SharePrefill: Accelerating Prefilling for Long-Context LLMs via Sparse Pattern Sharing

Anonymous ACL submission

Abstract

Sparse attention methods exploit the inherent sparsity in attention to speed up the prefilling phase of long-context inference, mitigating the quadratic complexity of full attention computation. While existing sparse attention methods rely on predefined patterns or inaccurate estimations to approximate attention behavior, they often fail to fully capture the true dynamics of attention, resulting in reduced efficiency and compromised accuracy. Instead, we propose a highly accurate sparse attention mechanism that shares similar yet precise attention patterns across heads, enabling a more realistic capture of the dynamic behavior of attention. Our approach is grounded in two key observations: (1) attention patterns demonstrate strong inter-head similarity, and (2) this similarity remains remarkably consistent across diverse inputs. By strategically sharing computed accurate patterns across attention heads, our method effectively captures actual patterns while requiring full attention computation for only a small subset of heads. Comprehensive evaluations demonstrate that our approach achieves superior or comparable speedup relative to state-ofthe-art methods while delivering the best overall accuracy. The code will be made available upon publication.

1 Introduction

003

017

042

Long-context inference is essential for realworld applications of large language models (LLMs). Modern models like GPT-4.1 and Gemini 1.5 (Team et al., 2024) now support contexts up to one million tokens, advancing multi-document QA (Wang et al., 2024), code understanding (Bairi et al., 2024; Ziftci et al., 2025), and multi-turn dialogue (Zhang et al., 2025). Nonetheless, the prefilling phase of long-context inference remains time-consuming, as the vanilla attention mechanism entails quadratic computational complexity with respect to sequence length (Fu, 2024).



Figure 1: Comparison of our method with baselines across different models between latency (under 128K) and the average score on Infinitebench.

043

044

047

050

051

054

060

061

062

063

064

065

066

067

069

070

Sparse attention offers a promising solution by computing only significant attention scores, leveraging inherent sparsity in attention mechanisms. Many works have discovered various patterns with distinct characteristics and exploited them to perform sparse attention computations, such as the sink pattern in (Xiao et al., 2024), the A-shape, vertical-slash and block-sparse patterns in MInference (Jiang et al., 2024). However, these static patterns fail to generalize to varied inputs, as attention patterns inherently vary with different inputs as shown in Figure 2. To cope with the fundamental requirements of dynamic patterns, MInference dynamically adjusts the vertical-slash index, and FlexPrefill (Lai et al., 2025) further adapts the vertical-slash sparsity ratio dynamically and uses pooled queries and pooled keys to estimate queryaware block-wise patterns. However, we argue that pooling-based pattern estimation struggles to fully capture critical blocks due to the inaccuracies inherent in its approximations (detailed in Section 3).

Alternatively, we discover two interesting phenomena. Firstly, the sparse pattern of many attention heads tends to be highly similar. More importantly, the similarity relationships among these heads remain largely consistent, even though the sparse patterns themselves vary significantly across different inputs, as shown in Figure 2. Conse-

163

164

165

166

167

168

169

120

121

122

quently, we propose a highly accurate sparse attention mechanism that shares similar yet precise attention patterns across heads, mitigating reliance on
predefined patterns and avoiding inaccurate pattern
estimation. By computing dense attention using
only a subset of heads, the prefilling is accelerated
while preserving its high accuracy.

Our contributions are summarized as follows:

- We empirically demonstrate two fundamental properties of sparse attention patterns: (1) similarity across attention heads and (2) similarity consistency across different inputs.
- We propose SharePrefill, a novel highly accuracy-preserving sparse attention method to accelerate the prefilling phase by dynamically generating accurate sparse patterns and sharing them across heads.
 - We conduct extensive experiments on several different models and tasks and show that our proposed method achieves superior or comparable speedup to state-of-the-art methods while achieving the best overall accuracy.

2 Related Work

087

090

091

093

097

101

102

103

104

106

108

109

110

111

112

113

114 115

116

117

118 119 Existing sparse attention methods for accelerating model inference can be categorized into two types: training-free sparse attention and training-based sparse attention. The former relies on predefined sparse patterns or pattern estimation, while the latter involves training sparse models to dynamically predict sparse patterns during inference.

Training-free Sparse Attention Several methods focus on predefined attention patterns, such as shifted sparse attention (Chen et al., 2024), sink attention (Xiao et al., 2024) and the A-shape, verticalslash and block-sparse patterns used in MInference (Jiang et al., 2024). However, these patterns, often derived from limited cases, lack the flexibility to effectively adapt to varying input demands. MInference introduced partially dynamic patterns, by adjusting vertical-slash indexes based on inputs. FlexPrefill (Lai et al., 2025) adapts sparsity ratios via cumulative thresholds and incorporates queryaware sparse patterns to enhance flexibility. However, query-aware sparse patterns rely on pooled query and key representations for pattern selection, which may cause information loss and lead to less accurate pattern estimation. Our method aligns with this line of work but further enhances pattern modeling by dynamically providing more precise

sparse patterns through pattern sharing, thereby achieving better accuracy preservation.

Training-based Sparse Attention Trainingbased sparse attention methods introduce attention gates, train the gate-associated network, and automatically predict important sequence segments during inference. In this series of works, approaches like MoBA (Lu et al., 2025) and NSA (Yuan et al., 2025) continue training the entire model, while SeerAttention (Gao et al., 2024) employs a linear layer as a learnable gate, training only the attention gate. Even though training-based sparse attention methods show promising acceleration while maintaining accuracy, the cost of resource-intensive and time-consuming training hinders their widespread practical applicability.

3 Static Patterns and Pooling-based Pattern Estimation are Not Enough

Attention patterns are highly dynamic, showing substantial variation both across different heads and within the same head under different inputs, as shown in Figure 2. In particular, the staircase-like patterns in *En.Dia* and the highly irregular patterns in *Code.Debug* deviate significantly from previously proposed static patterns like the vertical-slash pattern (Jiang et al., 2024). The highly dynamic nature intrinsic to attention mechanisms exposes the limitations of fixed-pattern approaches and underscores the need for adaptive, dynamic attention modeling techniques.

FlexPrefill (Lai et al., 2025) leverages pooled queries (Q) and keys (K) to estimate the average attention scores within each block for identifying critical regions, thus alleviating the reliance on predefined patterns. However, we highlight that this pooling-based method struggles to fully capture important blocks, and we identify that this challenge is rooted in two critical aspects.

Disregard for Token Alignment: The pooling operation disregards token-level position alignment within the query (Q) and key (K) segments, while attention mechanisms are inherently sensitive to token-level position alignment. This discrepancy leads to pooled results $pool(Q) \cdot pool(K)$ that cannot accurately estimate the average of actual attention scores for the block. For example, consider two 1-dimensional Q, K for 3 tokens: Q=[0, 0,1], K=[0, 1, 0]. Due to ignoring position alignment, the $pool(Q) \cdot pool(K)$ appears slightly significant($\frac{1}{9}$). However, all attention scores within



(a) Visualization of attention patterns for different heads across various tasks. Each group of three columns corresponds to the heads within a specific task.



(b) Similarity matrices show the pattern similarity between each head and other heads across different tasks.

Figure 2: Attention patterns of different heads and their similarity matrices across various tasks.

the block are actually zero, leading to an *overestimation* of the block's importance.

170

171

172

174

176

178

179

180

183

190 191

193

195

Smoothing of High-/Low- Values: The pooling operation smooths out high and low values within Q and K, which often contribute to high and low attention scores, resulting in inaccurate importance estimation. For instance, Q=[0, 0, 1], K=[0, -1, 1]. During pooling, the the high-value and low-value elements in Q and K are diluted, resulting in $pool(Q) \cdot pool(K) = 0$, which is less than the actual average of attention scores $pool(Q \cdot K) = \frac{1}{9}$, leading to an *underestimation* of the importance of the block.

4 Observation: Dynamic Attention Heads Exhibit Similar Patterns and Static Similarity Relationships

We present the foundational observations motivating our method: different heads exhibit high similarity, and the similarity remains highly consistent across varying inputs and tasks. Specifically, these observations distill into two key properties:

 (1) *Inter-head Pattern Similarity*: We observe many similar sparse patterns across attention heads, both within and between layers, as shown in Figure 2(a). These patterns are derived from the Llama-3-8B-Instruct-262k model using samples from different tasks in InfiniteBench (Zhang et al., 2024), with each group of three columns corresponding to heads from a specific task. For example, heads such as (L18, H4), (L22, H2), and (L25, H7) in the *En.Dia* task exhibit highly consistent staircaselike patterns, where L is the layer index and H is the head index. Additionally, Figure 2(b) shows the statistical similarity matrix based on Jaccard similarity scores (# intersection / # union) between one head and all others. This measure avoids artificially high similarity values that could arise from the presence of many zeros in these sparse patterns. Notably, a large number of similarity scores exceed 0.5, indicating that each head has many similar counterparts among the others.

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

222

223

224

225

226

227

228

231

232

233

234

(2) *Cross-input Similarity Consistency*: More importantly, the similarity among attention heads remains consistent regardless of the specific input or task, even though the pattern of a given attention head varies across different inputs and tasks. Figure 2(a) shows that (L18, H4), (L22, H2), and (L25, H7) are highly similar in *Code.Debug*, consistent with their previously observed similarity in *En.Dia* (Property 1), though their patterns differ between the two tasks. This highlights the consistent similarity among attention heads across inputs and tasks, suggesting that sparse patterns can potentially transfer to similar heads, regardless of context.

5 Proposed Approach

5.1 Problem Formulation

Generally, our goal is to replace dense attention with sparse attention in attention layers to reduce computational costs during the pre-filling phase, while minimizing the output loss of each attention layer, thereby accelerating the pre-filling process while preserving accuracy. This can be formulated as a multi-objective optimization problem:

$$\min_{M} |\mathcal{A}(Q, K, V, M) - \mathcal{A}(Q, K, V)|,$$

$$\min_{M} |M|$$
(1)

where,

$$\mathcal{A}(\boldsymbol{Q}, \boldsymbol{K}, \boldsymbol{V}, \boldsymbol{M}) = \sigma \left(\frac{1}{\sqrt{d}}\boldsymbol{Q}\boldsymbol{K}^{T} - c(1-M)\right)\boldsymbol{V},$$
 23

$$\mathcal{A}(\boldsymbol{Q}, \boldsymbol{K}, \boldsymbol{V}) = \sigma \left(\frac{1}{\sqrt{d}} \boldsymbol{Q} \boldsymbol{K}^{T}\right) \boldsymbol{V}$$
 237

We define $\mathcal{A}(Q, K, V, M)$ as our sparse attention, where M is a binary mask indicating 239

332

333

334

335

336

337

338

340

291

the sparse pattern applied in the sparse attention 240 computation, where 1 means the block is com-241 puted, and 0 means it is skipped. The output is 242 $O = \mathcal{A}(Q, K, V, M)$ and a sufficiently large constant c ensures that attention score is approximately 244 zero, whenever $M_{ij} = 0$. Here, σ denotes the 245 softmax function. The primary optimization ob-246 jective is to identify a pattern M that minimizes 247 the variance between the attention matrices of full attention $\mathcal{A}(Q, K, V)$ and our proposed sparse at-249 tention $\mathcal{A}(Q, K, V, M)$, while also reducing the computational time required for sparse attention computation and sparse pattern generation.

5.2 Accelerating Prefilling via SharePrefill

253

256

258

259

262

264

269

272

273

274

281

Our main idea is to compute the full attention for a subset of heads, identify the actual sparse patterns, and share these patterns with other heads that are known to exhibit similar behavior. This approach enables the remaining heads to approximate the actual patterns without computing each one separately, thus maintaining the model's original accuracy while accelerating inference.

As depicted in the overview of SharePrefill in Figure 3, our approach involves two key components: (1) Offline clustering to group heads based on the similarity of their attention score maps. (2) Online inference, where pivotal attention is constructed dynamically and shared with other heads during the inference process. The overall algorithm of SharePrefill is detailed in Algorithm 1.

Offline Clustering of Similar Heads We cluster attention heads across layers into distinct groups based on the similarity of their attention score maps, performed in offline mode. The pre-computed clusters serve as the foundation for constructing and sharing sparse attention patterns within each cluster during inference.

Given the consistent similarity between heads, we perform clustering on their attention score maps using a sample from the *Retr.KV* task in InfiniteBench. We first obtain compressed lowdimensional representations of the attention scores by training an autoencoder network on these attention score maps (the network architecture is illustrated in Appendix C). Next, we normalize the representations and apply hierarchical clustering with a distance threshold to group similar heads into clusters, while isolating dissimilar heads as noise clusters. Notably, we only store the layer and head indices within clusters, rather than the sparse patterns themselves. The actual sparse patterns are dynamically generated during online inference, ensuring adaptability to varying inputs.

Dynamic Pattern Construction and Sharing during Inference During inference, we construct attention patterns dynamically and share adaptive yet accurate patterns among similar heads to facilitate sparse attention computation. This is achieved by computing the attention output for each layer while maintaining an evolving global pattern dictionary, which serves as the basis for sharing patterns across similar heads.

In general, the online inference process—based on dynamic pattern construction and sharing—comprises three key steps, corresponding to the three sub-algorithms outlined in Algorithm 1. The algorithm takes the query matrix Q, key matrix K, value matrix V, similarity threshold τ , sparsity threshold δ and cumulative attention threshold γ as input. For simplicity, the illustration focuses on a single head rather than a layer. However, in practice, we perform sparse attention computation layer-by-layer.

Pivotal Pattern Sharing (See Algorithm 4): Before performing sparse attention computation, we first query the global pivotal pattern dictionary to check if a pivotal pattern is available for reuse. If a pivotal pattern exists, it is shared with the corresponding head; otherwise, the head computes full attention using a dense pattern (i.e., a pattern with all ones).

Sparse Attention Computation (See Algorithm 1): We then perform sparse attention computation to obtain the output for the current layer while simultaneously computing the block-wise average of QK values, denoted as \hat{A} , which captures the average QK scores within each block (line 8 in overall Algorithm 1). The sparse attention kernel is implemented in Triton (Tillet et al., 2019), following the block-wise strategy from FlashAttention 2 (Dao, 2024) and incorporating a block-wise sparse pattern to determine computation regions. Only blocks labeled as 1 in the sparse pattern are computed, while those labeled as 0 are skipped. During the computation for the final output, for each block where the pattern value is one, we compute the average QK value; for blocks where the pattern value is 0, we assign the average QK value as $-\infty$.

Pivotal Pattern Construction (See Algorithm 2): Subsequently, we use the obtained block-



Figure 3: Overview of proposed SharePrefill. Attention heads are clustered offline based on the similarity of their attention score maps to create a static head dictionary. During inference, each head retrieves its cluster index C_i . Pivotal Patterns are shared if available; otherwise, a dense pattern is assigned. The sparse attention output O is computed using M, and \tilde{A} updates the dynamic pivotal pattern dictionary.

wise average QK values to compute the block-wise average attention scores after applying softmax. These scores are then used to construct new pivotal patterns by applying a cumulative score threshold γ , which selects the minimal number of blocks required to cover the target cumulative attention score, as detailed in Algorithm 2. The resulting patterns are then updated into the pivotal pattern dictionary.

341 342

343

344

347

351

352

364

367

370

371

To ensure safe dynamic pattern sharing, we verify similarity before sharing patterns to prevent incorrect sharing that could adversely impact accuracy. Specifically, we compute the Jensen-Shannon (JS) distance between the block-wise average attention score of the last row block of the current head \hat{a} and the corresponding pivotal block-wise average attention score of the last row \tilde{a} , which is also stored in the pivotal pattern dictionary (line 6 in Algorithm 3). This distance serves as a measure to predict the similarity between the current head and its corresponding pivotal head. If the JS distance is below a given similarity threshold τ , we share the pivotal pattern with the current head. Otherwise, we fall back to a conservative verticalslash pattern (lines 7-11 in Algorithm 3) using a cumulative threshold-based vertical-slash pattern search algorithm (outlined in Algorithm 5), as proposed in FlexPrefill (Lai et al., 2025). Additionally, noisy clusters, which include dissimilar patterns, also revert to the vertical-slash pattern.

To enhance efficiency, we exclude highly sparse

heads from the pivotal pattern construction and sharing process, as we consider that computing full attention for these heads to derive pivotal patterns is not cost-effective in terms of acceleration. For highly sparse heads, we instead fall back to searching for a vertical-slash pattern for each head (see line 10 in Algorithm 3), as the pattern often serves as a suitable approximation for highly sparse heads (Jiang et al., 2024). To identify these highly sparse heads, we compute the Jensen-Shannon (JS) distance between the block-wise average attention score of the last row block of the current head \hat{a} and a uniform distribution u (see line 6 in Algorithm 3). We then compare this distance to a predefined sparsity threshold δ . If the JS distance is not less than the threshold, we classify the head as a highly sparse head.

373

374

375

377

379

380

381

385

386

391

392

393

395

396

398

400

6 Experiments

6.1 Settings

This section outlines the models, datasets, baselines, and implementation details of our method in comparison with baseline methods. Additional information is provided in Appendix A.

Models, Datasets, and Baselines We employ two cutting-edge, renowned long-context LLMs: (i) Llama-3-8B-Instruct-262k (Pekelis et al., 2024), (ii) Qwen2.5-7B-Instruct (Team, 2024). The models are evaluated on InfiniteBench (Zhang et al., 2024), a state-of-the-art public benchmark de-

Algorithm 1 Sparse Attention

1: Input: $Q, K, V \in \mathbb{R}^{N \times d_h}$; δ, τ, γ ; l, h

Decide the pattern type based on Q, K, sparsity threshold δ and similarity threshold τ

- *pattern* ← Determine Sparse Pattern (Q, K, δ, τ)
 # Decide the sparse pattern M based on pattern, Q, K and pattern threshold γ
- 3: **if** *pattern* == shared_pivot **then**
- 4: $M \leftarrow \text{Share Pivotal Pattern}(l, h)$
- 5: **else if** *pattern* == vertical_slash **then**
- 6: $M \leftarrow \text{Search Vertical Slash Pattern}(Q, K, \gamma)$
- 7: end if

Compute the output O and block-averaged QK values \tilde{A} by applying sparse pattern M

8: $O, \tilde{A} \leftarrow \mathcal{A}(Q, K, V, M)$

Construct and update global dynamic pivotal patterns via the newest block-averaged QK values \tilde{A} and pattern threshold γ

9: Construct Pivotal Pattern(Ã, γ, l, h) return O

Algorithm 3 Determine Sparse Pattern

1: **Input:** *Q*, *K*; δ, τ

Take a representative query subset

- 2: select \$\hat{Q} = Q_{[-block_size:]}\$
 # Compute estimated block-averaged average attention \$\hat{a}\$ and pivotal block-averaged attention \$\hat{a}\$
- 3: $\hat{a} \leftarrow \operatorname{softmax}(\operatorname{pool}(\hat{Q}K^T) / \sqrt{d})$ # Retrieve cluster index c in head_dict,
- 4: c ← lookup(l, h; head_dict)
 # Fetch the pivotal representative ã
- 5: ã ← lookup(c; pivotal_pattern_dict)
 # Compute sparsity and similarity divergence

6: $d_{sparse} \leftarrow \sqrt{JSD(\hat{a}||u)}, d_{sim} \leftarrow \sqrt{JSD(\hat{a}||\tilde{a})}$

- # Determine whether to use pattern sharing strategy
- 7: **if** $d_{sparse} < \delta$ and $d_{sim} < \tau$ **then**
- 8: $pattern \leftarrow shared_pivot$
- 9: else
- 10: $pattern \leftarrow vertical_slash$

```
11: end if
```

return pattern

Algorithm 2 Construct Pivotal Pattern

Input: \tilde{A} ; γ ; l, h

if \vec{A} is fully attention computed then

Compute block-averaged attention score $ilde{A}$ by applying softmax on block-averaged QK values

 $\tilde{A} = \operatorname{softmax}(\tilde{A})$

Take the last row of \tilde{A} as pivotal representative $\tilde{a} \leftarrow \tilde{A}_{[-1:]}$

Flatten and normalize attention score map

 $\tilde{A} \leftarrow \text{flatten}(\tilde{A} / \sum_{i,j} A[I,j])$

Sort attention <u>s</u>cores

 $I_a \leftarrow argsort(\tilde{A})$

Obtain the minimum computational budget making the sum of the scores exceeds γ

 $K \leftarrow \min \{ k: \sum_{i \in I_a[1:k]} \tilde{A}[i] \ge \gamma \}$

Select index set

 $S \leftarrow I_a[1:K]$

Convert index set S to mask pattern M $M \leftarrow index_to_mask(S)$

Lookup cluster index c in head_dict c \leftarrow lookup(l, h; head_dict)

Update M and \tilde{a} into pivotal_pattern_dict pivotal_pattern_dict.update({c: (\tilde{a} , M)})

end if

Algorithm 4 Share Pivotal Pattern

Input: *l*, *h*

Retrieve cluster index c in head_dict c \leftarrow lookup(l, h; head_dict)

Fetch the pivotal sparse pattern from the dynamic pivotal_pattern_dict M

$M \leftarrow \text{lookup}(c; pivotal_pattern_dict)$ if M not exist then

Assign a dense pattern to the first head within the cluster c for subsequent full attention computation $M \leftarrow ones$

else

Share existing pivotal pattern Mend if

return M

Models	Methods	En.Sum	En.QA	En.MC	En.Dia	Zh.QA	Code.Debug	Math.Find	Retr.PassKey	Retr.Number	Retr.KV	Avg
Llama-3-8B-Instruct-262k	FlashAttn	25.88	8.63	67.69	5.00	12.66	20.81	26.57	100.00	100.00	14.40	38.16
	FlexPrefill	19.91	12.60	57.21	5.50	11.63	22.84	20.86	100.00	100.00	13.80	36.44
	MInference	25.51	8.50	65.94	<u>8.00</u>	12.14	22.08	32.86	100.00	100.00	16.40	39.14
	Ours	<u>20.24</u>	8.00	63.32	11.84	11.94	<u>24.11</u>	<u>30.00</u>	100.00	100.00	<u>21.00</u>	39.05
	Ours ($\delta = 1.01$)	19.35	<u>11.74</u>	<u>64.63</u>	5.50	<u>11.96</u>	28.17	29.14	100.00	100.00	23.00	39.35
Qwen2.5-7B-Instruct	FlashAttn	15.53	3.18	35.81	10.50	3.95	14.47	38.57	100.00	93.56	0.00	31.56
	FlexPrefill	14.20	3.09	31.88	8.00	3.54	<u>15.99</u>	9.43	<u>97.29</u>	75.42	0.00	25.88
	MInference	<u>14.83</u>	2.86	34.93	9.00	<u>3.81</u>	14.97	<u>38.29</u>	96.78	<u>76.78</u>	0.00	29.23
	Ours	15.31	<u>2.88</u>	38.43	<u>8.50</u>	3.99	17.26	44.57	99.49	87.46	0.00	31.79

Table 1: Performance comparison of different methods on various models and tasks on InfiniteBench. The best and second-best results are highlighted in bold and underlined, respectively.

Methods	En.Sum	En.QA	En.MC	En.Dia	Zh.QA	Code.Debug	Math.Find	Retr.PassKey	Retr.Number	Retr.KV	Avg 🗄	128K Latency (s)
Our w/o Sharing (τ =0) Our w/o Exclusion (δ =1 01)	19.68	11.86 11.74	63.76 64.63	9.00 5.50	11.65 11.96	23.86 28.17	25.14 29.14	22.00 100.00	100.00	22.00 23.00	38.70 39.35	$\frac{17.01}{20.02}$
Ours	20.24	8.00	63.32	11.84	11.94	24.11	30.00	100.00	100.00	21.00	39.05	16.92

Table 2: Performance of ablation methods evaluated using LLaMA-3-8B-Instruct-262K on InfiniteBench.



Figure 4: Perplexity results on PG-19 (Rae et al., 2020) using different models and methods.



Implementation Details All our experiments 418 were conducted on a single NVIDIA A100 GPU 419 with 80GB of memory. For baseline implementa-420 tions, we use the official FlashAttention 2 pack-421 age 6, and adopt the official MInference reposi-422 423 tory 5 for both MInference and FlexPrefill, which includes the FlexPrefill implementation. For MIn-424 ference, we employ the default vertical-slash pat-425 tern configuration available in its code repository. 426 For FlexPrefill, we use the default parameters with 427



Figure 5: Latency comparison of different approaches across various context lengths using different models.

the sparse pattern threshold $\tau = 0.1$ and the cumulative pattern threshold $\gamma = 0.9$ consistently for all models. For a fair comparison, we also set the cumulative pattern threshold $\gamma = 0.9$ in our method; The similarity threshold τ is set to 0.2 and the sparsity threshold δ to 0.3, unless otherwise specified. Additionally, all the baseline methods employ sparse computation during prefilling and transition to dense computation during the decoding phase.

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

6.2 **Main Results**

We compare our method with baselines and present the main results on the aforementioned benchmarks and models. The results demonstrate that our approach achieves superior or comparable speedup while delivering the overall best accuracy.

InfiniteBench Table 1 shows that our method preserves most of the model's performance, achieving overall best accuracy maintenance. While our method with default parameters shows slightly lower accuracy compared to MInference, it achieves significantly lower latency, as shown in Figure 5 and Figure 1. However, our method outperforms MInference in both accuracy and efficiency by sharing all similar attention heads, including highly sparse ones.



Figure 6: Distribution of three sparse attention head patterns in LLaMA-3-8B-Instruct-262K.

Language Modeling We evaluate our method against baselines on the language modeling task based on the PG-19 dataset (Rae et al., 2020). As shown in Figure 4, the perplexity of our method closely approaches the performance of MInference and FlashAttention 2, with the gap between them being within about 1.0. Moreover, the perplexity score of our method is significantly lower than that of FlexPrefill, with reductions of approximately 1.0~4.0 in Qwen2.5-7B-Instruct and over 1.0 in Llama-3-8B-Instruct-262k. These results demonstrate the strong language modeling capabilities of our approach.

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

Performance vs. Latency Figure 5 shows the latency across different context windows under Llama-3-8B-Instruct-262k and Qwen2.5-7B-Instruct on a single A100. The results demonstrate that our method achieves better or comparable speedup compared to the baselines. Moreover, Figure 1 compares our method with baselines under different models in terms of model performance on InfiniteBench and average latency under 128K on the latency benchmark. The results indicate that our method achieves a favorable tradeoff between accuracy preservation and inference speedup.

7 Ablation Study

Analysis of Different Components To evaluate the contributions of different components in SharePrefill, we introduce two variants for the ablation study: (1) Ours w/o sharing, which uses only the vertical-slash pattern without pivotal pattern sharing mechanism, corresponding to a similarity threshold $\tau = 0$; (2) Ours w/o exclusion, where removing highly sparse heads strategy and all similar heads participate in the pattern sharing mechanism, corresponding to a sparsity threshold

 $\delta = 1.01$ (selected to account for boundary conditions, ensuring that patterns with $\delta = 1$ meet the sharing-consideration criterion in line 7 of Algorithm 3). Table 2 presents the ablation results on 493 Llama-3-8B-Instruct-262k. It first demonstrates 494 that removing the pattern sharing mechanism leads 495 to performance degradation, confirming the necessity of our pattern sharing strategy in preserving accuracy. Additionally, removing the strategy of excluding highly sparse heads-where all similar heads, including highly sparse ones, are allowed 500 to share patterns-results in reduced speedup but 501 improved performance. This demonstrates that the strategy of excluding highly sparse heads enhances 503 efficiency while potentially degrading the model's 504 accuracy maintenance potential. The observed ac-505 curacy improvement when removing the exclusion 506 strategy can be attributed to more similar heads 507 participating in pattern sharing, rather than being 508 forced into predefined vertical-slash patterns. This 509 further validates the effectiveness of our pattern 510 sharing mechanism in maintaining accuracy. 511

490

491

492

496

497

498

499

502

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

Pattern Distribution Figure 6 shows the distribution of dense, shared, and conservative verticalslash patterns used in Llama-3-8B-Instruct-262k. The majority of attention heads adopt the verticalslash pattern, while only a small number require the full-attention dense pattern—typically just 1 to 4 heads in total. Although the number of shared patterns is limited, they play a significant role in maintaining model accuracy as shown in Table 2.

8 **Conclusion and Future Work**

In this paper, we observe that attention heads exhibit similarity, and this similarity remains consistent across different inputs. Built on these observations, we propose a novel sparse attention method that highly preserves accuracy while accelerating the prefilling phase. Our method achieves this by dynamically generating accurate patterns and sharing them with other similar heads, thereby capturing more realistic attention dynamics. We conduct extensive experiments across several models and tasks, demonstrating that our proposed method achieves superior or comparable speedups to stateof-the-art approaches while delivering the highest accuracy maintenance. The principle of similarity between heads and the proposed pattern-sharing mechanism holds the potential for accelerating the decoding phase and extending to multi-modular systems, which will be explored in future work.

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

591

592

Limitations

540

558

562

563

564

565

567

568

569

570

571

572

573

574

575

576

577

578

580

581

583

585

586 587

588

590

Although we provide observational and statistical 541 evidence on the similarity properties between at-542 tention heads, the underlying explanation for the 543 highly consistent similarity relationships among 544 heads across different inputs remains unclear. This 545 546 open question requires further investigation. Additionally, while our approach demonstrates effec-547 tiveness in LLM prefilling on single devices, its scalability to larger-scale scenarios requires further study. Future work will focus on evaluating and 550 551 further enhancing the scalability of the proposed approach. This includes exploring efficient patternsharing mechanisms in scaled scenarios, such as 553 allowing each device to maintain a local partial 554 pivotal pattern dictionary or enabling a global dic-555 556 tionary to be shared across devices through interdevice communication. 557

References

- Ramakrishna Bairi, Atharv Sonwane, Aditya Kanade, Arun Iyer, Suresh Parthasarathy, Sriram Rajamani, B Ashok, and Shashank Shet. 2024. Codeplan: Repository-level coding using llms and planning. *Proceedings of the ACM on Software Engineering*, 1(FSE):675–698.
- Yukang Chen, Shengju Qian, Haotian Tang, Xin Lai, Zhijian Liu, Song Han, and Jiaya Jia. 2024. Longlora: Efficient fine-tuning of long-context large language models. In *The Twelfth International Conference* on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024.
- Tri Dao. 2024. Flashattention-2: Faster attention with better parallelism and work partitioning. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11,* 2024.
- Yao Fu. 2024. Challenges in deploying long-context transformers: A theoretical peak performance analysis. *Preprint*, arXiv:2405.08944.
- Yao Fu, Litu Ou, Mingyu Chen, Yuhao Wan, Hao Peng, and Tushar Khot. 2023. Chain-of-thought hub: A continuous effort to measure large language models' reasoning performance. *arXiv preprint arXiv:2305.17306*.
- Yizhao Gao, Zhichen Zeng, Dayou Du, Shijie Cao, Peiyuan Zhou, Jiaxing Qi, Junjie Lai, Hayden Kwok-Hay So, Ting Cao, Fan Yang, and 1 others. 2024. Seerattention: Learning intrinsic sparse attention in your llms. *arXiv preprint arXiv:2410.13276*.
- Huiqiang Jiang, Yucheng Li, Chengruidong Zhang, Qianhui Wu, Xufang Luo, Surin Ahn, Zhenhua Han,

Amir Abdi, Dongsheng Li, Chin-Yew Lin, and 1 others. 2024. Minference 1.0: Accelerating pre-filling for long-context llms via dynamic sparse attention. *Advances in Neural Information Processing Systems*, 37:52481–52515.

- Xunhao Lai, Jianqiao Lu, Yao Luo, Yiyuan Ma, and Xun Zhou. 2025. Flexprefill: A context-aware sparse attention mechanism for efficient long-sequence inference. In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025.*
- Enzhe Lu, Zhejun Jiang, Jingyuan Liu, Yulun Du, Tao Jiang, Chao Hong, Shaowei Liu, Weiran He, Enming Yuan, Yuzhi Wang, and 1 others. 2025. Moba: Mixture of block attention for long-context llms. *arXiv preprint arXiv:2502.13189*.
- Leonid Pekelis, Michael Feil, Forrest Moret, Mark Huang, and Tiffany Peng. 2024. Llama 3 gradient: A series of long context models.
- Jack W. Rae, Anna Potapenko, Siddhant M. Jayakumar, Chloe Hillier, and Timothy P. Lillicrap. 2020. Compressive transformers for long-range sequence modelling. In 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020.
- Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, and 1 others. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*.
- Qwen Team. 2024. Qwen2.5: A party of foundation models.
- Philippe Tillet, Hsiang-Tsung Kung, and David Cox. 2019. Triton: an intermediate language and compiler for tiled neural network computations. In *Proceedings of the 3rd ACM SIGPLAN International Workshop on Machine Learning and Programming Languages*, pages 10–19.
- Yu Wang, Nedim Lipka, Ryan A Rossi, Alexa Siu, Ruiyi Zhang, and Tyler Derr. 2024. Knowledge graph prompting for multi-document question answering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 19206–19214.
- 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, ICLR 2024, Vienna, Austria, May 7-11, 2024.*
- Jingyang Yuan, Huazuo Gao, Damai Dai, Junyu Luo, Liang Zhao, Zhengyan Zhang, Zhenda Xie, YX Wei, Lean Wang, Zhiping Xiao, and 1 others. 2025. Native sparse attention: Hardware-aligned and natively trainable sparse attention. *arXiv preprint arXiv:2502.11089*.

- 646 647
- 648 649

652

654

655

659

667

672

674

675

679

682

685

- Guibin Zhang, Luyang Niu, Junfeng Fang, Kun Wang, Lei Bai, and Xiang Wang. 2025. Multi-agent architecture search via agentic supernet. *arXiv preprint arXiv:2502.04180*.
- Xinrong Zhang, Yingfa Chen, Shengding Hu, Zihang Xu, Junhao Chen, Moo Hao, Xu Han, Zhen Thai, Shuo Wang, Zhiyuan Liu, and 1 others. 2024. ∞ 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.
- Celal Ziftci, Stoyan Nikolov, Anna Sjövall, Bo Kim, Daniele Codecasa, and Max Kim. 2025. Migrating code at scale with llms at google. *arXiv preprint arXiv:2504.09691*.

A More Details on Experimental Settings

A.1 Models

We employ two state-of-the-art long-context language models: Llama-3-8B-Instruct-262k⁻¹(released under the Meta Llama License) and Owen2.5-7B-Instruct² (released under the Apache 2.0 License). These models were selected due to their strong capabilities in handling long-context understanding tasks, with Llama-3-8B-Instruct-262k supporting contexts of up to 262K tokens and Qwen2.5-7B-Instruct supporting contexts of up to 128K tokens. Both models support multiple languages, primarily English, with Qwen2.5-7B-Instruct also demonstrating excellent performance in Chinese. Additionally, Qwen2.5-7B-Instruct supports up to 128K tokens and demonstrates excellent multilingual performance, with particular strength in Chinese. For further details, refer to the model repositories, as listed in 1 and 2

A.2 Datasets & Benchmarks

• InfiniteBench InfiniteBench³ (Zhang et al., 2024) is publicly released under the Apache-2.0 License. It is a state-of-the-art benchmark designed to evaluate long-context language models with context lengths exceeding 100K tokens. The benchmark consists of 12 unique tasks, each carefully crafted to assess different

aspects of language processing and comprehension in extended contexts. These tasks encompass a mix of real-world scenarios and synthetic constructs, including novels, dialogues, code, and math, ensuring a comprehensive evaluation of model capabilities. In our experiments, we compare our method's long-context performance against baselines across 10 tasks, using all available samples. Consistent with MInference and FlexPrefill, we excluded *Code.Run* and *Math.Calc* because they are highly challenging, with fullattention models often scoring near 0. 690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

- PG-19 Language Modeling Benchmark (Rae et al., 2020) proposed a long-context language modeling benchmark ⁴ that evaluates perplexity on the PG-19 dataset and whose repository is released under the Apache 2.0 license. Perplexity quantifies how well a model predicts the next token in a sequence and is commonly used to assess the language modeling performance of long-context LLMs on extended texts. PG-19 contains books with lengths of up to 500K tokens, making it wellsuited for long-context evaluation. To assess language modeling performance across different context lengths, we conduct experiments using 100 randomly selected samples from the PG-19 dataset, truncating them to various lengths ranging from 1K to 104K tokens. We then report the average perplexity based on these truncated samples. Due to high memory usage during perplexity computation, we evaluate contexts with lengths of up to 104K tokens.
- Latency Benchmark We follow the latency benchmarks provided in MInference ⁵ (Jiang et al., 2024), which is released under the MIT License. The prompts, sourced from the Chain-of-Thought Hub (Fu et al., 2023) (also released under the MIT License), were trimmed to varying token lengths to measure the prefilling stage latency. To ensure reliable measurements, we conduct ten repeated experiments after a warm-up phase and report the average latency.

¹https://huggingface.co/gradientai/ Llama-3-8B-Instruct-Gradient-262k

²https://huggingface.co/Qwen/Qwen2.

⁵⁻⁷B-Instruct

³https://huggingface.co/datasets/ xinrongzhang2022/InfiniteBench

⁴https://github.com/google-deepmind/pg19

⁵https://github.com/microsoft/MInference.git

A.3 Baselines

736

737

738

740

741

743

744

745

747

748

751

752

753

755

757

762

763

764

765

770

772

775

777

778 779

780

FlashAttention 2⁶ (Dao, 2024): Flash Attention 2 is an I/O-aware exact attention algorithm designed to improve the efficiency of dense attention computation. It leverages tiling techniques to minimize the number of memory read and write operations between GPU high-bandwidth memory (HBM) and on-chip SRAM, thereby significantly reducing memory overhead and improving computational throughput. As an optimized implementation of dense attention, Flash Attention 2 enables faster and more scalable transformer inference and training, especially in long-sequence scenarios.

- 2. MInference (Jiang et al., 2024): MInference is a state-of-the-art sparse attention mechanism that exploits the static patterns observed in the attention mechanisms of LLMs, aiming to accelerate the prefilling phase for long-context inputs. It first determines offline which sparse pattern each attention head belongs to. During inference, it approximates the sparse indices online and dynamically computes attention using optimized custom kernels. This design enables significant speedup while maintaining strong accuracy.
 - 3. **FlexPrefill** (Lai et al., 2025): FlexPrefill is another state-of-the-art sparse attention mechanism that enhances flexibility by incorporating cumulative-attention-based index selection and query-aware sparse patterns, enabling more adaptive sparse attention during the prefilling phase of LLM inference.

A.4 Implementation Details

For offline clustering, we train an autoencoder on the attention score map with a latent dimension of 64. The model is trained for 1000 epochs with early stopping and a learning rate of 1e-3. We then apply the hierarchy clustering method fcluster from scipy ⁷ package on the normalized compressed representation using a distance threshold of 10, assigning clusters with fewer than 5 samples to a noise cluster.

B Detailed Algorithms

Algorithm 5 Search Vertical Slash Pattern

Input: $Q, K; \gamma$

Compute a subset of the full attention map $\hat{A} \leftarrow \operatorname{softmax}(\hat{Q}K^T / \sqrt{d}), \text{ where } \hat{Q} \subset Q$ # Sum and normalize attention scores along the vertical and slash directions $a_{v} \leftarrow \operatorname{sum_vertical}(\hat{A}) / \sum_{i,j} \hat{A}[I,j])$ $a_{s} \leftarrow \operatorname{sum_slash}(\hat{A}) / \sum_{i,j} \hat{A}[I,j])$ # Sort vertical and slash attention scores $I_v \leftarrow argsort(a_v)$ $I_s \leftarrow argsort(a_s)$ # Obtain the minimum computational budget making the sum of the scores exceeds γ $K_{v} \leftarrow \min \{ k: \sum_{i \in I_{v}[1:k]} a_{v}[i] \ge \gamma \}$ $K_s \leftarrow \min \{ k: \sum_{i \in I_s[1:k]} a_s[i] \ge \gamma \}$ # Select vertical and slash index $S_v \leftarrow I_v[1:K_v], S_s \leftarrow I_s[1:K_s]$ $S \leftarrow S_v \cup S_s$ $M \leftarrow index_to_mask S$ return M

C Autoencoder Architecture

782

781

Autoencoder architecture is shown in Table 3. 783

⁶https://pypi.org/project/flash-attn
⁷https://scipy.org

Encoder								
Layer	Туре	Parameters						
Conv2d	Conv2D	out: 16, kernel: 3×3 , padding: 1						
ReLU	Activation							
MaxPool2d	Pooling	kernel: 4×4 , stride: 4						
Conv2d	Conv2D	in: 16, out: 32, kernel: 3×3 , padding: 1						
ReLU	Activation							
MaxPool2d	Pooling	kernel: 4×4 , stride: 4						
Flatten	Transformation	start_dim=1, end_dim=-1						
Linear	Fully Connected	in: 468512, out: 64						
Decoder								
Layer	Туре	Parameters						
Linear	Fully Connected	in: 64, out: 468512						
ReLU	Activation							
Unflatten	Transformation	shape=(32, 121, 121)						
ConvTranspose2d	Transposed Convolution	in: 32, out: 16, kernel: 4×4 , stride: 2, padding: 1						
ReLU	Activation							
ConvTranspose2d	Transposed Convolution	in: 16, out: 8, kernel: 4 × 4, stride: 2, padding: 1						
ReLU	Activation							
ConvTranspose2d	Transposed Convolution	in: 8, out: 1, kernel: 4×4 , stride: 4						
Sigmoid	Activation							

Table 3: Network Architecture of the Autoencoder