

# Extending Context Window of Large Language Models via Positional Interpolation

---

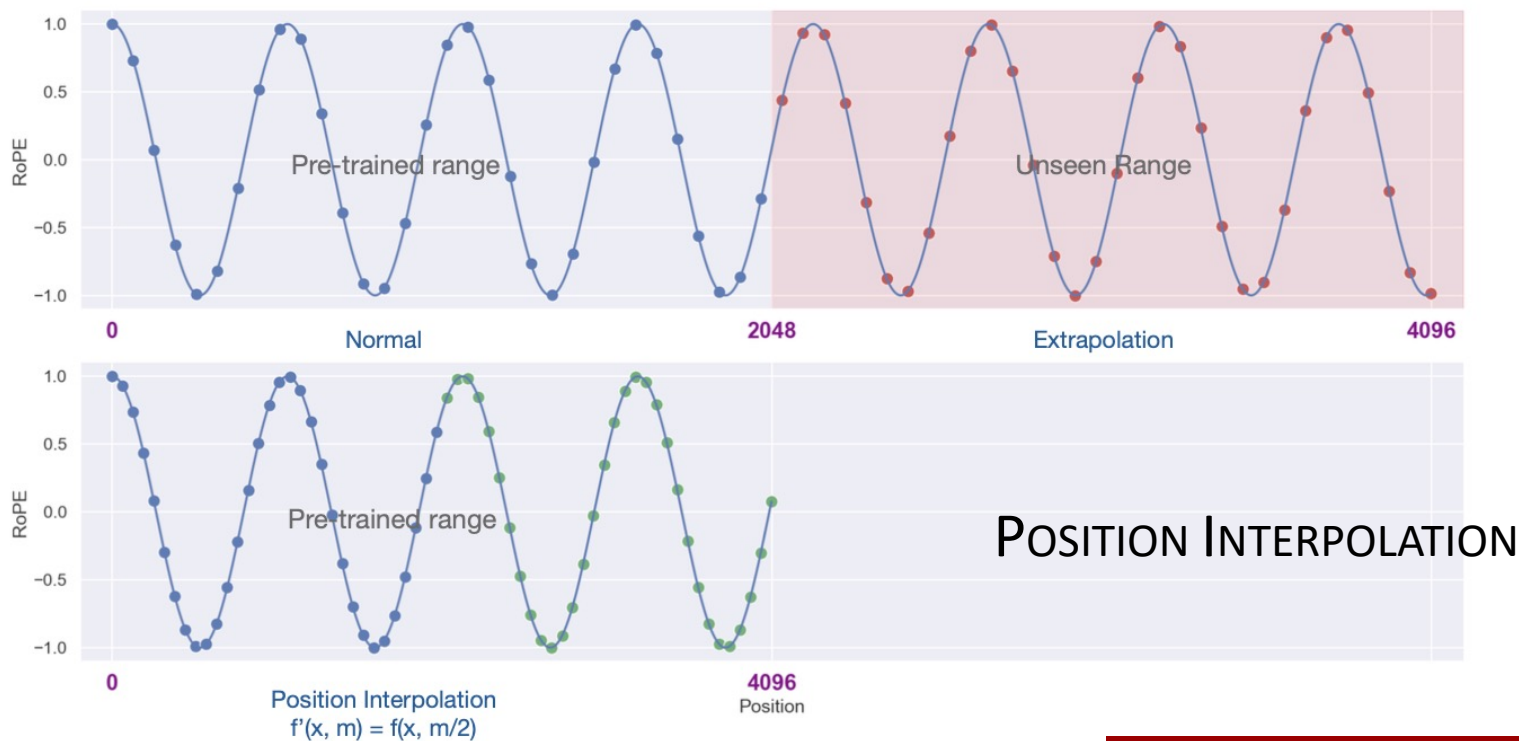
**Shouyuan Chen   Sherman Wong   Liangjian Chen   Yuandong Tian**  
Meta Platforms Inc.  
`{chenshouyuan, shermanwong, clj, yuandong}@meta.com`

author={Chen, Shouyuan and Wong, Sherman and Chen, Liangjian and Tian, Yuandong},  
journal={arXiv preprint arXiv:2306.15595},  
year={2023.06.28}  
Cites={113}

# Method

Size	Model		Fine-tuning steps					
	Context Window	Method	200	400	600	800	1000	10000
7B	8192	FT	1792	2048	2048	2048	2304	2560
33B	8192	FT	1792	2048	1792	2048	2304	-
7B	8192	PI	8192	8192	8192	8192	8192	-
7B	16384	PI	16384	16384	16384	16384	16384	-
7B	32768	PI	32768	32768	18432	32768	32768	-
33B	8192	PI	8192	8192	8192	8192	8192	-
33B	16384	PI	16384	16384	16384	16384	16384	-

Table 4: Effective context window sizes after fine-tuning. FT: Direct fine-tuning. PI: Position Interpolation.



## 位置编码只依赖于相对位置

Transformer models require explicit positional information to be injected, typically in the form of positional encodings, to represent the order of inputs. We consider **Rotary Position Embedding (RoPE)** (Su et al., 2021), which is the position encoding used in the LLaMA model (Touvron et al., 2023). Given a position index  $m \in [0, c]$  and an embedding vector  $\mathbf{x} := [x_0, x_1, \dots, x_{d-1}]^\top$ , where  $d$  is the dimension of the attention head, RoPE defines a vector-valued complex function  $\mathbf{f}(\mathbf{x}, m)$  as follows

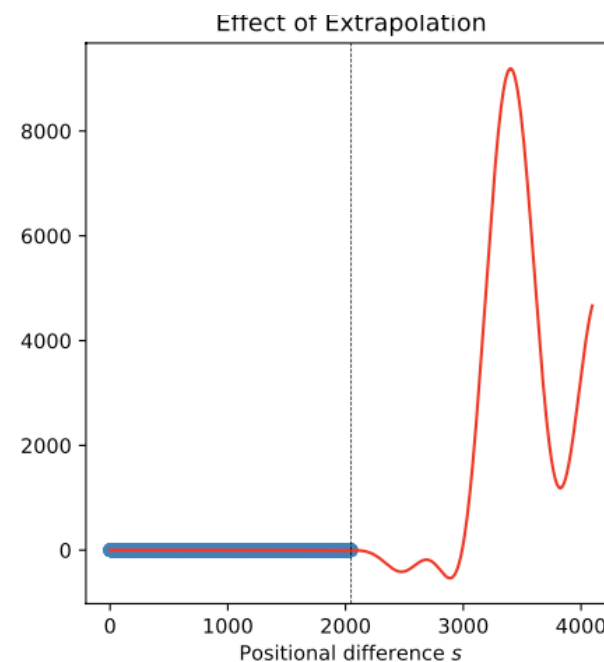
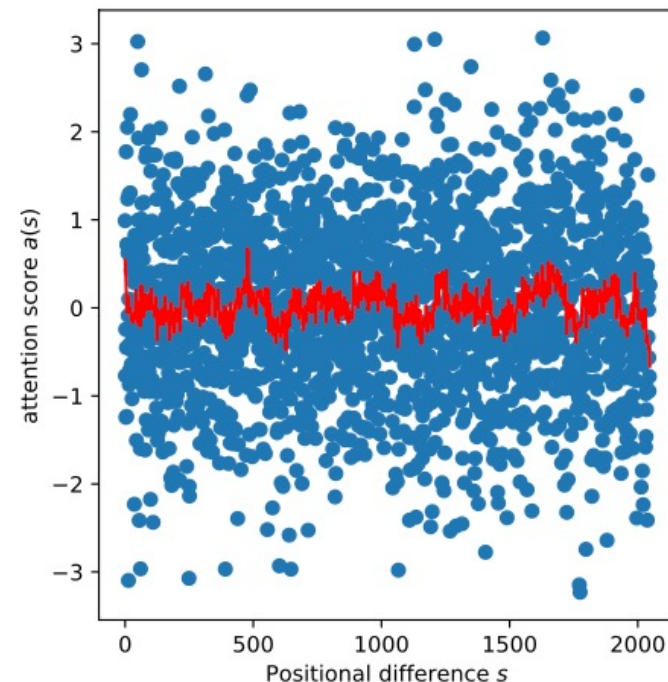
$$\mathbf{f}(\mathbf{x}, m) = [(x_0 + ix_1)e^{im\theta_0}, (x_2 + ix_3)e^{im\theta_1}, \dots, (x_{d-2} + ix_{d-1})e^{im\theta_{d/2-1}}]^\top \quad (1)$$

where  $i := \sqrt{-1}$  is the imaginary unit and  $\theta_j = 10000^{-2j/d}$ . Using RoPE, the self-attention score

$$\begin{aligned} a(m, n) &= \text{Re}\langle \mathbf{f}(\mathbf{q}, m), \mathbf{f}(\mathbf{k}, n) \rangle \\ &= \text{Re} \left[ \sum_{j=0}^{d/2-1} (q_{2j} + iq_{2j+1})(k_{2j} - ik_{2j+1})e^{i(m-n)\theta_j} \right] \\ &= \sum_{j=0}^{d/2-1} (q_{2j}k_{2j} + q_{2j+1}k_{2j+1}) \cos((m-n)\theta_j) + (q_{2j}k_{2j+1} - q_{2j+1}k_{2j}) \sin((m-n)\theta_j) \\ &=: a(m-n) \end{aligned} \quad (2)$$

is **only dependent on relative position  $m-n$**  through trigonometric functions. Here  $\mathbf{q}$  and  $\mathbf{k}$  are the query and key vector for a specific attention head. At each layer, RoPE is applied on both query and key embeddings for computing attention scores.

$$a(s) = \text{Re} \left[ \sum_{j=0}^{d/2-1} h_j e^{is\theta_j} \right] \quad h_j := (q_{2j} + iq_{2j+1})(k_{2j} - ik_{2j+1})$$



## Method

证明Positional Interpolation的上界比Extrapolation小很多

In Fig. 2, thanks to the smoothness of bases functions  $\phi_j$  interpolation is much more stable and will not lead to wild values. Therefore, instead of extrapolate the attention score in Eqn. 3 to  $s > L$ , how about we define an attention score  $\tilde{a}(s) = a(Ls/L')$  where  $L'$  is the longer context window? Formally, we replace RoPE  $\mathbf{f}$  by  $\mathbf{f}'$  defined as follows

$$\mathbf{f}'(\mathbf{x}, m) = \mathbf{f}\left(\mathbf{x}, \frac{mL}{L'}\right). \quad (4)$$

Positional Interpolation的上界

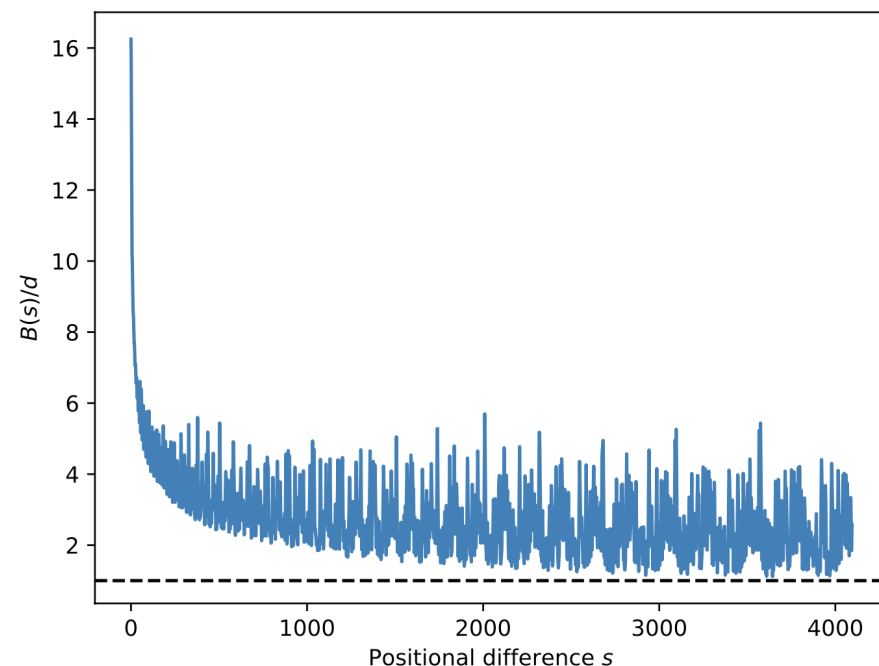
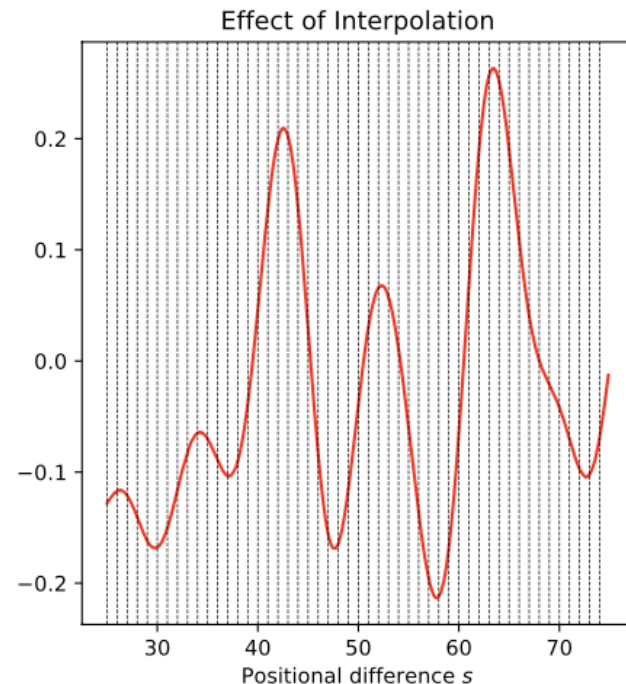
$$a(s) = \text{Re} \left[ \sum_{j=0}^{d/2-1} h_j e^{is\theta_j} \right] \quad s \in [s_1, s_2]$$

$$|a(s) - a_{\text{linear}}(s)| \leq d \left( \max_j |h_j| \right) \frac{(s - s_1)(s_2 - s)}{8 \ln c}$$

$$|a(s) - a_{\text{linear}}(s)| \leq \frac{d}{32 \ln c} \max_j |h_j| \approx \frac{d \max_j |h_j|}{294.73}$$

Extrapolation的上界

$$|a(s)| \leq \left( \max_j |h_j - h_{j+1}| \right) \sum_{k=0}^{d/2-1} |A_{k+1}(s)| \leq 2 \left( \max_j |h_j| \right) \sum_{k=0}^{d/2-1} |A_{k+1}(s)|,$$



# Method

Size	Model		Evaluation Context Window Size				
	Context Window	Method	2048	4096	8192	16384	32768
7B	2048	None	7.20	$> 10^3$	$> 10^3$	$> 10^3$	$> 10^3$
7B	8192	FT	7.21	7.34	7.69	-	-
7B	8192	PI	7.13	6.96	6.95	-	-
7B	16384	PI	7.11	6.93	6.82	6.83	-
7B	32768	PI	7.23	7.04	6.91	6.80	6.77
13B	2048	None	6.59	-	-	-	-
13B	8192	FT	6.56	6.57	6.69	-	-
13B	8192	PI	6.55	6.42	6.42	-	-
13B	16384	PI	6.56	6.42	6.31	6.32	-
13B	32768	PI	6.54	6.40	6.28	6.18	6.09
33B	2048	None	5.82	-	-	-	-
33B	8192	FT	5.88	5.99	6.21	-	-
33B	8192	PI	5.82	5.69	5.71	-	-
33B	16384	PI	5.87	5.74	5.67	5.68	-
65B	2048	None	5.49	-	-	-	-
65B	8192	PI	5.42	5.32	5.37	-	-

Size	Model		Evaluation Context Window Size				
	Context Window	Method	2048	4096	8192	16384	32768
7B	2048	None	2.77	-	-	-	-
7B	8192	FT	2.85	2.74	2.73	-	-
7B	8192	PI	2.79	2.57	2.39	-	-
7B	16384	PI	2.79	2.57	2.37	2.25	-
7B	32768	PI	2.82	2.59	2.39	2.24	2.48
13B	2048	None	2.66	-	-	-	-
13B	8192	FT	2.71	2.56	2.50	-	-
13B	8192	PI	2.67	2.47	2.30	-	-
13B	16384	PI	2.68	2.47	2.29	2.18	-
13B	32768	PI	2.68	2.46	2.28	2.15	2.35
33B	2048	None	2.49	-	-	-	-
33B	8192	FT	2.56	2.48	2.47	-	-
33B	8192	PI	2.50	2.32	2.18	-	-
33B	16384	PI	2.53	2.34	2.18	2.07	-
65B	2048	None	2.42	-	-	-	-
65B	8192	PI	2.43	2.26	2.12	-	-

# Longlora: Efficient fine-tuning of long-context large language models

---

**Yukang Chen**<sup>1</sup>   **Shengju Qian**<sup>1</sup>   **Haotian Tang**<sup>2</sup>   **Xin Lai**<sup>1</sup>  
**Zhijian Liu**<sup>2</sup>   **Song Han**<sup>2,3</sup>   **Jiaya Jia**<sup>1</sup>  
<sup>1</sup>CUHK   <sup>2</sup>MIT   <sup>3</sup>NVIDIA

@article{  
Conference={ICLR 2024}  
Cites={57}  
**Oral**  
**8866**



# Method

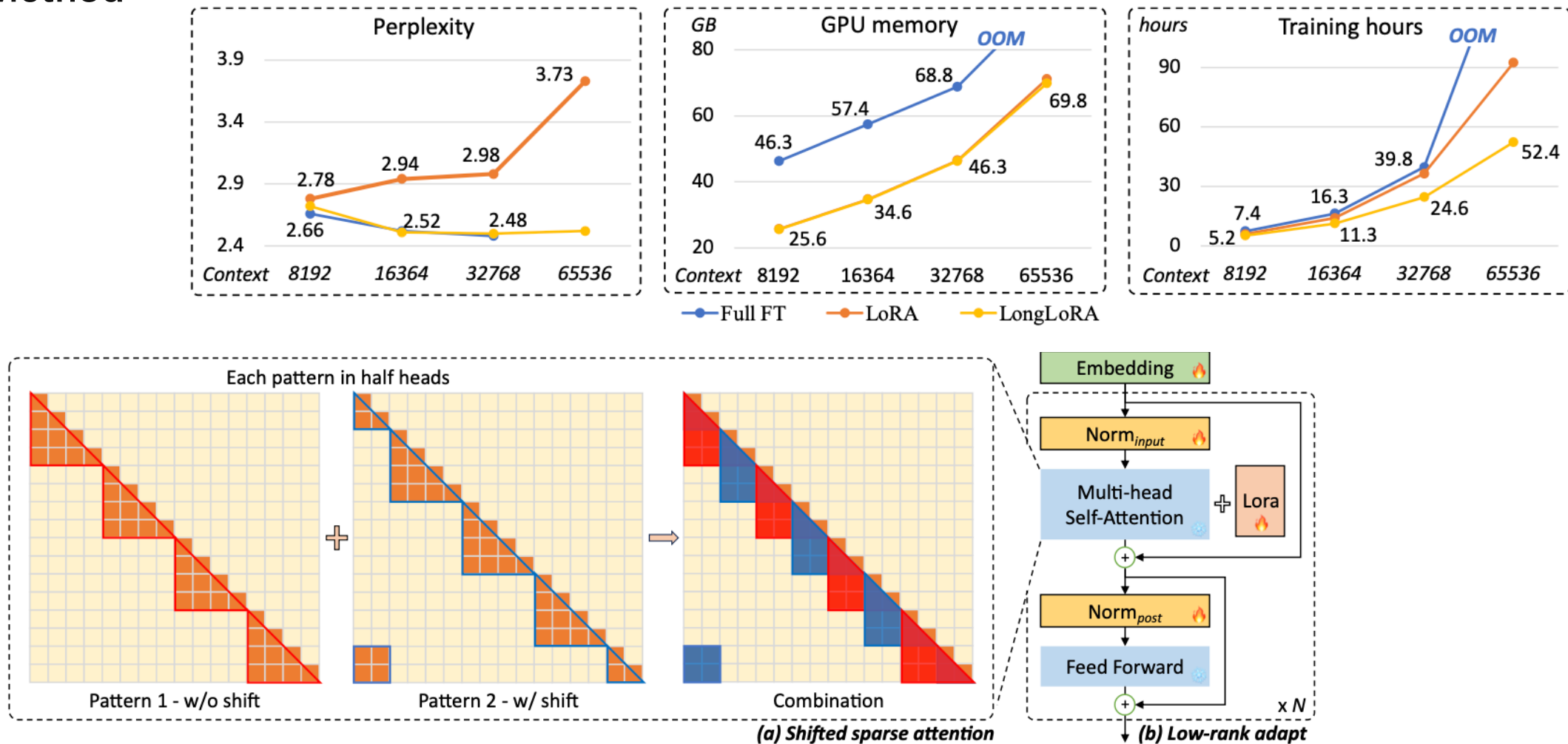
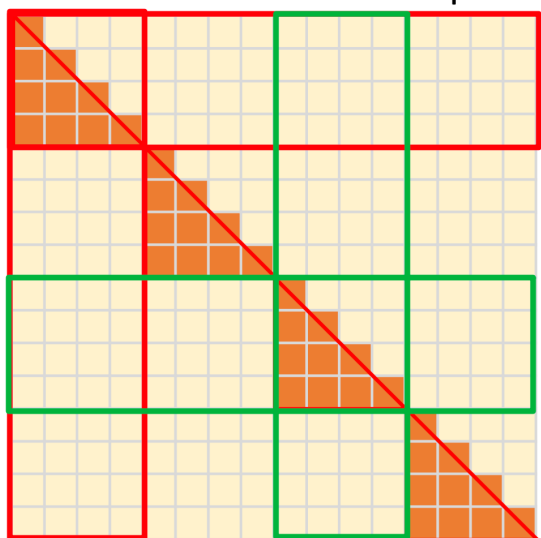
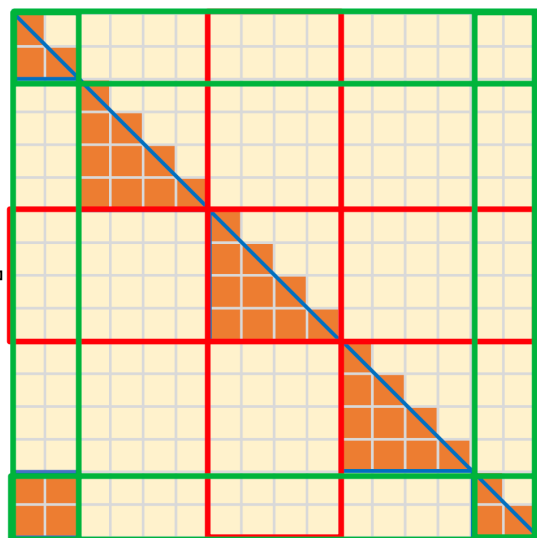


Figure 2: **Overview of LongLoRA.** We introduce Shifted Sparse Attention ( $S^2$ -Attn) during fine-tuning. The trained model retains original standard self-attention at inference time. In addition to training LoRA weights in linear layers, LongLoRA further makes embedding and normalization layers trainable. This extension is pivotal for context extension, and only introduces a minimal number of additional trainable parameters.

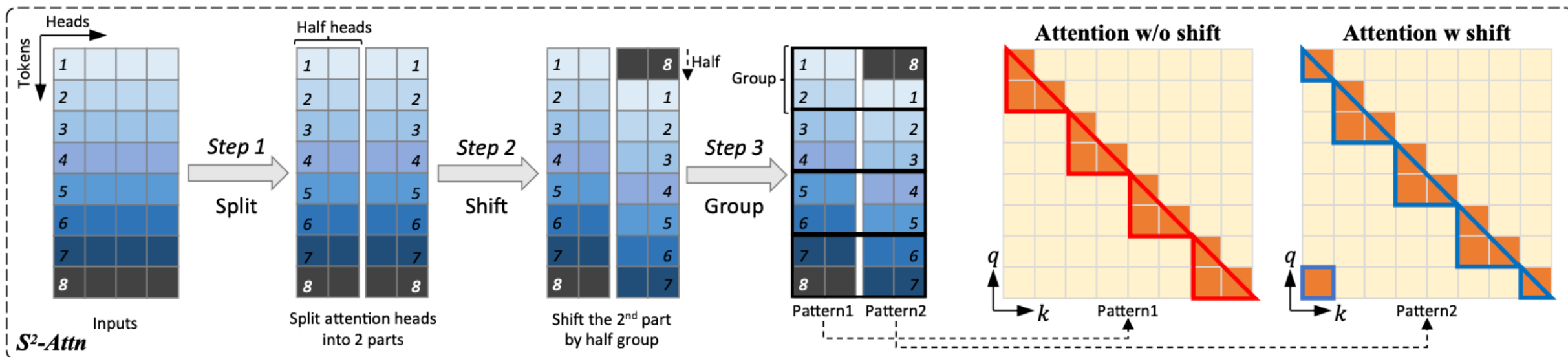
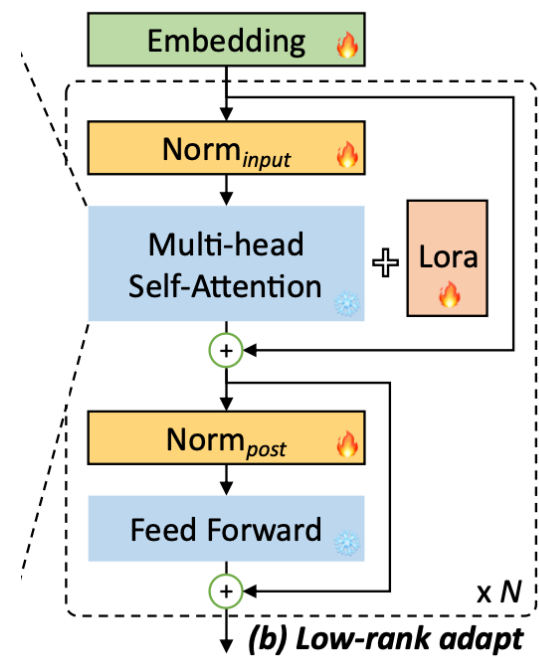
# Method



Pattern 1 - w/o shift



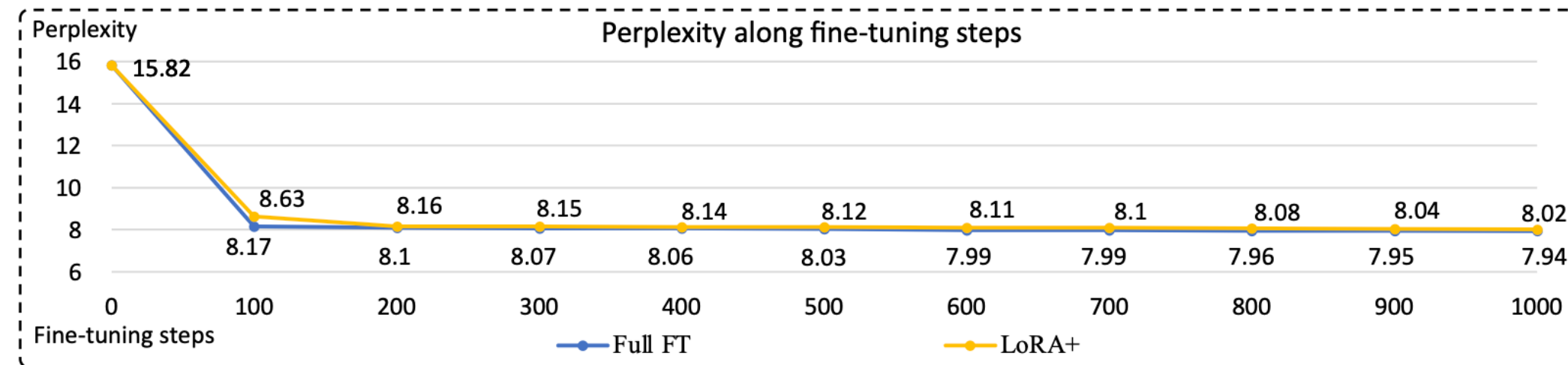
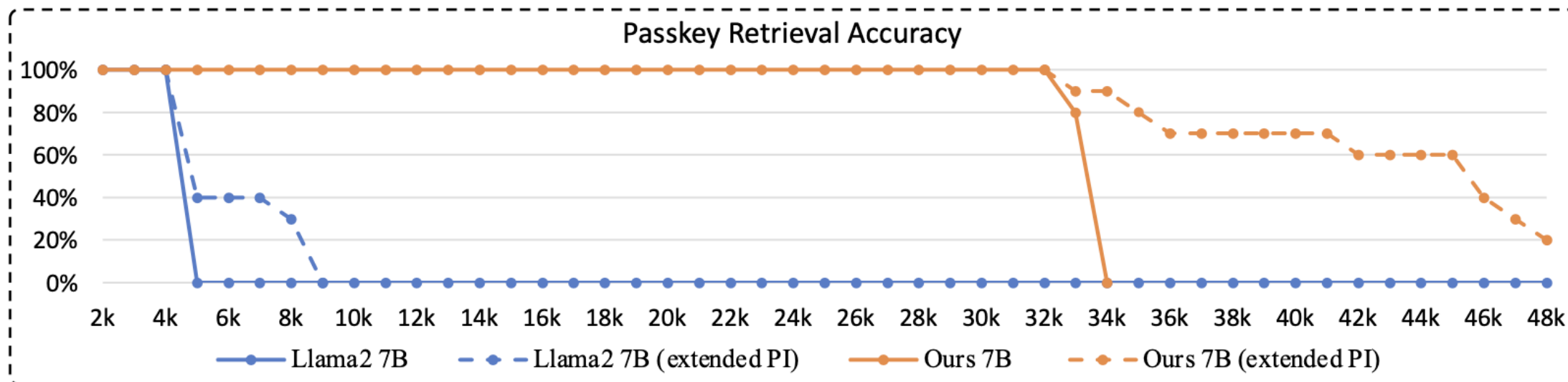
Pattern 2 - w/ shift





# EXPFRIMENT

Method	Full FT	LoRA (rank)						LoRA (rank = 8)		
		8	16	32	64	128	256	+ Norm	+ Embed	+ Norm & Embed
PPL	8.08	11.44	11.82	11.92	11.96	11.97	11.98	10.49	8.29	8.12



EXPERIMENT

在 8 \* A100上LoRA

Size	Training Context Length	LongLoRA		Evaluation Context Length				
		S <sup>2</sup> -Attn	LoRA <sup>+</sup>	2048	4096	8192	16384	32768
7B	8192	✓		3.14	2.85	2.66	-	-
		✓	✓	3.15	2.86	2.68	-	-
				3.20	2.91	2.72	-	-
	16384	✓		3.17	2.87	2.68	2.55	-
		✓	✓	3.17	2.87	2.66	2.51	-
	32768	✓		3.20	2.90	2.69	2.54	2.49
13B	8192	✓		2.96	2.69	2.53	-	-
		✓	✓	3.01	2.74	2.57	-	-
				3.04	2.77	2.60	-	-
	16384	✓		2.99	2.72	2.53	2.40	-
		✓	✓	3.03	2.74	2.55	2.41	-
	32768	✓		3.04	2.75	2.56	2.42	2.33
		✓	✓	3.05	2.76	2.57	2.42	2.32

Evaluation Context	3k	6k	10k	13k	16k
ChatGLM2-6B (Du et al., 2022)	0.88	0.46	0.02	0.02	0.02
MPT-30B-chat (Team, 2023a)	0.96	<b>1.0</b>	0.76	-	-
MPT-7B-storywriter (Team, 2023b)	0.46	0.46	0.28	0.34	0.36
LongChat-13B (Li et al., 2023)	<b>1.0</b>	<b>1.0</b>	<b>1.0</b>	<b>0.98</b>	0.9
Ours-13B	<b>1.0</b>	0.98	0.98	<b>0.98</b>	<b>0.94</b>

Size	Training Context Length	Evaluation Context Length						
		2048	4096	8192	16384	32768	65536	100,000
7B	100,000	3.36	3.01	2.78	2.60	2.58	2.57	2.52
13B	65536	3.20	2.88	2.66	2.50	2.39	2.38	-
70B	32768	2.84	2.57	2.39	2.26	2.17	-	-

# Model tells you what to discard: Adaptive kv cache compression for llms

---

**Suyu Ge<sup>1\*</sup>, Yunan Zhang<sup>1\*</sup>, Liyuan Liu<sup>2\*</sup>, Minjia Zhang<sup>2</sup>, Jiawei Han<sup>1</sup>, Jianfeng Gao<sup>2</sup>**

<sup>1</sup>University of Illinois Urbana-Champaign, <sup>2</sup>Microsoft

{suyuge2, yunanz2, hanj}@illinois.edu

{lucliu, minjiaz, jfgao}@microsoft.com

@article{

Conference={ICLR 2024}

Cites={6}

**Oral**

**8888**

# Method

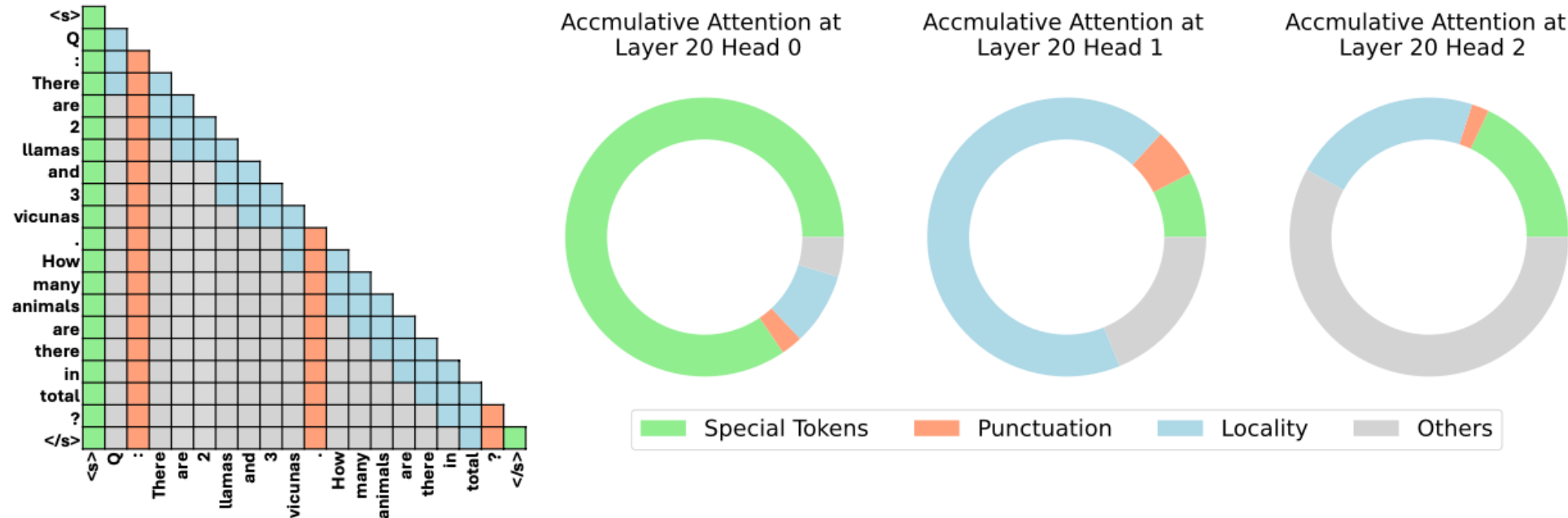


Figure 1: Different attention heads usually have different structures. Left: Four common attention structures (more details are elaborated in Section 3 and Section 4). Right: Attention map compositions of three attention heads that are in the same layer.

et al., 2023; Zhang et al., 2023; Liu et al., 2023a). We monitor for each token its cumulative sum of attention score, then treat these scores as token *frequency* and only **keep the most frequent tokens in the KV cache**. The length budget of frequent tokens over the current sequence length is controlled by a ratio  $r_f$ . This policy is referred to  $C_{\text{frequent}}$ .

## Method

对于初始层和最终层，它们有更多分配给全KV缓存的注意力头，表明这些层的注意力头可能关注所有令牌。同时，对于中间层，注意力图集中在特殊令牌上，这表明这些层的大多数注意力头主要关注特殊令牌

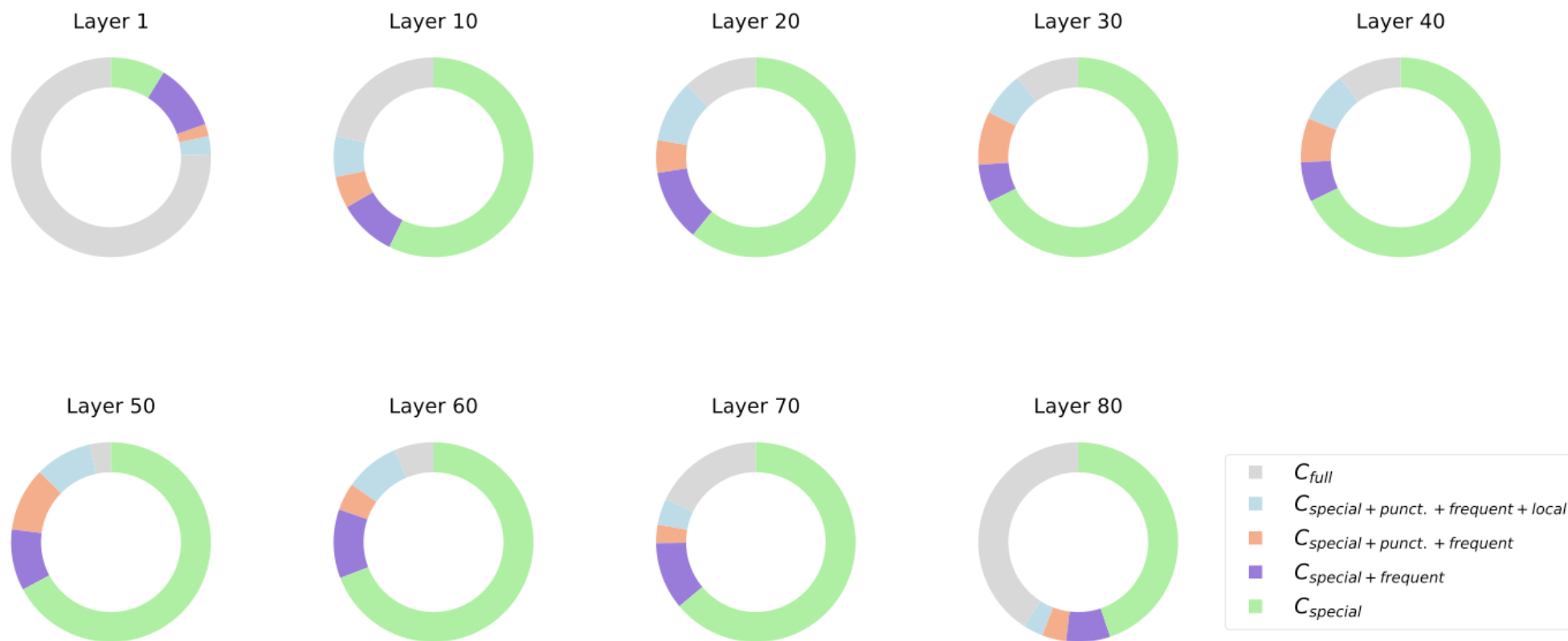


Figure 3: Attention profiling result distribution across different layers.

# Method

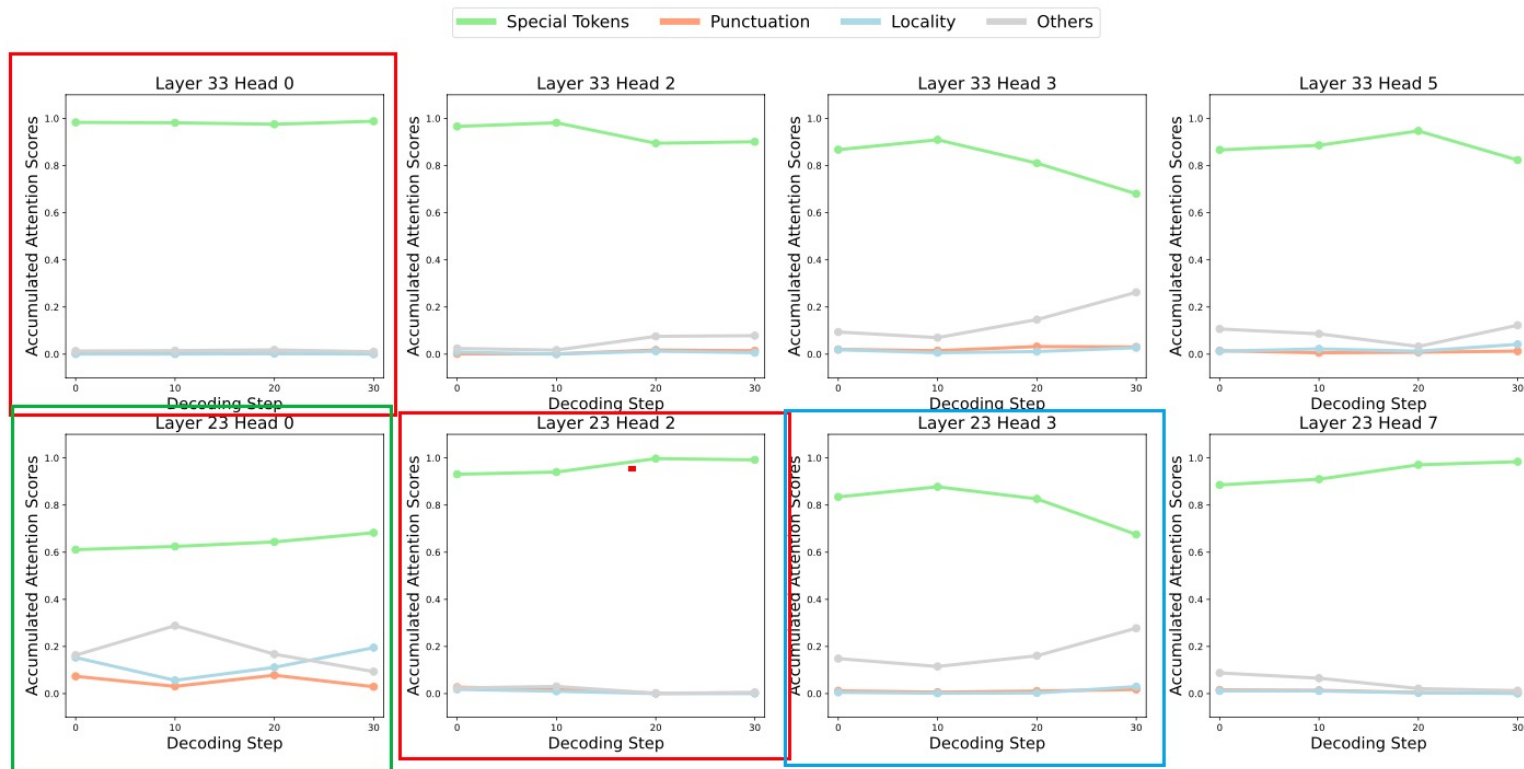


Figure 4: Accumulated attention score at 1st (prompt encoding), 10th, 20th, 30th decoding steps.

**Observation.** Despite some fluctuations of accumulated attention scores across time steps, the pattern of the attention maps remains relatively stable. For example, Layer 33 Head 0 and Layer 23 Head 2 almost only attend to the special token, while the locality and punctuation plays an important role in Layer 23 Head 0. As to Layer 23 Head 3, more than 10% of the attention score is allocated to the others portion, making it suitable for a uncompressed KV cache  $C_{full}$ .



**Algorithm 1:** FastGen–Prompt Encoding.**Input:** Feasible Policy Set ( $\mathcal{C}$ ), Prompt**Output:** Adaptive KV Cache

---

```

1 for Attention Head  $H_i$  in LLM do
2    $\mathbf{K}^i, \mathbf{Q}^i, \mathbf{V}^i \leftarrow H_i(\text{Prompt})$ 
3    $\mathbf{A}^i \leftarrow \text{softmax}(\mathbf{Q}^i \mathbf{K}^{iT})$ 
4    $\mathbf{C}^i \leftarrow \text{apply Equation 1 to } \mathbf{A}^i$ 
   /*  $\mathbf{C}^i$ : optimal policy */
5    $\mathbf{K}_{\mathbf{C}^i}^i, \mathbf{V}_{\mathbf{C}^i}^i \leftarrow f(\mathbf{K}^i, \mathbf{V}^i, \mathbf{C}^i)$ 
6    $\hat{\mathbf{K}}^i, \hat{\mathbf{V}}^i \leftarrow \mathbf{K}_{\mathbf{C}^i}^i \mathbf{V}_{\mathbf{C}^i}^i$ 
7 return  $\{\mathbf{C}^i, \hat{\mathbf{K}}^i, \hat{\mathbf{V}}^i\}$ 

```

---

**Algorithm 2:** FastGen–Token Generation.**Input:** Adaptive KV cache ( $\{\mathbf{C}^i, \hat{\mathbf{K}}^i, \hat{\mathbf{V}}^i\}$ )**Output:** Generated Text

---

```

1  $\mathbf{z}_0 \leftarrow \text{last prompt token}$ 
2 for  $j \in \{1, \dots, \text{Max Generate Length}\}$  do
3   for Attention Head  $H_i$  in LLM do
4      $\mathbf{K}^i, \mathbf{Q}^i, \mathbf{V}^i \leftarrow H_i(\mathbf{z}_{j-1}, \hat{\mathbf{K}}^i, \hat{\mathbf{V}}^i)$ 
5      $\mathbf{K}_{\mathbf{C}^i}^i, \mathbf{V}_{\mathbf{C}^i}^i \leftarrow f(\mathbf{K}^i, \mathbf{V}^i, \mathbf{C}^i)$ 
6      $\hat{\mathbf{K}}^i, \hat{\mathbf{V}}^i \leftarrow \mathbf{K}_{\mathbf{C}^i}^i, \mathbf{V}_{\mathbf{C}^i}^i$ 
7    $\mathbf{z}_j \leftarrow \text{sample from LLM prediction}$ 
8 return  $\{\mathbf{z}_j\}$ 

```

---

$$\mathbf{C}^* = \arg \min_{\mathbf{C} \in \mathcal{C}} \text{CacheMemoryCost}(\mathbf{C}) \quad \text{s.t.} \quad |\mathbf{A} - \text{softmax}(\mathbf{Q} \mathbf{K}_{\mathbf{C}}^T)| \leq 1 - T,$$

EXPERIMENT

从内存，时间分析

Model Size	KV Cache			$T$	Win rate
	Full	FastGen	Pruned ratio		
7B	32Gb	14Gb	56.6%	91%	30.8%
		20Gb	39.8%	95%	37.7%
		27Gb	16.9%	98%	47.4%
13B	63Gb	29Gb	53.4%	91%	32.0%
		38Gb	39.0%	95%	39.9%
		51Gb	18.3%	98%	48.7%
30B	158Gb	69Gb	56.7%	93%	37.0%
		81Gb	48.8%	95%	42.5%
		115Gb	27.4%	98%	47.5%
65B	320Gb	140Gb	56.3%	93%	40.9%
		176Gb	44.9%	95%	44.2%
		205Gb	36.0%	98%	49.8%

Generation Length	Overall Generation Duration (s)	Profiling Duration (s)	Decoding Time Per Token (s)	Profiling/Overall (%)
128	30.98	0.11	0.10	0.35%
256	50.1	0.11	0.10	0.21%
512	94.98	0.11	0.10	0.12%
1024	157.43	0.11	0.10	0.07%

Batch size	1				2			8		16
[prompt len, gen len]	[32,512]	[32,2048]	[32,8192]	[32,16384]	[512,32]	[512,512]	[4096,4096]	[512,512]	[4096,4096]	[512,512]
HF	13.35	57.37	299	799.14	1.12	19.16	167.64	23.44	OOM	OOM
DS	11.58	47.12	201.23	435.74	0.79	10.45	91.04	12.93	127.94	OOM
FastGen	11.21	44.6	179.43	359.83	0.73	9.71	76.93	10.57	82.16	OOM
Speed-up(%) over HF	16.03%	22.30%	40.00%	55.00%	34.80%	49.30%	54.10%	54.90%	-	OOM
Speed-up(%) over DS	3.20%	5.35%	10.83%	17.42%	7.59%	7.08%	15.50%	18.25%	35.78%	OOM

## EXPERIMENT

	Feasible Policy Set	Pruned KV Ratio	Win Rate
	$\mathcal{C}$	36.04%	49.75%
special	$\{C_{\text{punct.}}, C_{\text{punct.}+\text{frequent}}, C_{\text{punct.}+\text{frequent}+\text{local}}, C_{\text{full}}\}$	31.16%	47.64%
punct	$\{C_{\text{special}}, C_{\text{special}+\text{frequent}}, C_{\text{special}+\text{frequent}+\text{local}}, C_{\text{full}}\}$	34.23%	49.56%
-local	$\{C_{\text{special}}, C_{\text{special}+\text{punct.}}, C_{\text{special}+\text{punct.}+\text{frequent}}, C_{\text{full}}\}$	30.18%	49.06%
frequent	$\{C_{\text{special}}, C_{\text{special}+\text{punct.}}, C_{\text{special}+\text{punct.}+\text{local}}, C_{\text{full}}\}$	21.26%	46.08%

frequent和special是最重要的策略，移除它们将分别导致3.67%和2.11%的胜率下降

Cache Order	Pruned KV Ratio	Win Rate
$C_{\text{special}} \rightarrow C_{\text{punct.}} \rightarrow C_{\text{frequent}} \rightarrow C_{\text{local}}$	36.04%	49.75%
$C_{\text{special}} \rightarrow C_{\text{frequent}} \rightarrow C_{\text{local}} \rightarrow C_{\text{punct.}}$	36.40%	47.64%

Table 5: Policy order ablation on fine-tuned Llama 1-65B with AlpacaEval.

# EXPERIMENT

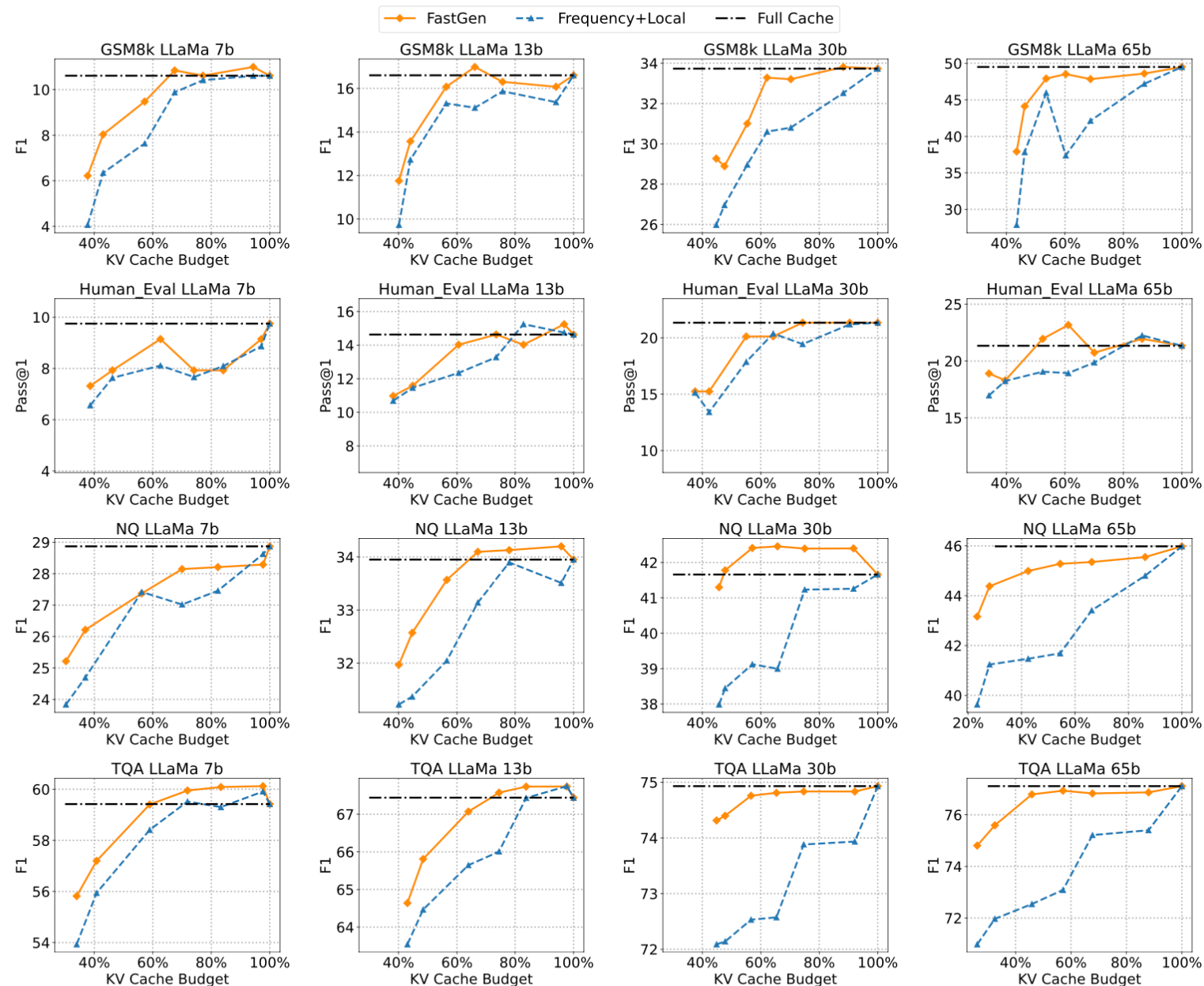


Figure 5: Performance of Adaptive KV Cache (FastGen) and Fixed KV Cache (Frequency+Local; Zhang et al., 2023 and Liu et al., 2023a) of Llama 1 on GSM8k, HumanEval, NQ, and TQA.