen Thai, Shuo Wang, Zhiyuan Liu, and Maosong Sun. 2024a.  $\infty$ Bench: Extending long context evaluation beyond 100K tokens. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15262–15277, Bangkok, Thailand. Association for Computational Linguistics.
- Zhenyu Zhang, Ying Sheng, Tianyi Zhou, Tianlong Chen, Lianmin Zheng, Ruisi Cai, Zhao Song, Yuandong Tian, Christopher Ré, Clark Barrett, et al. 2024b. H2o: Heavy-hitter oracle for efficient generative inference of large language models. *Advances in Neural Information Processing Systems*, 36.
- Yu Zhao, Yuanbin Qu, Konrad Staniszewski, Szymon Tworkowski, Wei Liu, Piotr Miłoś, Yuxiang Wu, and Pasquale Minervini. 2024. Analysing the impact of sequence composition on language model pre-training. *arXiv preprint arXiv:2402.13991*.
- Ming Zhong, Da Yin, Tao Yu, Ahmad Zaidi, Mutethia Mutuma, Rahul Jha, Ahmed Hassan Awadallah, Asli Celikyilmaz, Yang Liu, Xipeng Qiu, and Dragomir R. Radev. 2021. Qmsum: A new benchmark for query-based multi-domain meeting summarization. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021*, pages 5905–5921. Association for Computational Linguistics.



(a) Layer-7 Head-10, high correlation between attention score and  $L_2$ -Norm.



(b) Layer-0 Head-0, low correlation between attention score and  $L_2$ -Norm.

Figure 8: Attention loss of ideal compression and  $L_2$  norm-based compression in Llama-2-7b-80k. The  $x$ -axis represents the compression ratio; the  $y$ -axis represents the attention loss (defined by Equation (1)) The results average over 1024 chunks on Wikipedia, with a length of 1024.

## A Attention score loss when compressing the KV cache

We discuss further the correlation between  $L_2$  norm and attention scores. We already displayed in Figure 1 the  $L_2$  norm and attention correlation across heads and layers using the original Llama2-7b and the long context Llama2-7b-32k and Llama2-7b-80k. We can see that patterns are quite consistent across all the models. To better visualise how correlation varies across different heads, in Figure 8, we only consider two heads from layer 10 and layer 0 and show the ALR from Equation (1). As expected, we see that in layer 0, the difference is larger due to a lower correlation.

## B More results on Language modelling task

In the following, we show results when performing compression only on layers that show a lower correlation between  $L_2$  norm and attention score. We show in Fig. 10 that for language modelling tasks, the different layer drop has little impact on final accuracy and perplexity. The difference becomes significant only when the KV cache is pruned to retain only one thousand pairs. All experiments are averaged over 50 chunks from English Wikipedia.

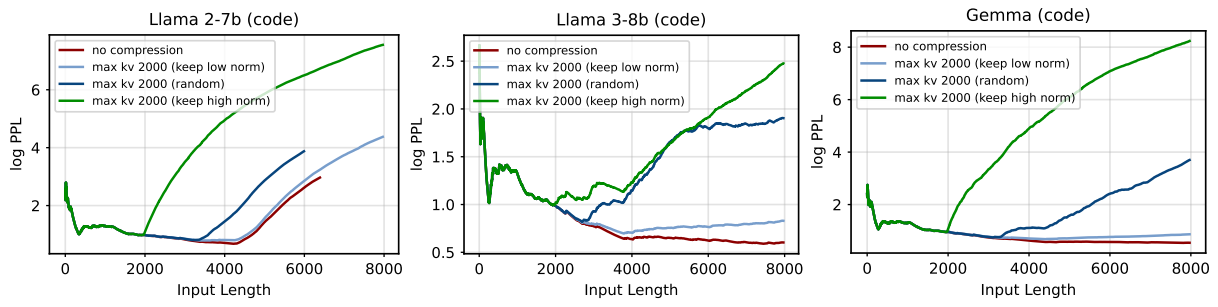
## C More Results on Long-Context Modelling Tasks

In addition to llama-2-7b-80k (Fu et al., 2024), we test the compression method using llama-2-7b-longlora-32k-ft (Chen et al., 2023) on the needle-in-a-haystack and passkey retrieval tasks. As shown in Fig. 11a, we can see that compressing 30% of KV cache only results in a slight performance degradation on the needle-in-a-haystack task. We also observe that the performance even increases slightly when we compress 10% of KV cache. In figure Fig. 11b, we observe that the llama-2-7b-longlora-32k-ft maintains 100% performance when compressing 80% of KV cache and only as a slight decrease when compressing 90% of KV cache. Furthermore, the model fails to generate correct answers if we compress KV pairs with low  $L_2$  norm and keep high  $L_2$  norm ones. The evaluation results of llama-2-7b-longlora-32k-ft are consistent with the llama-2-7b-80k, which further indicates the effectiveness of compressing KV cache using  $L_2$  norm.

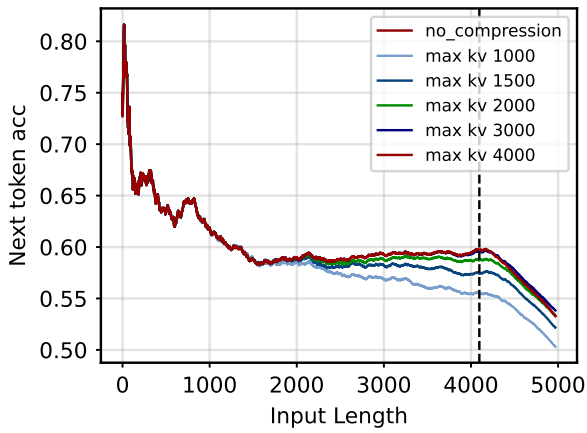
### C.1 Analysis of Skipped Layers

As shown in Fig. 1, we find heads in the first two layers and the middle layers have a relatively low correlation between attention scores and  $L_2$  norm. Thus, we conduct experiments to analyse the impact of skipping layers that have a low correlation for compression. As shown in Fig. 12a and Fig. 12c, we observe that only skipping the first layer (layer-0) decreases the performance on the needle-in-a-haystack task significantly. We can see that skipping the first two layers (layer-0,1) has a similar performance compared to skipping the first three layers (layer-0,1,2). Furthermore, as shown in Fig. 12b and Fig. 12d, only skipping the first layer can result in significant performance degradation. We also find that the compression

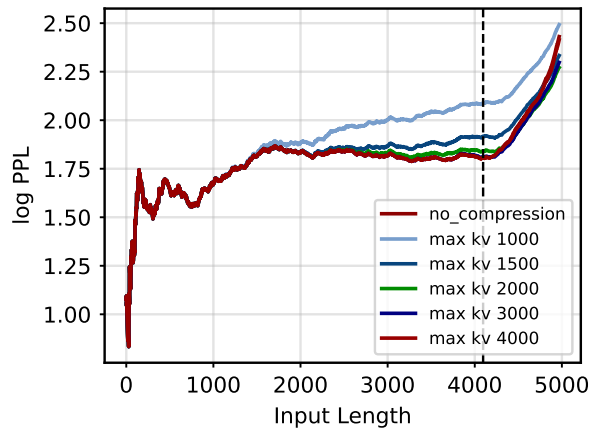




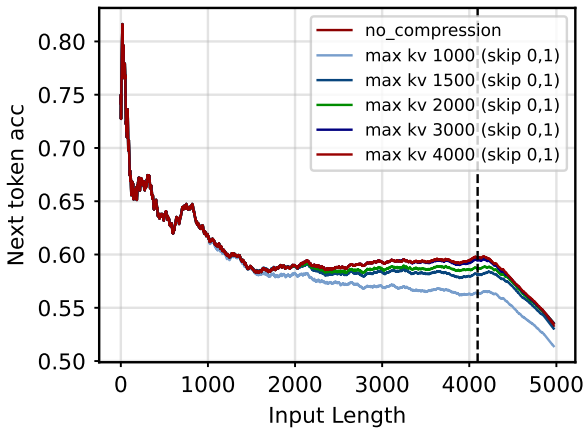
ratio is not proportional to the overall accuracy of models in the passkey retrieval task when we compress the first layer, where the accuracy shows a U-shape curve regarding the compression ratio.



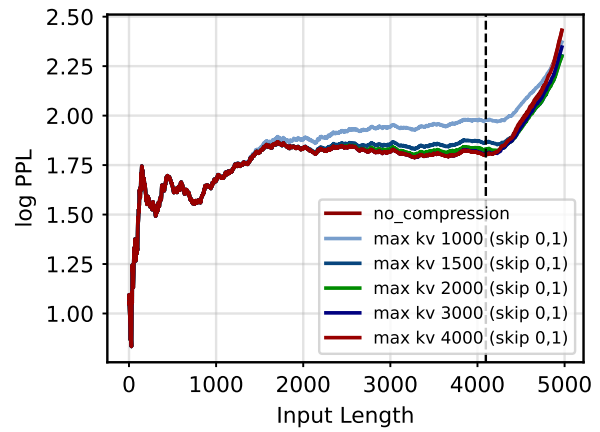
(a) Accuracy on language modelling task when not skipping any layers.



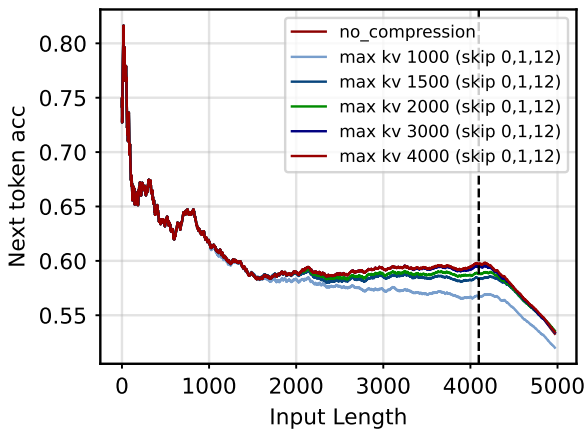
(b) Perplexity on language modelling task when not skipping any layers.



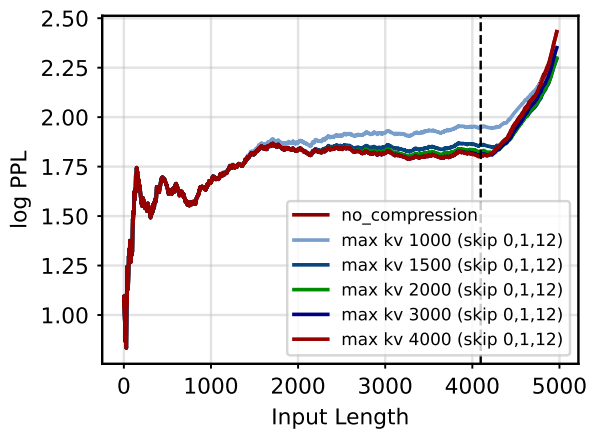
(c) Accuracy on language modelling task when skipping the first two layers.



(d) Perplexity on language modelling task when skipping the first two layers.

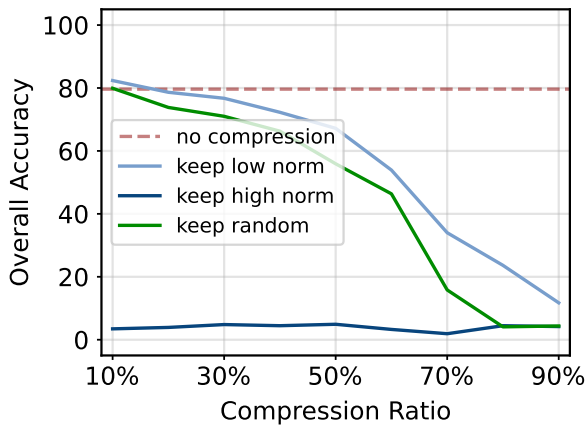


(e) Accuracy on language modelling task when skipping layers 0,1 and 12.

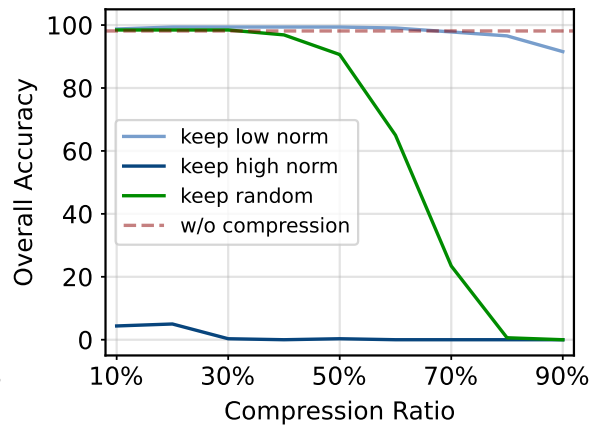


(f) Perplexity on language modelling task when skipping layers 0,1 and 12.

Figure 10: Skipping compression at different layers with Llama2-7b

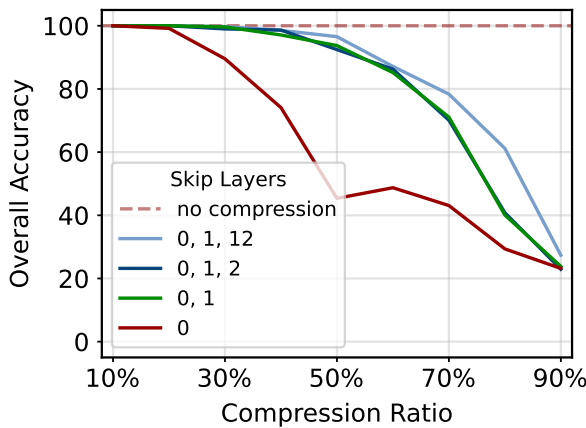


(a) Overall accuracy of Llama-2-7b-longlora-32k-ft on the needle-in-a-haystack task.

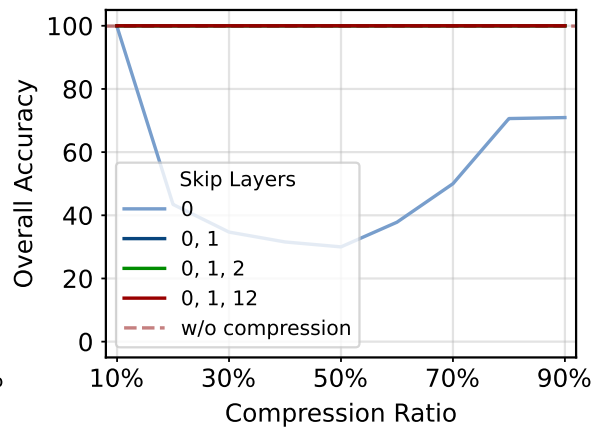


(b) Overall accuracy of Llama-2-7b-longlora-32k-ft on the passkey retrieval task.

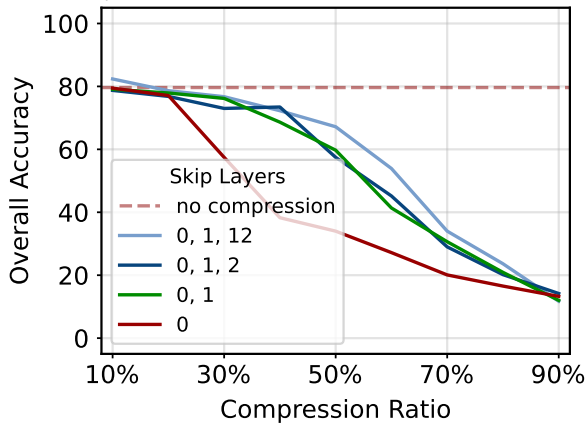
Figure 11: Evaluation results of Llama-2-7b-longlora-32k-ft on the needle-in-a-haystack and passkey retrieval tasks.



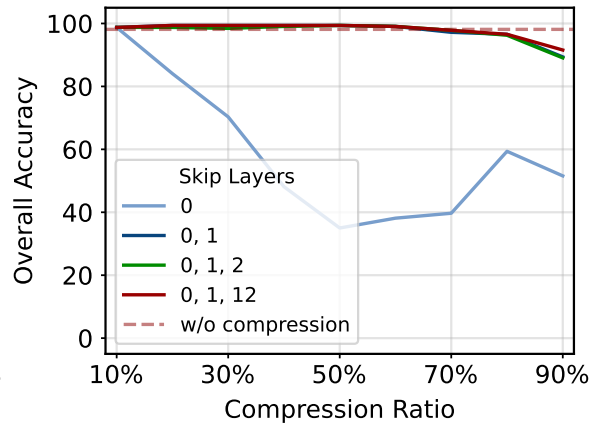
(a) Overall accuracy of Llama-2-7b-80k on the needle-in-a-haystack task.



(b) Overall accuracy of Llama-2-7b-80k on the passkey retrieval task.

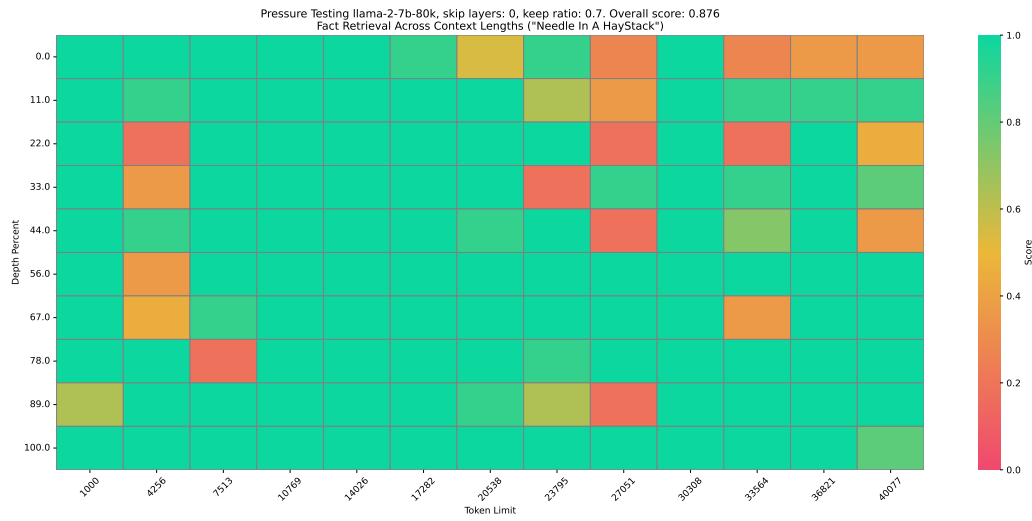


(c) Overall accuracy of Llama-2-7b-longlora-32k-ft on the needle-in-a-haystack task.

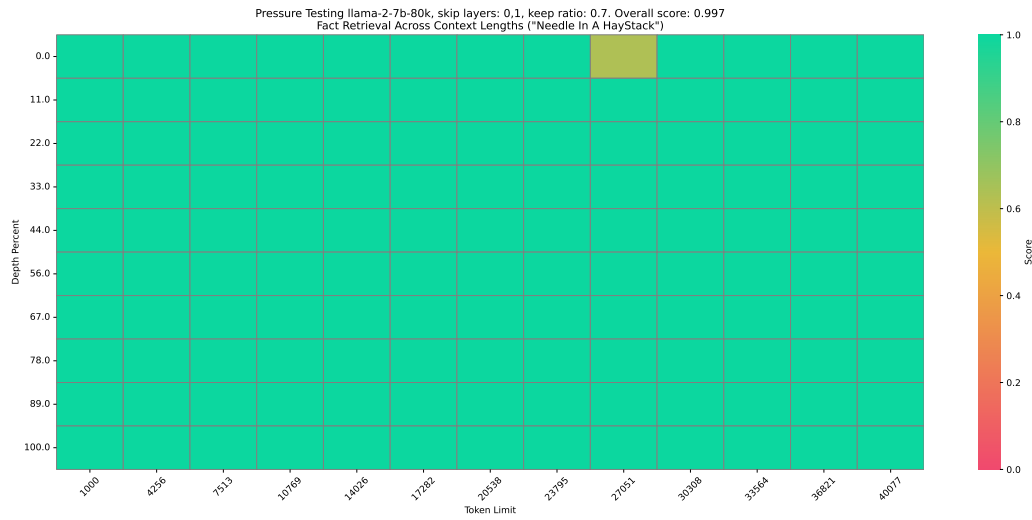


(d) Overall accuracy of Llama-2-7b-longlora-32k-ft on the passkey retrieval task.

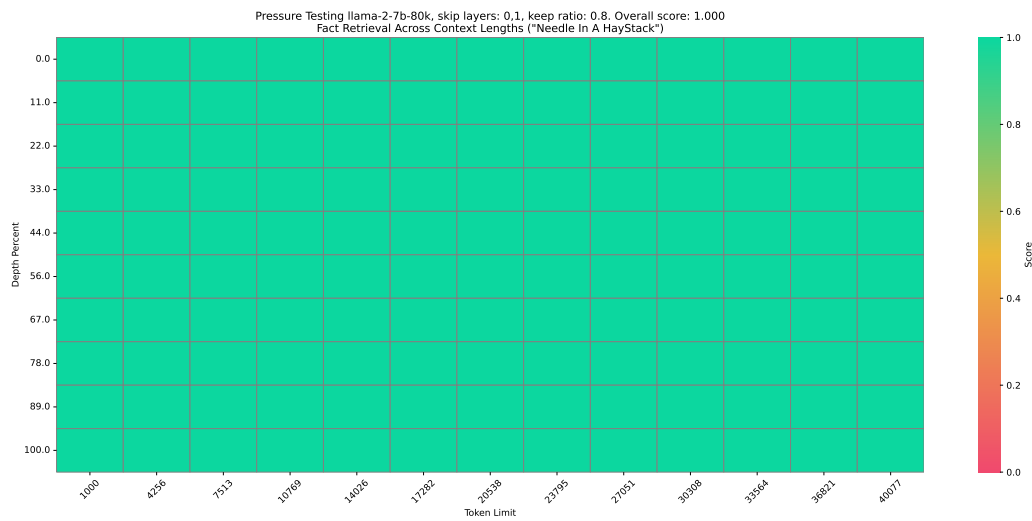
Figure 12: Analysing of skipping different layers for compression.



(a) Llama-2-7b-80k, skip layer-0, compression ratio 30%



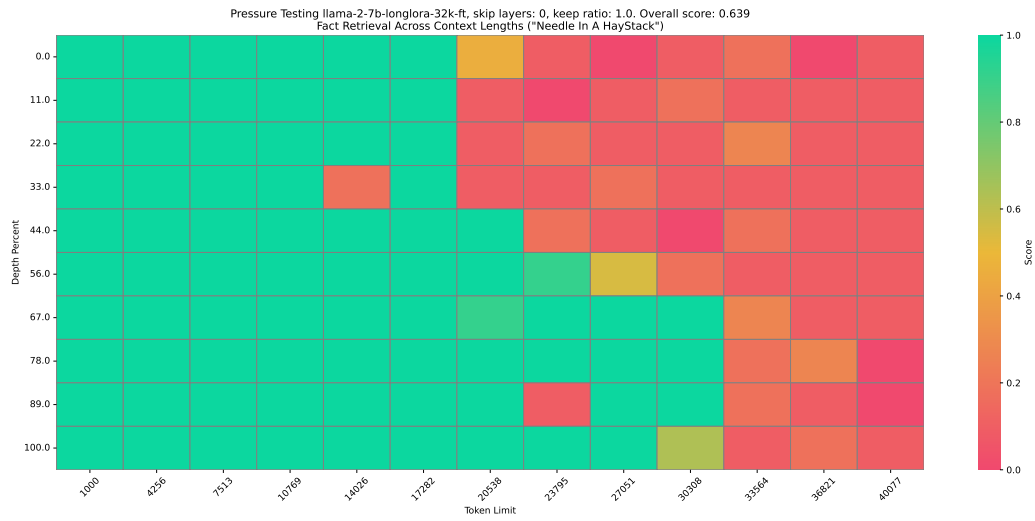
(b) Llama-2-7b-80k, skip layer-0 and layer-1, compression ratio 30%



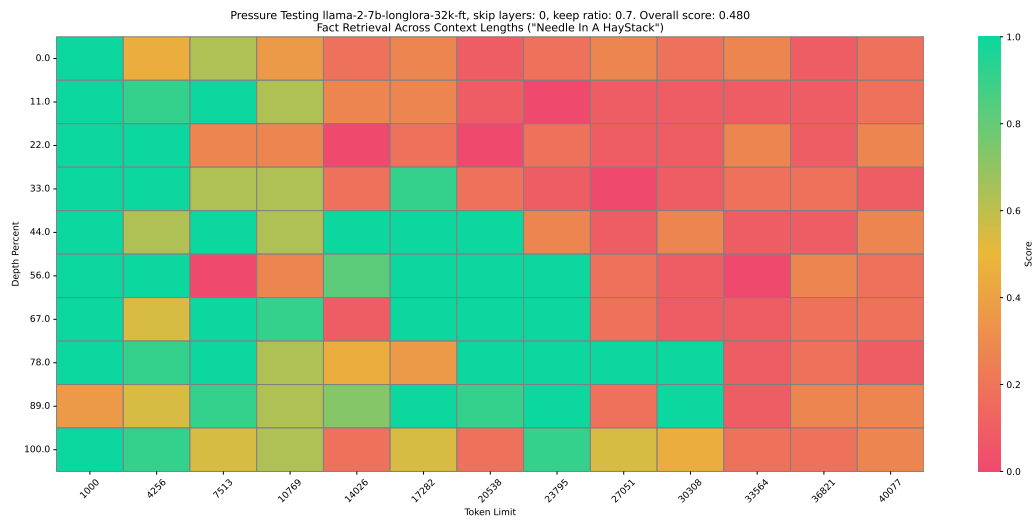
(c) Llama-2-7b-80k, skip layer-0 and layer-1, compression ratio 20%

Figure 13: Detailed results of Llama-2-7b-80k on the needle-in-a-haystack task.

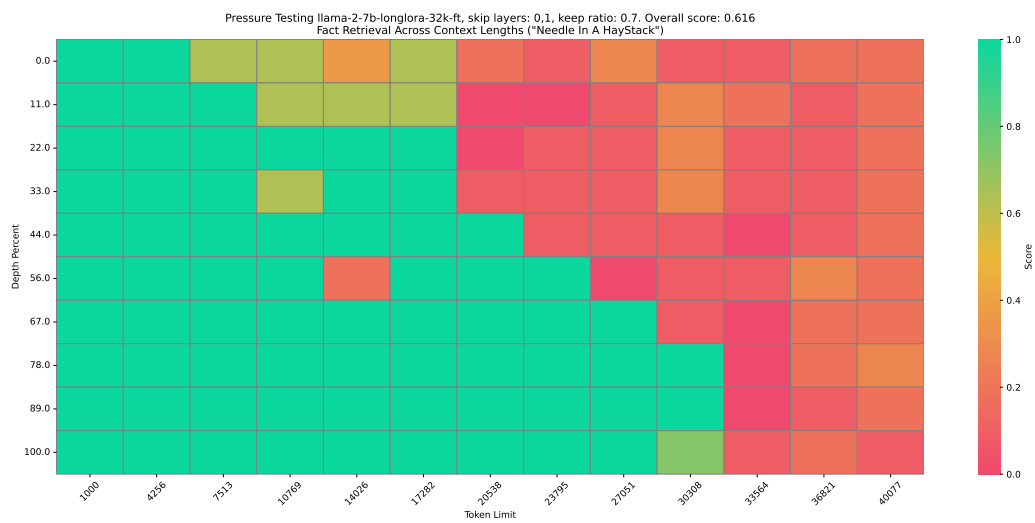




(a) Llama-2-7b-longlora-32k-ft, without compression



(b) Llama-2-7b-longlora-32k-ft, skip layer-0, compression ratio 30%



(c) Llama-2-7b-longlora-32k-ft, skip layer-0 and layer-1, compression ratio 30%

Figure 14: Detailed results of Llama-2-7b-longlora-32k-ft on the needle-in-a-haystack task.

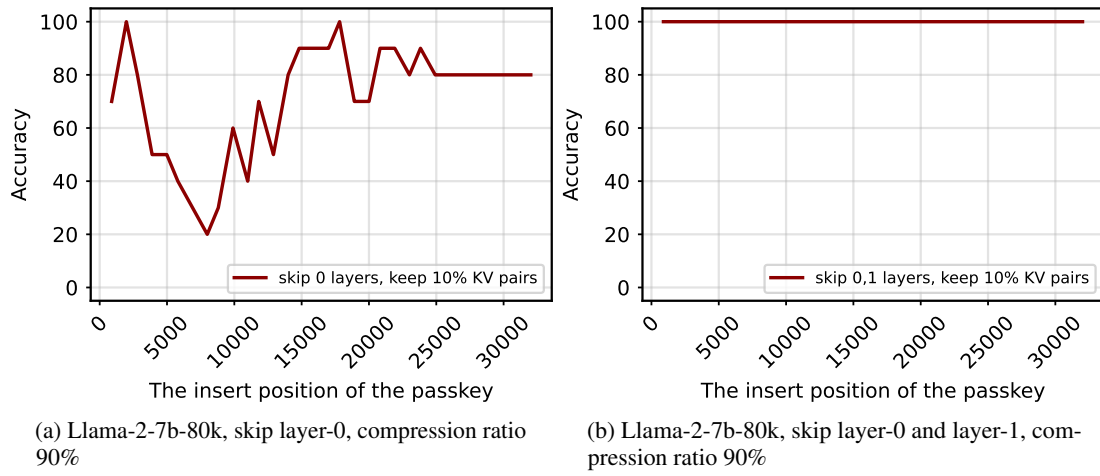


Figure 15: Accuracy on the passkey retrieval. The  $x$ -axis presents the position of the passkey, and the  $y$ -axis presents the accuracy.

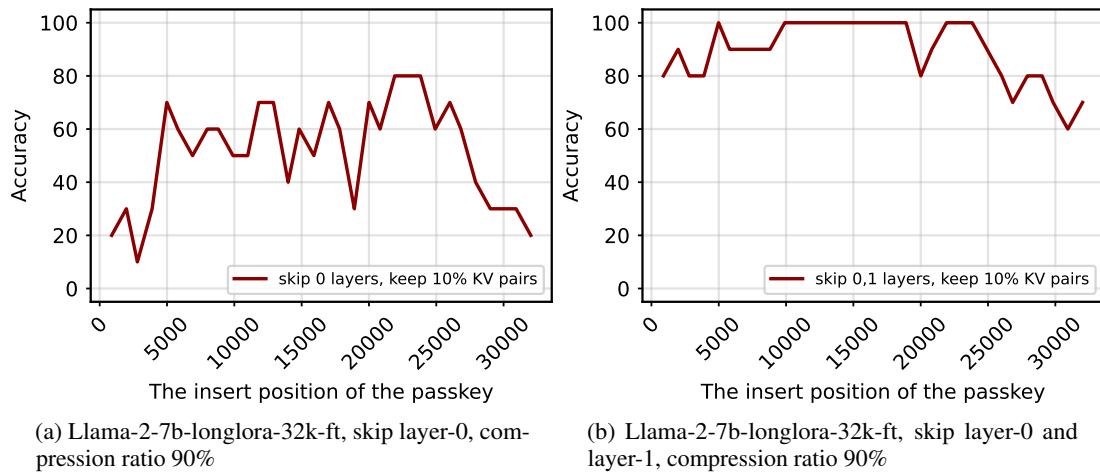


Figure 16: Accuracy on the passkey retrieval. The  $x$ -axis presents the position of the passkey, and the  $y$ -axis presents the accuracy.

## **C.2 Longbench Evaluation**

In this section we show detailed results from the LongBench dataset ([Zhang et al., 2024a](#)). In Figure 17 we show results for Llama2-80k, while in Figure 18 we show results for the long context model Llama3.1-8b.

## **D More Visualizations**

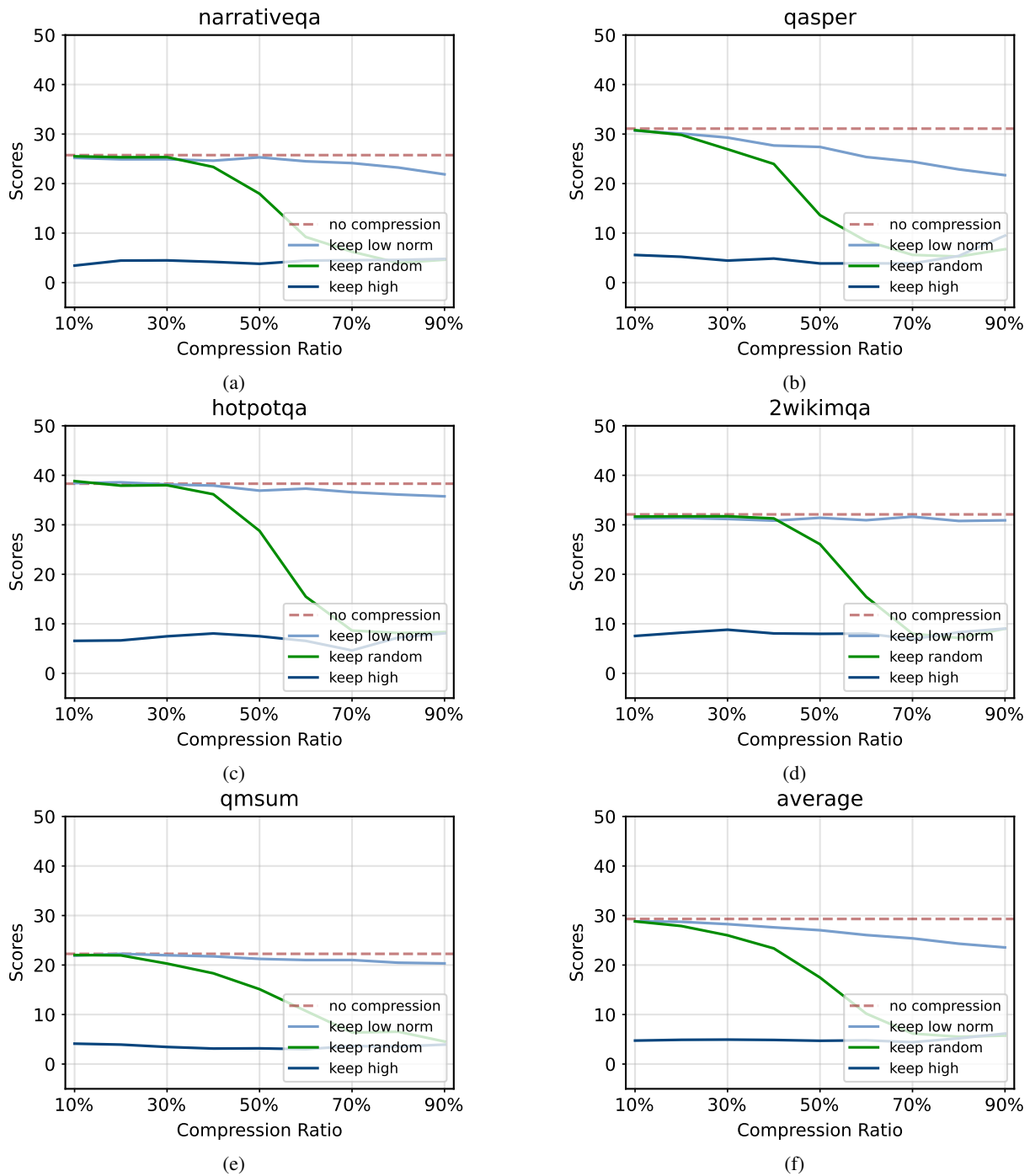


Figure 17: Evaluation results of Llama-2-7b-80k on long context tasks from Longbench, including narrativeqa and qasper, hotpotqa, 2wikimqa, and qmsum.



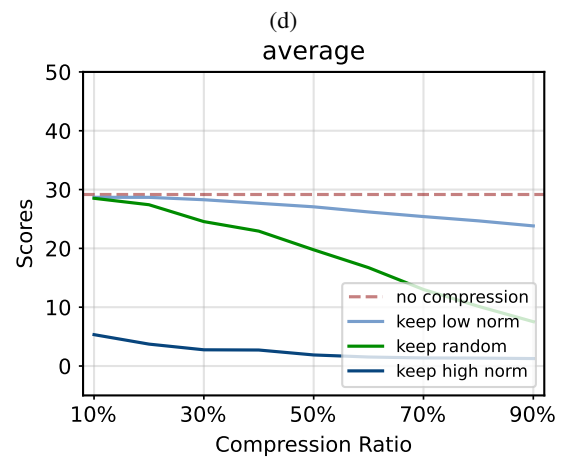
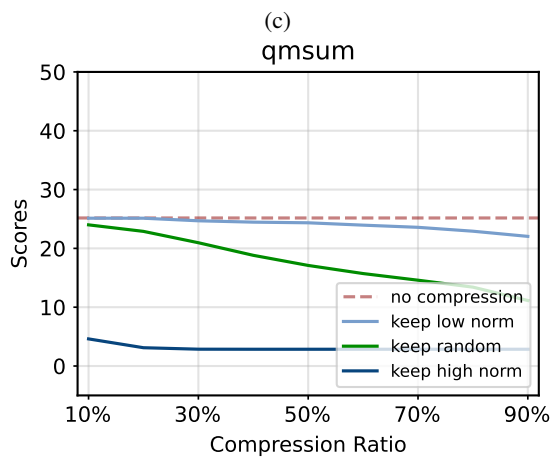
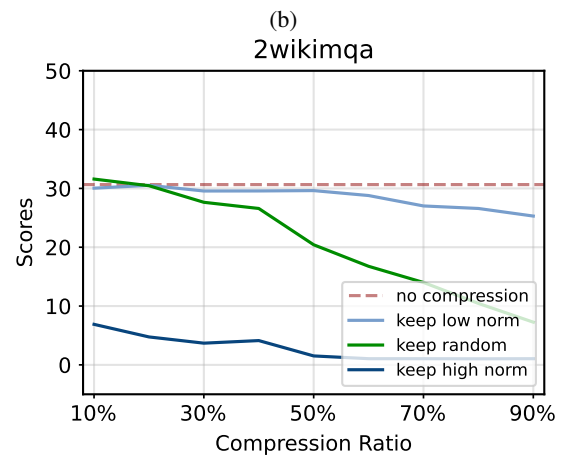
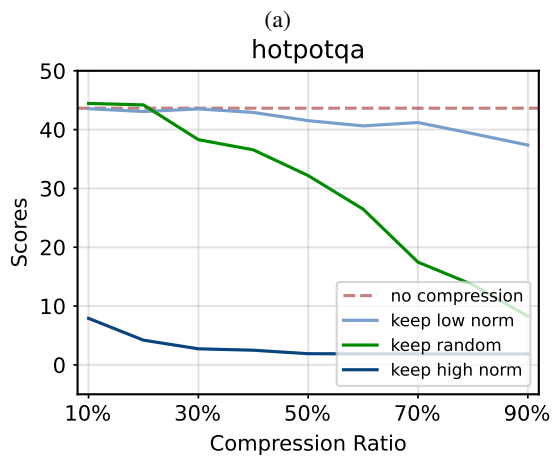
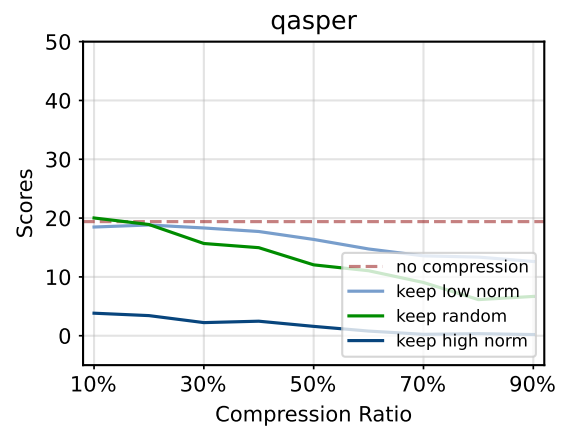
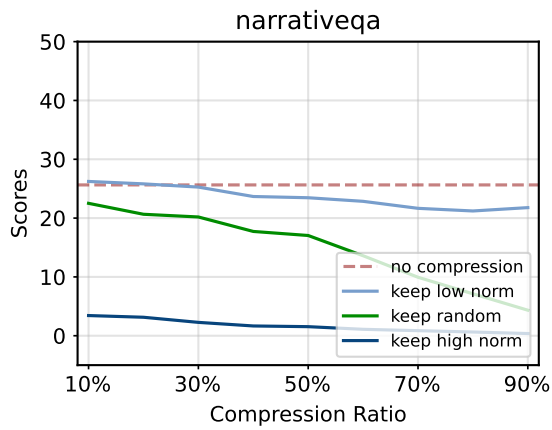


Figure 18: Evaluation results of Llama-3.1-8B on long context tasks from Longbench, including narrativeqa and qasper, hotpotqa, 2wikimqa, and qmsum.



Figure 19: Attention maps in Llama2-7B

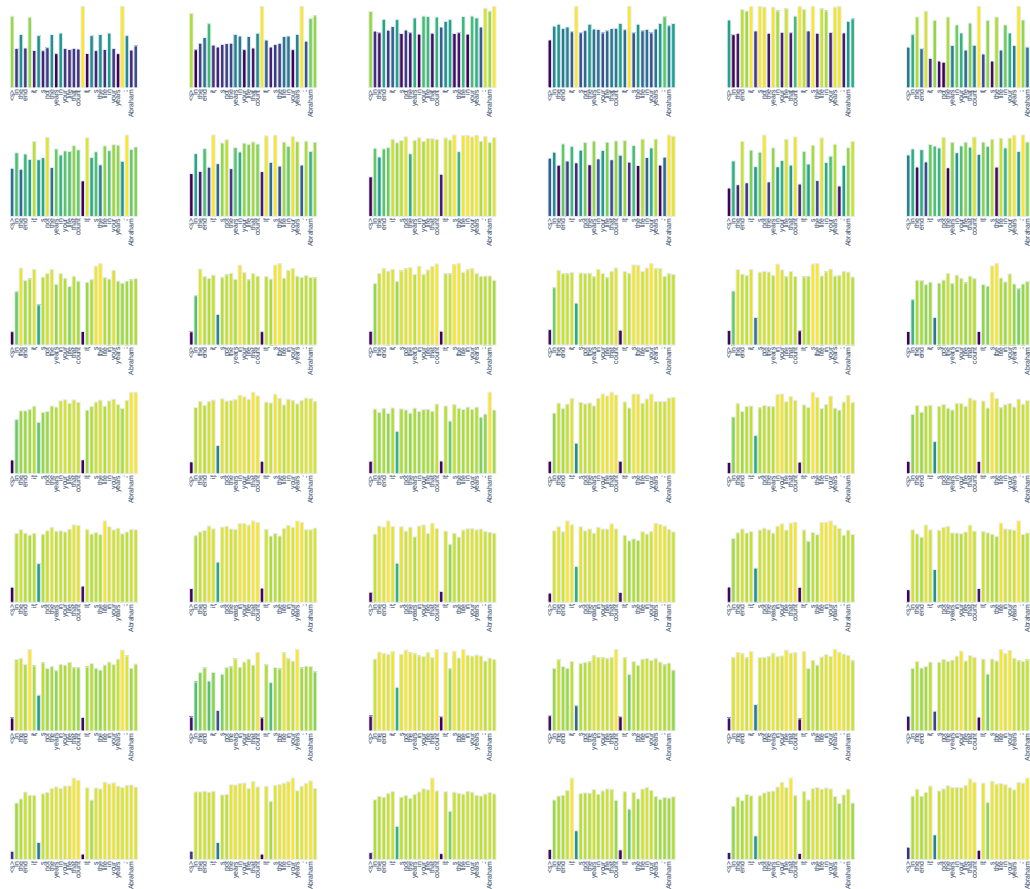


Figure 20: Norms of KV cache tokens in Llama2-7B





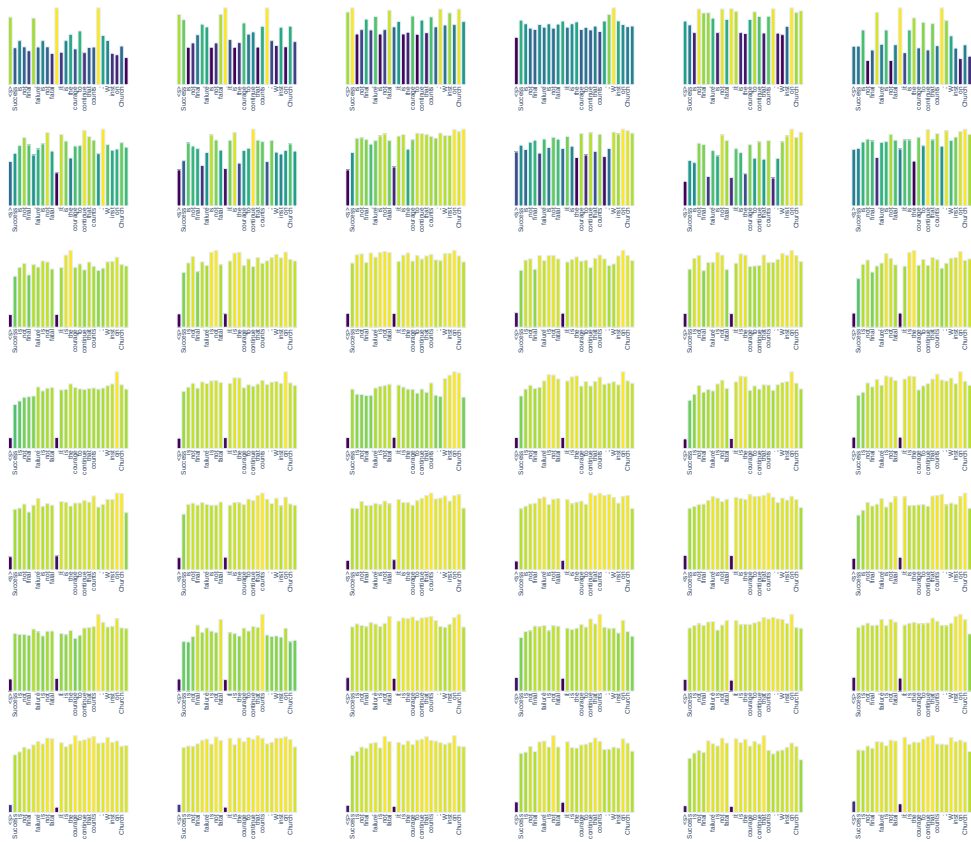


Figure 22: Norms of KV cache tokens in Llama2-7B



Figure 23: Attention maps in Llama2-7B



Figure 24: Norms of KV cache tokens in Llama2-7B

### E Additional token embeddings plots

We show in Figure 25 some additional figure that represent Llama3-8b token embeddings sparsity.

### F Experimental setup

In all experiments, we used the HuggingFace library and did not change the model’s default hyperparameters. For language modelling, results are averaged across 50 samples. The Figure 8 and Figure 1 are the average results of 1024 examples with a chunk size of 1024 using Wikipedia.

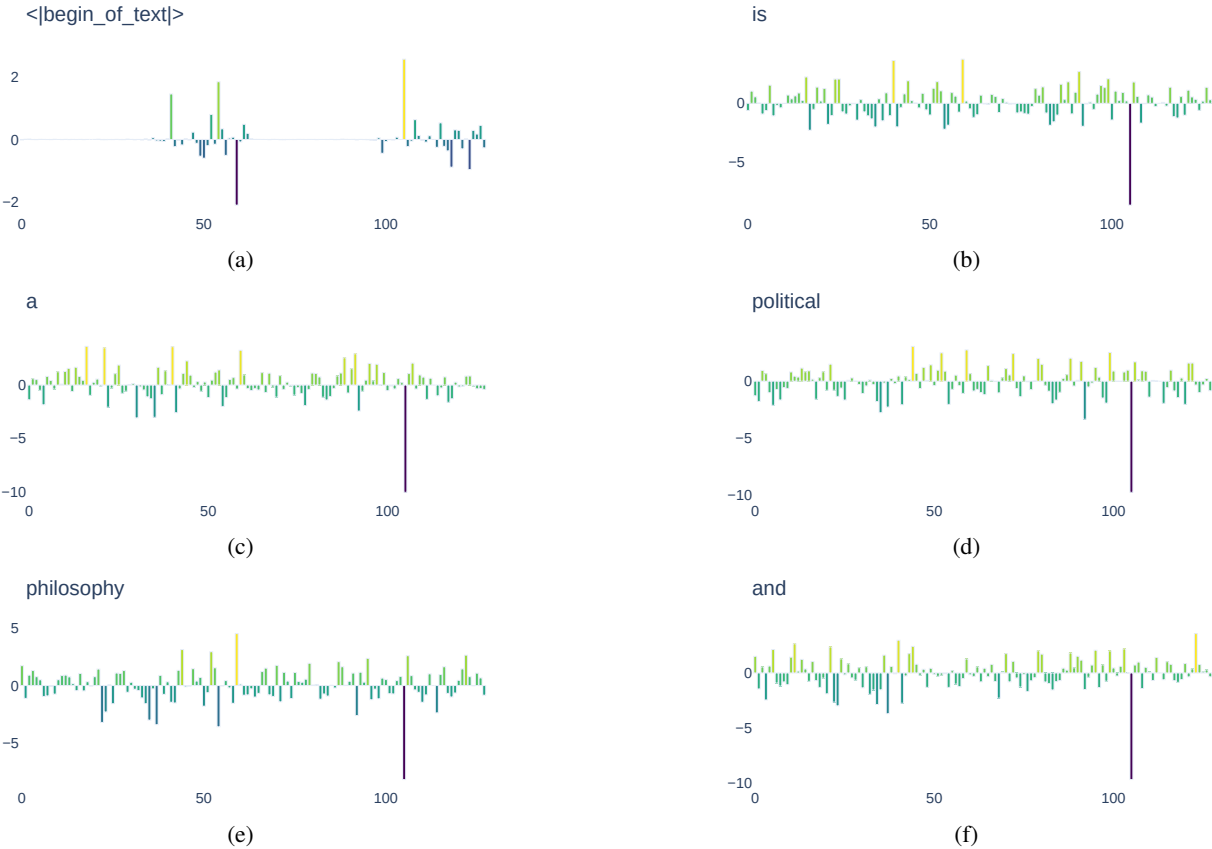


Figure 25: Key projections of Llama3-8b of the bos *|beginoftext|* token vs other tokens. Each value represents the activation in a specific dimension for the embedding of the key projection. We found similar patterns across almost all heads and layers and in multiple texts.