# TokenSelect: Efficient Long-Context Inference and Length Extrapolation for LLMs via Dynamic Token-Level KV Cache Selection

Anonymous ACL submission

#### Abstract

The rapid advancement of Large Language Models (LLMs) has driven growing demand for processing extended context sequences in contemporary applications. However, this progress faces two major challenges: performance degradation due to sequence lengths outof-distribution, and excessively long inference times caused by the quadratic computational complexity of attention. These issues hinder the application of LLMs in long-context scenarios. In this paper, we propose Dynamic Token-Level KV Cache Selection (TokenSelect), a training-free method for efficient and accurate long-context inference. TokenSelect builds upon the observation of non-contiguous attention sparsity, using Query-Key dot products to measure per-head KV Cache criticality at token-level. By per-head soft voting mechanism, TokenSelect selectively involves a few critical KV cache tokens in attention calculation without sacrificing accuracy. To further accelerate TokenSelect, we design the Selection Cache based on observations of consecutive Query similarity and implemented efficient dot product kernel, significantly reducing the overhead. A comprehensive evaluation of TokenS*elect* demonstrates up to  $23.84 \times$  speedup in attention computation and up to  $2.28 \times$  acceleration in end-to-end latency, while providing superior performance compared to state-of-theart long-context inference methods.

#### 1 Introduction

003

011

014

027

031

042

With the rapid development of large language models (LLMs), the number of parameters is no longer the sole factor significantly affecting model performance. The ability to effectively process longer context information has become one of the key metrics for evaluating LLMs' capabilities. The latest applications such as cross-document understanding (Bai et al., 2024), LLM-powered search systems (Sharma et al., 2024), and complex reasoning (OpenAI) have all placed higher demands



Figure 1: Distribution of tokens participating in attention computation under different sparsity patterns (indicated by blue dots). *TokenSelect* can more accurately select critical tokens for attention computation.

on the long-context abilities of LLMs. There are two main difficulties in using pre-trained LLMs for long-context inference. On one hand, LLMs are limited by their context length during pre-training (*e.g.* Llama 3 only has 8192 tokens). Directly inferencing on longer sequences can lead to severe performance degradation due to reasons including sequence lengths out-of-distribution (Xiao et al., 2024b; Han et al., 2024). On the other hand, even if LLMs possess sufficiently large context lengths, the quadratic computational complexity of attention with respect to sequence length makes the response time for long-context inference unbearable.

Previous works have made numerous attempts to address these difficulties. To extend the context length of LLMs, the current common practice is to perform post-training on long texts (Team et al., 2024; Yang et al., 2024a; GLM et al., 2024). However, this approach comes with significant computational costs, particularly in two aspects: the synthesis of high-quality long-text data and the training process on extended sequences. To accelerate long-context inference, many studies focus on the sparsity of attention, attempting to reduce the scale of KV Cache involved in computation. The key to this type of method lies in designing sparse patterns for attention, which can be mainly divided into two categories: one uses predefined sparse patterns (Wang et al., 2019; Zaheer et al., 2020; Xiao et al., 2024b; Han et al., 2024), while the other estimates the potential importance of KV Cache

073

043

during the inference process (Zhang et al., 2024c; Oren et al., 2024; Xiao et al., 2024a; Tang et al., 2024b; Jiang et al., 2024), attempting to select relevant KV Cache tokens into attention calculations. However, the design of these sparse patterns is often heuristically based on historical criticality or coarse-grained criticality estimation of tokens, making it difficult to ensure that the selected tokens are truly critical, thus resulting in sub-optimal performance, as shown in *Fig.* 1.

075

076

079

100

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

In this paper, we further observe the noncontiguous sparsity of attention, revealing the importance of designing more fine-grained dynamic sparse patterns. To this end, we propose TokenSelect, a training-free approach that utilizes tokenlevel selective sparse attention for efficient longcontext inference and length extrapolation. Specifically, for each Query, TokenSelect dynamically calculates token-level per-head criticality for the past KV Cache and selects the k most critical tokens through our head soft vote mechanism, involving them in the attention calculation. This reduces the scale of attention calculation to a constant length familiar to the model, while maintaining almost all of the long-context information, thereby simultaneously addressing the two main difficulties for longcontext inference. To reduce the overhead of token selection, TokenSelect manages the KV Cache in token-level pages (Zheng et al., 2024) and design efficient kernel for token selection based on paged KV Cache management through Triton (Tillet et al., 2019). Furthermore, based on our observation of high similarity between consecutive queries, we have designed the Selection Cache, which allows consecutive similar queries to share token selection results, thereby reducing the selection frequency while ensuring its effectiveness.

We evaluate the performance and efficiency of *TokenSelect* on three representative long-context benchmarks using three open-source LLMs. The experimental results demonstrate that our *TokenSelect* can achieve up to  $23.84 \times$  speedup in attention computation compared to FlashInfer (flashinfer ai), and up to  $2.28 \times$  acceleration in end-to-end inference latency compared to state-of-the-art long-context inference method (Xiao et al., 2024a). Simultaneously, it provides superior performance on three long-text benchmarks. In summary, we make the following contributions:

• An observation on the non-contiguous sparsity of attention that highlights the importance of token-level KV Cache selection.

• *TokenSelect*, a training-free method that achieves accurate and efficient long-context inference and length extrapolation, which is compatible with mainstream LLM serving systems.

126

127

128

129

130

131

132

133

134

135

136

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

161

162

163

164

165

166

167

168

169

170

171

172

173

174

• Comprehensive evaluations of our method, showing up to 23.84× speedup in attention computation and up to 2.28× acceleration in end-to-end latency while exhibiting superior performance.

## 2 Preliminaries

In this section, we introduce the inference of LLMs and define the Selective Sparse Attention Problem.

## 2.1 LLMs Inference

Nowadays, mainstream LLMs are primarily based on the Decoder-only Transformer architecture. Each transformer layer includes a multi-head attention (MHA) and a feed-forward networks (FFN). The inference process of LLMs can be divided into two stages: the Prefill Stage and the Decode Stage.

The Prefill Stage is the preparatory phase of the inference process. In this stage, the user's input is processed layer by layer through a single forward pass of LLMs, generating KV Cache for each layer. The generation of KV Cache is completed by the MHA module. Assuming  $\mathbf{X}_{\text{prefill}} \in \mathbb{R}^{n_{\text{in}} \times d}$  is the input of a transformer layer, where  $n_{\text{in}}$  is the number of tokens in user's input sequence and d is the hidden size. The MHA computation in the Prefill Stage is as follows (simplified to single head):

$$\begin{bmatrix} \mathbf{Q}_{\text{prefill}}, \mathbf{K}_{\text{prefill}}, \mathbf{V}_{\text{prefill}} \end{bmatrix} = \mathbf{X}_{\text{prefill}} \cdot \begin{bmatrix} \mathbf{W}_{q}, \mathbf{W}_{k}, \mathbf{W}_{v} \end{bmatrix}, \quad (1)$$
$$\mathbf{Q}_{\text{prefill}} = \operatorname{softmax} \left( \frac{\mathbf{Q}_{\text{prefill}} \cdot \mathbf{K}_{\text{prefill}}}{2} \right) \cdot \mathbf{V}_{\text{prefill}} \quad (2)$$

$$\mathbf{O}_{\text{prefill}} = \text{softmax}\left(\frac{\mathbf{Q}_{\text{prefill}} \cdot \mathbf{K}_{\text{prefill}}}{\sqrt{d}}\right) \cdot \mathbf{V}_{\text{prefill}}, \quad (2)$$

where  $\mathbf{W}_q$ ,  $\mathbf{W}_k$ ,  $\mathbf{W}_v$  are linear projections, [·] represents tensor concatenation operation, and Eq.(2) is also known as Scaled Dot-Product Attention (SDPA). After these computation,  $\mathbf{K}_{\text{prefill}}$  and  $\mathbf{V}_{\text{prefill}}$  are stored as the KV Cache for current layer  $\mathbf{K}_{\text{cache}}$  and  $\mathbf{V}_{\text{cache}}$ , and  $\mathbf{O}_{\text{prefill}}$  is used for subsequent calculations.

The Decode Stage is the phase where LLMs actually generate the response. In the Decode Stage, LLMs load the KV Cache and generate  $n_{out}$  output tokens autoregressively through  $n_{out}$  forward passes. Assuming  $\mathbf{X}_{decode} \in \mathbb{R}^{1 \times d}$  is the input of a transformer layer in a forward pass, the computation of MHA in the Decode Stage is as follows (The calculation of  $\mathbf{Q}_{prefill}$  and  $\mathbf{O}_{prefill}$  is consistent with that in the Prefill Stage):

where  $\mathbf{K}_{decode}$ ,  $\mathbf{V}_{decode}$  are composed of the KV Cache and the KV corresponding to the current



(a) Attention is sparse in token-level. (b) Block-level selection is sub-optimal. (c) Attention logits is head-distinctive. Figure 2: Motivations for token-level selection. (a) Visualization of attention scores sparsity. (b) Attention scores and critical token recalled by 1K token budget. (c) The  $L_1$  norm of attention logits in each attention head.

input, which are then used to update the KV Cache of the current layer for use in the next forward pass.

LLMs inference, unlike training, is memorybound, necessitating frequent GPU I/O operations between HBM and SRAM while underutilizing processing units. This bottleneck is particularly evident in SDPA computation. Optimizing for I/O is crucial for enhancing LLMs inference efficiency, especially in long-context scenarios.

#### 2.2 Selective Sparse Attention

175

176

178

179

183

184

185

187

190

191

193

194

195

196

198

199

201

203

205

As discussed in the *Sec.* 1, the high attention sparsity in LLMs suggests sparse attention as a promising solution for long-context inference challenges. Sparse attention can keep the number of tokens participating in attention computations at a constant scale, rather than increasing with sequence length. Given that predefined sparse patterns are detrimental to performance, we aim to dynamically select crucial tokens at each step during the inference process. Therefore, we formalize this problem according to the following definition.

**Definition 1** (Selective Sparse Attention Problem, informal). For current input of length C (C = 1in the Decode Stage) and KV Cache of length N, assuming there are H attention heads with size of  $d_h$ , let **O** be the output of the SDPA:

$$\mathbf{O} = \left[ \sigma \left( \frac{\mathbf{Q}^{h} \cdot \left[ \mathbf{K}_{\text{cache}}^{h}, \mathbf{K}_{\text{current}}^{h} \right]^{\top}}{\sqrt{d}} \right) \cdot \left[ \mathbf{V}_{\text{cache}}^{h}, \mathbf{V}_{\text{current}}^{h} \right] \right]_{h=1}^{H},$$
(4)

where  $\sigma$  denotes softmax,  $\mathbf{Q}^h, \mathbf{K}^h_{\text{current}}, \mathbf{V}^h_{\text{current}} \in \mathbb{R}^{C \times d_h}$  are Query, Key, Value matrices of current input for head h and  $\mathbf{K}^h_{\text{cache}}, \mathbf{V}^h_{\text{cache}} \in \mathbb{R}^{N \times d_h}$  represent the KV Cache. Let  $\hat{\mathbf{O}}$  be the output of the Selective Sparse Attention:

$$\hat{\mathbf{O}} = \left[ \sigma \left( \frac{\mathbf{Q}^{h} \cdot [\mathbf{K}_{\text{select}}^{h}, \mathbf{K}_{\text{current}}^{h}]^{\top}}{\sqrt{d}} \right) \cdot [\mathbf{V}_{\text{select}}^{h}, \mathbf{V}_{\text{current}}^{h}] \right]_{h=1}^{H},$$
(5)

208 where  $\mathbf{K}_{\text{select}}^{h}, \mathbf{V}_{\text{select}}^{h} \in \mathbb{R}^{k \times d_{h}}$  are k selected KV 209 Cache ( $k \ll N$ ). The selection of  $\mathbf{K}_{\text{select}}, \mathbf{V}_{\text{select}}$  is 210 performed by selection function S:

$$\mathcal{S}(\mathbf{Q}, \mathbf{K}_{\text{cache}}) = \mathcal{I}, \text{ where } \mathcal{I} \in \mathcal{P}(\{1, \cdots, N\}),$$
  

$$\mathbf{K}_{\text{select}} = [(\mathbf{K}_{\text{cache}})_i]_{i \in \mathcal{I}}, \quad \mathbf{V}_{\text{select}} = [(\mathbf{V}_{\text{cache}})_i]_{i \in \mathcal{I}}, \quad (6)$$

where  $\mathcal{I}$  is the set of selected indices. The objective is to find an appropriate selection function S that minimizes the difference between the outputs of the SDPA and the selective sparse attention:

$$\min_{S} \left\| \mathbf{O} - \hat{\mathbf{O}} \right\|_{2}^{2}.$$
 (7)

211

212

213

214

215

216

217

218

219

220

221

224

225

226

228

230

231

232

233

234

235

236

237

239

240

241

242

243

244

245

246

247

248

Existing works on long-context inference can be categorized under the Selective Sparse Attention Problem, with variations in the design of the selection function S. Zaheer et al. (2020); Xiao et al. (2024b) have developed input-independent selection functions S(), while Zhang et al. (2024c); Oren et al. (2024); Li et al. (2024) propose Query-independent functions  $S(\mathbf{K}_{cache})$  for improved performance. Current state-of-the-art methods (Xiao et al., 2024a; Tang et al., 2024b; Jiang et al., 2024) utilize Query-aware selection functions  $S(\mathbf{Q}, \mathbf{K}_{cache})$ . However, these approaches typically operate at a block-level, which limits their effectiveness and overall performance.

#### **3** Motivations and Observations

Attention is Sparse, Non-contiguous and Head-Distinctive. Previous works on long-context inference have demonstrated the sparsity of attention scores in LLMs, particularly when processing long texts. Recent approaches (Xiao et al., 2024a; Jiang et al., 2024; Tang et al., 2024b) partition the KV Cache into non-overlapping blocks, estimating block criticality for sparse attention calculations. These methods assume that tokens with higher attention scores tend to be contiguous. However, our further observations reveal that this assumption does not always hold true in practice. As illustrated in Fig. 2a, attention scores are sparsely distributed at the token-level, with critical tokens not necessarily contiguous. This non-contiguity leads to significant omissions in block-level token selection. Fig. 2b demonstrates that finer selection granularity



(a) Consecutive queries show consistent similarity patterns across datasets. (b) Selection overlaps with similar queries. Figure 3: Observations on similarity of consecutive queries. (a) Cosine similarity distribution between consecutive queries. (b) The token selection overlap rate  $(\frac{|\mathcal{I}_i \cap \mathcal{I}_{i+1}|}{|\mathcal{I}_{i+1}|})$  with respect to consecutive Query similarity.

improves recall of critical tokens, motivating us to perform token-level selection. For token-level selection, an intuitive approach would be to directly select the top-k tokens with the highest attention logits. However, observation in *Fig.* 2c reveals considerable disparity in the  $L_1$  norm of attention logits across attention heads. As a result, the selection result tends to be dominated by a few heads with disproportionately large attention logits, driving us to design a more robust selection function that maintains the independence of heads.

250

251

252

254

258

259

260

261

262

263

267

268

270

271

272

274

275

276

277

278

279

287

**Consecutive Queries are Similar.** As sparsity of attention is dynamic (Jiang et al., 2024), token selection should be performed for every Query, which inevitably increases the computational overhead of selective sparse attention. Fortunately, we observe that consecutive Queries exhibit high similarity, as shown in *Fig.* 3a. Intuitively, when two consecutive Queries are highly similar, their dot products with the Keys will also be similar, leading to substantial overlap in the token selection results. Due to space constraints, we provide an informal lemma about this below. The formal version and corresponding proof can be found in the Appendix D.

*Lemma 1* (Informal). Consider Queries  $\mathbf{Q}_1, \mathbf{Q}_2 \in \mathbb{R}^{1 \times d}$  that are consecutive and a Key set  $\{\mathbf{K}_i\}_{i=1}^N$ . Let  $\mathcal{I}_1$ , and  $\mathcal{I}_2$  be the sets of indices of the topk Keys selected by dot product for  $\mathbf{Q}_1$ , and  $\mathbf{Q}_2$ respectively. If  $\cos(\mathbf{Q}_1, \mathbf{Q}_2) > \epsilon$ , where  $\epsilon$  is a threshold, then  $\mathcal{I}_1 = \mathcal{I}_2$ .

*Fig.* 3b illustrates this lemma experimentally. It shows that the overlap rate of token selection tends to increase with Query similarity. This key insight motivates us to reuse selection results for similar queries, improving computational efficiency. Moreover, the similarity distribution of consecutive Queries remains consistent across different tasks, as demonstrated in *Fig.* 3a, allowing us to apply a global similarity threshold across all scenarios.

### 4 Designs of TokenSelect

In this section, we will introduce the design details of *TokenSelect*, primarily encompassing the Selection Function, the Selection Cache, and efficient implementation of *TokenSelect*. The overall workflow of *TokenSelect* is illustrated in the appendix (*Fig.* 9) due to space limitations. 289

290

291

292

293

294

296

297

298

299

300

301

302

303

304

306

307

308

309

310

311

312

313

314

315

316

317

318

319

#### 4.1 Selection Function

The simplest selection function is to determine the criticality of the tokens through the dot product of  $\mathbf{Q}$  and  $\mathbf{K}_{cache}$ , then select the top-k critical ones as  $\mathbf{K}_{select}$ . The selected indices  $\mathcal{I}$  are calculated as follow:

$$\mathcal{I}_{\text{topk}} = \text{TopK}\left(\mathbf{Q} \cdot \mathbf{K}_{\text{cache}}^{h^{\top}}\right).$$
(8)

However, as discussed in Sec. 3, this approach is prone to inaccuracies due to disparities in norm of attention logits between heads. To maintain independence between heads, a better approach is to have each head select the top-k most critical tokens, and then determine the final selection through voting among the heads:

$$\mathcal{I}_{\text{head-vote}} = \text{TopK}\left(\sum_{h=1}^{H} \mathbb{I}\left(i \in \text{TopK}\left(\mathbf{Q}^{h} \cdot \mathbf{K}_{\text{cache}}^{h^{\top}}\right)\right)\right),\tag{9}$$

where  $\mathbb{I}$  is the indicator function. Unfortunately, despite better performance, this method relies on scatter\_add and multiple topk operations, resulting in low efficiency on GPUs. Additionally, the 0/1 voting ignores the relative importance of tokens for each head. Therefore, we propose a head soft vote approach that offers better performance and efficiency. Specifically, we first calculate the perhead criticality, then normalize through softmax, and sum the results for all heads:

$$\mathcal{I}_{\text{head-soft-vote}} = \text{TopK}\left(\sum_{h=1}^{H} \sigma\left(\mathbf{Q}^{h} \cdot \mathbf{K}_{\text{cache}}^{h^{\top}}\right)\right). \quad (10)$$

#### 4.2 Optimizing Selection Frequency

321

322

329

334

336

341

342

345

351

355

361

367

371

Although the aforementioned selection function can reduce the complexity of attention from  $O(N^2)$ to  $O(k^2), k \ll N$ , while maintaining performance, the execution time of the selection function itself still affects the latency of inference. To further accelerate long-context inference, based on our observations of the similarity of consecutive queries, we design optimization strategies for both the Prefill Stage and the Decode Stage to reduce the selection frequency while ensuring its effectiveness.

In the Prefill Stage,  $\mathbf{Q}_{\text{prefill}} \in \mathbb{R}^{n_{\text{in}} \times d}$  is inputed. In long-context scenarios, the number of tokens in the user's input sequence  $n_{\text{in}}$  may reach up to 1M, making it impractical to perform selection for each Query token. Considering the similarity of consecutive Queries, we use chunk-wise token selection, inputting  $\frac{1}{c} \sum_{i=1}^{c} (\mathbf{Q}_{C})_{i}$  into the selection function, where  $\mathbf{Q}_{C} \in \mathbb{R}^{c \times d}$  is the Query chunk and c is the chunk size. This method helps maintain the compute-intensive nature of the Prefill Stage, preventing it from becoming memory bound.

In the Decode Stage, due to the auto-regressive characteristic of LLMs, we need to frequently perform selection for  $Q_{decode}$ , and this process cannot be executed chunk-wise like in the Prefill Stage. To reduce the frequency of token selection in the Decode Stage, we propose the Selection Cache. Consecutive similar Queries will hit the cache, thereby directly loading the cached selection results for the previous Query. The Selection Cache allows us to reduce decode latency while maintaining the performance. The formal formulation of the Selection Cache is detailed in Appendix C (Algorithm 1).

#### 4.3 Efficient Implementation

To ensure that our proposed *TokenSelect* is ready for real-world applications, efficient implementation is crucial. We first analyze the time breakdown of representative block-level selective sparse attention method, InfLLM (Xiao et al., 2024a). From (1)(2)(3) in Fig. 4, we can observe that although selective sparse attention can significantly reduce the complexity of attention calculations, the actual computation time is still highly dependent on the implementation. The incompatibility with efficient attention implementations such as Flash Attention has resulted in methods requiring historical attention scores (Zhang et al., 2024c; Oren et al., 2024; Li et al., 2024; Xiao et al., 2024a) being difficult to be applied in real-world Web applications. Through the analysis of InfLLM's Flash Attention-



Figure 4: Time breakdown for single chunk prefill step under different attention implementations (chunk size: 512, KV Cache length: 128K, attended tokens: 4K).

372

373

374

375

376

377

378

379

381

382

383

384

386

387

388

390

391

392

393

394

395

396

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

compatible version, we make several discoveries. The initial motivation for estimating token criticality at the block-level is to reduce the overhead of selection function (mainly considering dot product calculation). However, we find that dot product is not the primary performance bottleneck. Instead, a significant portion of the overhead comes from indexing the KV Cache using selected indices and making them contiguous in GPU memory, which frequently occurs during the updating of KV blocks and the concatenation of selected KV Cache. The extensive I/O required for this operation further exacerbates the memory-bound in LLMs inference. Based on this, we propose that Paged Attention is a more suitable implementation for selective sparse attention. Using Paged KV Cache management (with page size=1 for *TokenSelect*), we can reduce the I/O volume for selection results from the scale of all selected KV Caches O(2kd) to the scale of their indices O(k). However, by observing (4) in Fig. 4, we find that we encounter another bottleneck under Paged KV Cache management. Since logically contiguous KV Cache is not entirely contiguous in GPU memory, it also needs to be made contiguous before performing computational operations. To address this issue, we draw inspiration from the concept of Paged Attention and implement a Paged Dot Product Kernel using Triton (Tillet et al., 2019), which significantly improves the overall efficiency of TokenSelect.

#### **5** Experiments

In this section, we introduce the experimental setup and evaluate the performance and efficiency of our *TokenSelect* on long-context inference benchmarks.

### 5.1 Experimental Settings

**Datasets.** To evaluate *TokenSelect*'s performance on long-context inference, we use three representative datasets: InfiniteBench (Zhang et al., 2024a), RULER (Hsieh et al., 2024), and LongBench (Bai et al., 2024). Detailed descriptions and the evaluation metrics used are provided in Appendix F.

Methods	En.Sum	En.QA	En.MC	En.Dia	Code.D	Math.F	R.PK	R.Num	R.KV	Avg.
Qwen2-7B	23.80	14.92	54.59	8.50	28.17	19.71	28.81	28.64	19.00	25.13
NTK	18.73	15.34	41.28	<b>7.50</b>	24.87	<b>27.71</b>	99.15	97.46	59.80	43.54
SelfExtend	3.76	4.44	20.09	5.00	8.12	2.29	0.00	0.00	0.00	4.86
StreamingLLM	19.60	13.61	48.03	3.50	27.92	19.43	5.08	5.08	2.40	16.07
InfLLM	19.65	15.71	46.29	<b>7.50</b>	27.41	24.00	70.34	72.20	5.40	32.06
<b>TokenSelect</b>	<b>22.62</b>	<b>18.86</b>	<b>54.31</b>	<b>7.50</b>	<b>30.20</b>	21.71	<b>100.00</b>	<b>100.00</b>	<b>86.60</b>	<b>49.08</b>
Llama-3-8B	24.70	15.50	44.10	7.50	27.92	21.70	8.50	7.80	6.20	18.21
NTK	6.40	0.40	0.00	0.00	0.50	2.60	0.00	0.00	0.00	1.10
SelfExtend	14.70	8.60	19.70	0.00	0.00	22.60	<b>100.00</b>	<b>100.00</b>	0.20	29.53
StreamingLLM	20.40	14.30	40.60	5.00	<b>28.43</b>	21.40	8.50	8.30	0.40	16.37
InfLLM	24.30	19.50	43.70	<b>10.50</b>	27.41	23.70	<b>100.00</b>	99.00	5.00	39.23
TokenSelect	<b>26.99</b>	<b>21.32</b>	<b>45.85</b>	8.00	27.41	<b>28.29</b>	<b>100.00</b>	97.29	<b>48.40</b>	<b>43.90</b>
<i>Yi-1.5-6B</i>	18.78	10.48	39.74	5.00	29.95	16.00	5.08	5.08	$\begin{array}{c} 0.00 \\ 0.00 \\ 0.00 \\ 0.00 \\ 0.00 \\ 0.00 \\ 0.00 \end{array}$	14.45
NTK	4.66	0.58	0.87	0.00	0.00	1.43	0.00	0.00		0.83
SelfExtend	5.62	1.07	1.31	0.00	0.00	1.14	0.00	0.00		1.01
StreamingLLM	15.35	9.26	35.81	5.00	27.41	14.29	5.08	4.92		13.01
InfLLM	16.98	8.93	34.06	3.00	27.41	16.86	<b>100.00</b>	96.61		33.76
<b>TokenSelect</b>	<b>21.13</b>	<b>12.32</b>	<b>40.61</b>	<b>5.50</b>	<b>30.71</b>	<b>20.86</b>	<b>100.00</b>	<b>99.83</b>		<b>36.77</b>

Table 1: Comparison of different methods with different origin models on InfiniteBench.

**Baselines.** To conduct a comprehensive eval-413 uation of TokenSelect's performance, we carry 414 415 out benchmarks on three mainstream open-source LLMs - Qwen2-7B-Instruct (Yang et al., 2024a), 416 Llama-3-8B-Instruct (Dubey et al., 2024), and 417 Yi-1.5-6B-Chat (AI et al., 2024) - comparing 418 419 against the following state-of-the-art long-context inference methods: NTK-aware scaled RoPE, Self-420 Extend, StreamingLLM, InfLLM and MInference. 421 Detailed descriptions of these methods are pro-422 vided in Appendix E. It's worth noting that since 423 MInference doesn't support length extrapolation, 424 we use an alternative evaluation method, applying 425 it to Llama-3-8B-Instruct-262k (Llama3 after 426 long-text post-training). Additionally, we do not 427 include another state-of-the-art method, QUEST, 428 as it does not support Grouped Query Attention. 429

430 **Implementation details.** In all experiments in this paper, we employ greedy decoding to ensure 431 432 the reliability of the results. For our *TokenSelect*, we implement it on SGLang (Zheng et al., 2024), 433 which is a fast serving framework based on Flasher-434 infer (flashinfer ai). We implement our method us-435 ing PyTorch (Paszke et al., 2019) and Triton (Tillet 436 et al., 2019). We follow the baseline approach, in-437 cluding 128 initial tokens and  $n_{\text{local}}$  most recent 438 tokens in the attention computation in addition to 439 the k selected tokens. For NTK and SelfExtend, 440 we extend the model's context length to 128K. For 441 StreamLLM, we set  $n_{\text{local}} = 4$ K. For InfLLM, we 442 set k = 4K,  $n_{\text{local}} = 4K$ . For our *TokenSelect*, we 443 set k = 2K,  $n_{\text{local}} = 512$  to demonstrate our token-444 445 level KV Cache selection allows us to achieve better performance with a smaller token budget. Due 446 to the need to demonstrate the method under dif-447 ferent  $n_{\text{local}}$  and k, we denote the specific token 448 budgets in the form of  $k + n_{\text{local}}$  if they differ from 449

the aforementioned settings. For InfiniteBench and LongBench, we set the threshold  $\theta$  of the Selection Cache to 0.9. We use NVIDIA A100 to conduct all experiments. When inferencing sequences over 1M tokens, we additionally employee tensor parallelism, which is transparent to our *TokenSelect*.

450

451

452

453

454

455

456

#### 5.2 Performance Comparisons

InfiniteBench. As shown in Table 1, our TokenS-457 elect achieves significantly superior overall per-458 formance on InfiniteBench compared to all base-459 line methods, even though TokenSelect uses the 460 smallest token budget (<3K). The fact that it sig-461 nificantly outperforms the original models demon-462 strates TokenSelect's strong length extrapolation 463 capability. We analyze that this is due to our adop-464 tion of a fine-grained KV Cache selection strat-465 egy, while considering the equal contribution of 466 each head to selection, which ensures that we can 467 select most critical tokens. Observing the perfor-468 mance of other methods, we find that RoPE in-469 terpolation methods (NTK, SelfExtend) generally 470 perform poorly unless used on specially trained 471 models such as Qwen2-7B-Instruct. The better 472 performance of Qwen2-7B-Instruct on the origi-473 nal model can also be attributed to this. The sparse 474 attention method StreamingLLM, based on fixed 475 sparse patterns, can guarantee some of the model's 476 capabilities, but due to discarding a large amount 477 of long-context information, it performs poorly 478 on retrieval-related tasks (R.PK, R.Num, R.KV). 479 The block-level selection method InfLLM can re-480 tain more long-context information compared to 481 StreamingLLM. However, due to its sub-optimal 482 block-level selection, it results in lower perfor-483 mance on most tasks compared to TokenSelect, 484 even though we set a larger token budget for In-485 fLLM. It is worth noting that Yi-1.5-6B does not 486

Methods	4K	8K	16K	32K	64K	128K	Avg.
Qwen2-7B	90.74	84.03	80.87	79.44	74.37	64.13	78.93
StreamingLLM	94.41	54.59	33.54	22.40	15.38	10.88	38.53
InfLLM (2K+512)	52.85	36.09	29.36	23.52	18.81	18.29	29.82
InfLLM (4K+4K)	55.22	52.10	40.53	29.77	21.56	18.64	36.30
Ours (2K+512)	94.11	81.81	68.68	60.62	51.81	42.75	66.63
Ours (4K+4K)	94.42	90.22	82.06	70.40	59.66	54.28	75.17
Llama-3-8B	93.79	90.23	0.09	0.00	0.00	0.00	30.69
StreamingLLM	93.68	54.48	33.77	20.35	14.88	11.47	38.11
InfLLM (2K+512)	79.79	52.43	40.12	33.60	25.68	23.39	42.50
InfLLM (4K+4K)	93.79	86.11	64.33	45.39	33.13	27.81	58.43
Ours (2K+512)	93.73	82.92	71.92	65.38	59.35	33.39	67.78
Ours (4K+4K)	93.88	90.29	70.13	57.72	48.36	39.38	66.63
Yi-1.5-6B	73.12	9.09	0.37	0.01	0.00	0.01	13.77
StreamingLLM	72.10	33.03	21.69	15.39	12.58	12.61	27.90
InfLLM (2K+512)	59.66	36.77	27.41	24.49	21.49	21.17	31.83
InfLLM (4K+4K)	74.81	52.57	27.65	22.83	20.19	19.48	36.26
Ours (2K+512)	75.93	59.55	49.69	42.36	34.68	31.36	48.93

Table 2: Performance comparison on RULER.

Mathada		Infi	niteBenc	h		LongBench
Wiethous	En.Sum	En.QA	Code.D	Math.F	R.KV	Avg.
Llama-3-8B-262K	20.2	12.4	22.1	26.6	14.4	33.9
+ MInference	20.5	12.9	22.3	33.1	12.8	38.4
Ours (w/ Llama-8K)	26.9	21.3	27.4	28.2	48.4	44.0

Table 3: Comparison of different methods on post-<br/>trained models on InfiniteBench and LongBench.

perform normally on the R.KV task, as it is unable to correctly recite strings like the UUID.

487

488

489

490

491

492

493

494

495

496

497

498

499

504

505

**RULER.** To further demonstrate the capability of TokenSelect, we conduct evaluation on the more challenging long-context benchmark RULER. Considering the increased difficulty of RULER and its substantial computational requirements, we include only comparable baseline methods. As shown in Table 2, our *TokenSelect* maintains significantly superior overall performance compared to other long-context inference methods. For all models, TokenSelect achieves length extrapolation while preserving the model's original capabilities, benefiting from our efficient utilization of the model's limited context length. Notably, due to the constraints of model's context length, TokenSelect experiences performance degradation with larger token budgets (4K+4K) on Llama and Yi. However, its performance with smaller token budgets still significantly surpasses other baseline methods.

LongBench. Due to space constraints, the results
of LongBench are presented in the Appendix H.
Although its relatively shorter text length makes
it less suitable for evaluating state-of-the-art longcontext inference methods, our *TokenSelect* still
demonstrates superior overall performance compared to most baseline methods.

Comparing to methods based-on post-trained
model. In Table 3, we present the performance of
the post-trained model and long-context inference
method (Jiang et al., 2024) based on it. It shows
that even compared to length extrapolation meth-

S	En.QA	En.MC	Code.D	R.PK	R.Num	R.KV
$\mathcal{I}_{topk}$	15.15	45.85	28.43	100.00	98.47	16.60
$\mathcal{I}_{\text{head-vote}}$	17.01	45.85	28.68	100.00	100.00	22.40
$\mathcal{I}_{head-soft-vote}$	18.86	54.31	30.20	100.00	100.00	86.60

Table 4: Ablation study of the Selection Function S on InfiniteBench using Qwen2-7B-Instruct.

k	En.Sum	En.QA	En.Mc	Math.F	R.Num	R.KV
128	21.23	10.46	41.48	18.00	100.00	13.40
256	22.01	11.66	41.92	19.71	100.00	20.00
512	21.60	13.31	40.17	21.71	100.00	45.60
1K	21.35	15.13	44.10	21.71	100.00	73.00
2K	22.62	18.86	54.31	21.71	100.00	86.60
4K	24.09	21.11	51.53	21.71	100.00	88.00
8K	25.32	22.93	58.52	23.71	100.00	85.40
16K	26.54	23.04	62.88	28.16	100.00	72.00

Table 5: Performance vs. Number of selected tokens k on InfiniteBench using Qwen2-7B-Instruct.

ods requiring additional training, the training-free *TokenSelect* still exhibits superior performance on most tasks. Although Minference can improve the performance of the original model, it fails to reverse the negative impact of long-text post-training on shorter text tasks (LongBench).

#### 5.3 Ablation Studies

In ablation studies, we primarily analyze the impact of different Selection Functions S on performance. To compare the performance of different Selection Functions S under low token budgets (*i.e.*, token efficiency), we maintain the 2K+512 configuration. From Table 4, we can observe that our proposed head soft vote mechanism performs significantly better across all tasks. This indicates that using the head soft vote mechanism to balance each head's contribution to token selection results can help us avoid the domination of selection by few heads with large attention logits.

#### 5.4 Hyper-parameter Analysis

Number of selected tokens k. As shown in Table 5, we fix  $n_{local}$  to a small value (512) to compare the performance when selecting different numbers of tokens. First, we observe that even selecting a very small number of tokens (*e.g.*, 128, 256), our *TokenSelect* still demonstrates very comparable performance. Then, as k increases, the effectiveness of *TokenSelect* further improves, indicating that more moderately critical tokens also contribute to the retention of long-context information. Finally, we find that when k is set to larger values (*e.g.*, 16K), our *TokenSelect* shows significant improvements in most tasks, further advancing the performance landscape of long-context inference methods.

Similarity threshold of the Selection Cache  $\theta$ . *Fig.* 5 shows that the Selection Cache hit rate in552

553

554

519

520

521

522





Figure 6: Computation time vs. KV Cache lengths for single chunk prefill step using Qwen2-7B-Instruct. The vertical axis represents the number of attended tokens. SDPA denotes full attention by Flashinfer (chunk size: 512).



Figure 7: End to end latency per sample with different methods on InfiniteBench using Qwen2-7B-Instruct.

creases significantly as the similarity threshold  $\theta$ decreases, converging around  $\theta = 0.5$ . This suggests potential for further acceleration of *TokenSelect*'s Decode Stage by reducing  $\theta$ . Performance sensitivity to  $\theta$  varies across tasks. While most tasks exhibit slight performance degradation with decreasing  $\theta$ , and R.PK in InfiniteBench shows no degradation, more challenging retrieval tasks like R.KV demonstrate significant performance deterioration. This indicates higher dynamicity requirements for token selection in these tasks.

#### 5.5 Efficiency Comparisons

555

556

557

558

561

562

567

571

573

574

**Efficiency of selective sparse attention.** *Fig.* 6 demonstrates the significant acceleration of attention computation achieved by *TokenSelect* during long-context inference. With a KV Cache length of 1M, *TokenSelect* can provide up to 23.84× speedup compared to FlashInfer, which is the inference kernel library we based on. This substantial improvement is attributed to our efficient kernel design.

575End-to-end efficiency.Fig. 7 compares the end-576to-end latency of TokenSelect, InfLLM, and SDPA577across various tasks.578celerates long-context inference in real-world sce-579narios, achieving a maximum speedup of  $4.70 \times$ 580over SDPA and  $2.28 \times$  over the state-of-the-art



Figure 8: Performance comparison on extended R.PK and R.KV using Qwen2-7B-Instruct.

long-context inference method while also delivering superior overall performance.

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

597

598

599

600

601

602

603

604

605

#### 5.6 Scaling Beyond 1 Million Context Length

To further explore *TokenSelect*'s performance in extreme long-context scenarios, we design an extended benchmark with different text lengths following InfiniteBench. As illustrated in the *Fig.* 8, our *TokenSelect* demonstrates the ability to accurately capture critical information with a small token budget in contexts up to 2M tokens, underscoring its potential in more application scenarios.

### 6 Conclusion

In this paper, we introduces *TokenSelect*, a trainingfree approach for efficient long-context inference and length extrapolation. *TokenSelect* addresses the two major challenges faced by LLMs in processing long texts: the context length limitation from pre-training and the computational complexity of attention. This is achieved through a novel tokenlevel selective sparse attention mechanism. Experimental results demonstrate that *TokenSelect* can achieve up to  $23.84 \times$  speedup in attention computation and up to  $2.28 \times$  acceleration in end-to-end inference latency, while exhibiting superior performance across multiple long-context benchmarks.

#### 7 Limitations

606

Our approach has inherent limitations that present opportunities for future work. A primary limitation of our method is that its training-free design-a significant advantage-acts as a doubleedged sword, as its absolute performance is inher-611 612 ently tied to the quality of the underlying LLMs. Although our experiments demonstrate robustness 613 of TokenSelect across various LLMs, some inher-614 ent shortcomings—such as the misrecognition of UUID strings by Yi-1.5-6B-Chat—indicate that 616 certain issues may still require training to resolve. 617 Moreover, while our method currently achieves 618 state-of-the-art performance in long-context inference, recent long-text post-training techniques in the LLM community have shown impressive per-621 formance; notably, our TokenSelect is orthogonal to these approaches and can be employed dur-623 ing inference to trade a slight performance drop for significant efficiency gains. Finally, although our method achieves state-of-the-art efficiency improvements in long-context inference, the task remains inherently resource-intensive. For instance, even with a 8B-parameter model, complex bench-629 marks (e.g., RULER) can require approximately 8×A100 GPUs for nearly one day of runtime, and the computational cost is expected to increase substantially for larger models. We hope that our work, 633 together with the community's advances in model 634 design, algorithm development, and infrastructure 635 optimization, will help pave the way for further mitigating these computational challenges. 637

#### References

641

644

646

647

648

649

653

- 01. AI, :, Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, Kaidong Yu, Peng Liu, Qiang Liu, Shawn Yue, Senbin Yang, Shiming Yang, Tao Yu, Wen Xie, Wenhao Huang, Xiaohui Hu, Xiaoyi Ren, Xinyao Niu, Pengcheng Nie, Yuchi Xu, Yudong Liu, Yue Wang, Yuxuan Cai, Zhenyu Gu, Zhiyuan Liu, and Zonghong Dai. 2024. Yi: Open foundation models by 01.ai. *Preprint*, arXiv:2403.04652.
- Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao Liu, Aohan Zeng, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. 2024. LongBench: A bilingual, multitask benchmark for long context understanding. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 3119–3137, Bangkok, Thailand. Association for Computational Linguistics.

bloc97. 2023. Ntk-aware scaled rope allows llama models to have extended (8k+) context size without any fine-tuning and minimal perplexity degradation. Website. https://www.reddit.com/ r/LocalLLaMA/comments/14lz7j5/ntkaware\_ scaled\_rope\_allows\_llama\_models\_to\_have/.

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

- Aydar Bulatov, Yury Kuratov, and Mikhail Burtsev. 2022. Recurrent memory transformer. *Advances in Neural Information Processing Systems*, 35:11079– 11091.
- Shouyuan Chen, Sherman Wong, Liangjian Chen, and Yuandong Tian. 2023. Extending context window of large language models via positional interpolation. *Preprint*, arXiv:2306.15595.
- Zihang Dai\*, Zhilin Yang\*, Yiming Yang, William W. Cohen, Jaime Carbonell, Quoc V. Le, and Ruslan Salakhutdinov. 2019. Transformer-XL: Language modeling with longer-term dependency.
- Tri Dao. 2024. FlashAttention-2: Faster attention with better parallelism and work partitioning. In *International Conference on Learning Representations* (*ICLR*).
- 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).
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, et al. 2024. The llama 3 herd of models. *Preprint*, arXiv:2407.21783.
- emozilla. 2023. Dynamically scaled rope further increases performance of long context llama with zero fine-tuning. Website. https://www.reddit.com/r/LocalLLaMA/ comments/14mrgpr/dynamically\_scaled\_rope\_ further\_increases/.
- flashinfer ai. GitHub flashinfer-ai/flashinfer: FlashInfer: Kernel Library for LLM Serving — github.com. https://github.com/ flashinfer-ai/flashinfer. [Accessed 12-10-2024].
- Suyu Ge, Yunan Zhang, Liyuan Liu, Minjia Zhang, Jiawei Han, and Jianfeng Gao. 2024. Model tells you what to discard: Adaptive KV cache compression for LLMs. In *The Twelfth International Conference on Learning Representations*.
- Team GLM, :, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, et al. 2024. Chatglm: A family of large language models from glm-130b to glm-4 all tools. *Preprint*, arXiv:2406.12793.
- Chi Han, Qifan Wang, Hao Peng, Wenhan Xiong, Yu Chen, Heng Ji, and Sinong Wang. 2024. Lminfinite: Zero-shot extreme length generalization for large language models. *Preprint*, arXiv:2308.16137.

822

823

824

Yefei He, Luoming Zhang, Weijia Wu, Jing Liu, Hong Zhou, and Bohan Zhuang. 2024. Zipcache: Accurate and efficient kv cache quantization with salient token identification. *Preprint*, arXiv:2405.14256.

712

713

714

717

721

722

727

728

730

731

732

734

736

739

740

741

742

743

744

745

746

747

751

754

755

756

759

761

763

- Cheng-Ping Hsieh, Simeng Sun, Samuel Kriman, Shantanu Acharya, Dima Rekesh, Fei Jia, Yang Zhang, and Boris Ginsburg. 2024. Ruler: What's the real context size of your long-context language models? *arXiv preprint arXiv:2404.06654*.
  - Yunpeng Huang, Jingwei Xu, Junyu Lai, Zixu Jiang, Taolue Chen, Zenan Li, Yuan Yao, Xiaoxing Ma, Lijuan Yang, Hao Chen, Shupeng Li, and Penghao Zhao. 2024. Advancing transformer architecture in long-context large language models: A comprehensive survey. *Preprint*, arXiv:2311.12351.
- Huggingface. 2024. Huggingface text generation inference. Website. https://github.com/ huggingface/text-generation-inference.
- Sam Ade Jacobs, Masahiro Tanaka, Chengming Zhang, Minjia Zhang, Shuaiwen Leon Song, Samyam Rajbhandari, and Yuxiong He. 2023. Deepspeed ulysses: System optimizations for enabling training of extreme long sequence transformer models. *Preprint*, arXiv:2309.14509.
- Arthur Jacot, Franck Gabriel, and Clément Hongler. 2018. Neural tangent kernel: Convergence and generalization in neural networks. *Advances in neural information processing systems*, 31.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b. Preprint, arXiv:2310.06825.
- Huiqiang Jiang, Yucheng Li, Chengruidong Zhang, Qianhui Wu, Xufang Luo, Surin Ahn, Zhenhua Han, Amir H. Abdi, Dongsheng Li, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. 2024. Minference 1.0: Accelerating pre-filling for long-context llms via dynamic sparse attention. *Preprint*, arXiv:2407.02490.
- Hao Kang, Qingru Zhang, Souvik Kundu, Geonhwa Jeong, Zaoxing Liu, Tushar Krishna, and Tuo Zhao.
  2024. Gear: An efficient kv cache compression recipe for near-lossless generative inference of llm. *Preprint*, arXiv:2403.05527.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. *Preprint*, arXiv:2309.06180.
  - Wonbeom Lee, Jungi Lee, Junghwan Seo, and Jaewoong Sim. 2024. Infinigen: Efficient generative inference of large language models with dynamic kv cache management. *Preprint*, arXiv:2406.19707.

- 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. *Preprint*, arXiv:2404.14469.
- Akide Liu, Jing Liu, Zizheng Pan, Yefei He, Gholamreza Haffari, and Bohan Zhuang. 2024a. Minicache: Kv cache compression in depth dimension for large language models. *Preprint*, arXiv:2405.14366.
- Hao Liu, Matei Zaharia, and Pieter Abbeel. 2024b. Ringattention with blockwise transformers for nearinfinite context. In *The Twelfth International Conference on Learning Representations*.
- Zichang Liu, Aditya Desai, Fangshuo Liao, Weitao Wang, Victor Xie, Zhaozhuo Xu, Anastasios Kyrillidis, and Anshumali Shrivastava. 2023. Scissorhands: Exploiting the persistence of importance hypothesis for llm kv cache compression at test time. In Advances in Neural Information Processing Systems, volume 36, pages 52342–52364. Curran Associates, Inc.
- Zirui Liu, Jiayi Yuan, Hongye Jin, Shaochen Zhong, Zhaozhuo Xu, Vladimir Braverman, Beidi Chen, and Xia Hu. 2024c. KIVI: A tuning-free asymmetric 2bit quantization for KV cache. In *Forty-first International Conference on Machine Learning*.
- Tsendsuren Munkhdalai, Manaal Faruqui, and Siddharth Gopal. 2024. Leave no context behind: Efficient infinite context transformers with infiniattention. *Preprint*, arXiv:2404.07143.
- NVIDIA. 2024. Tensorrt-llm. Website. https://github.com/NVIDIA/TensorRT-LLM.
- OpenAI. Introducing OpenAI o1. https://openai. com/o1/. [Accessed 06-10-2024].
- Matanel Oren, Michael Hassid, Nir Yarden, Yossi Adi, and Roy Schwartz. 2024. Transformers are multistate rnns. *Preprint*, arXiv:2401.06104.
- Arka Pal, Deep Karkhanis, Manley Roberts, Samuel Dooley, Arvind Sundararajan, and Siddartha Naidu.
  2023. Giraffe: Adventures in expanding context lengths in llms. *Preprint*, arXiv:2308.10882.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. Pytorch: An imperative style, high-performance deep learning library. In Advances in Neural Information Processing Systems, volume 32. Curran Associates, Inc.
- Bowen Peng, Jeffrey Quesnelle, Honglu Fan, and Enrico Shippole. 2024. YaRN: Efficient context window extension of large language models. In *The Twelfth International Conference on Learning Representations*.

- 825 826
- oz 82
- 02
- 830
- 831 832
- 833
- 835
- 836 837
- 838 839
- 840 841
- 8
- 843 844
- 844 845
- 84
- 848 849
- 0
- 850 851
- 852 853
- 85
- 855 856 857
- 8

8

- 0
- 8
- 86
- 867
- 8
- 870 871

872 873

874

875 876

877

- Jack W. Rae, Anna Potapenko, Siddhant M. Jayakumar, Chloe Hillier, and Timothy P. Lillicrap. 2020. Compressive transformers for long-range sequence modelling. In *International Conference on Learning Representations*. *Compressive transformers for long-range sequence* modelling. In *International Conference on Learning*
- Luka Ribar, Ivan Chelombiev, Luke Hudlass-Galley, Charlie Blake, Carlo Luschi, and Douglas Orr. 2024. Sparq attention: Bandwidth-efficient LLM inference. In *Forty-first International Conference on Machine Learning*.
- Nikhil Sharma, Q. Vera Liao, and Ziang Xiao. 2024. Generative echo chamber? effect of llm-powered search systems on diverse information seeking. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, CHI '24, New York, NY, USA. Association for Computing Machinery.
  - Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan Catanzaro. 2020. Megatron-lm: Training multi-billion parameter language models using model parallelism. *Preprint*, arXiv:1909.08053.
- Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. 2024. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063.
- Hanlin Tang, Yang Lin, Jing Lin, Qingsen Han, Shikuan Hong, Yiwu Yao, and Gongyi Wang. 2024a. Razorattention: Efficient kv cache compression through retrieval heads. *Preprint*, arXiv:2407.15891.
- Jiaming Tang, Yilong Zhao, Kan Zhu, Guangxuan Xiao, Baris Kasikci, and Song Han. 2024b. QUEST: Query-aware sparsity for efficient long-context LLM inference. In *Forty-first International Conference on Machine Learning*.
- Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *Preprint*, arXiv:2403.05530.
- Junfeng Tian, Da Zheng, Yang Cheng, Rui Wang, Colin Zhang, and Debing Zhang. 2024. Untie the knots: An efficient data augmentation strategy for longcontext pre-training in language models. *Preprint*, arXiv:2409.04774.
- Philippe Tillet, H. T. 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, MAPL 2019, page 10–19, New York, NY, USA. Association for Computing Machinery.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *Preprint*, arXiv:2307.09288.

Zhongwei Wan, Xinjian Wu, Yu Zhang, Yi Xin, Chaofan Tao, Zhihong Zhu, Xin Wang, Siqi Luo, Jing Xiong, and Mi Zhang. 2024. D2o: Dynamic discriminative operations for efficient generative inference of large language models. *Preprint*, arXiv:2406.13035. 878

879

881

882

883

884

885

886

887

888

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

- Zheng Wang, Boxiao Jin, Zhongzhi Yu, and Minjia Zhang. 2024. Model tells you where to merge: Adaptive kv cache merging for llms on long-context tasks. *Preprint*, arXiv:2407.08454.
- Zhiguo Wang, Patrick Ng, Xiaofei Ma, Ramesh Nallapati, and Bing Xiang. 2019. Multi-passage bert: A globally normalized bert model for open-domain question answering. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5878–5882.
- Chaojun Xiao, Pengle Zhang, Xu Han, Guangxuan Xiao, Yankai Lin, Zhengyan Zhang, Zhiyuan Liu, and Maosong Sun. 2024a. Infllm: Training-free long-context extrapolation for llms with an efficient context memory. *Preprint*, arXiv:2402.04617.
- Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. 2024b. Efficient streaming language models with attention sinks. In *The Twelfth International Conference on Learning Representations*.
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, et al. 2024a. Qwen2 technical report. *Preprint*, arXiv:2407.10671.
- June Yong Yang, Byeongwook Kim, Jeongin Bae, Beomseok Kwon, Gunho Park, Eunho Yang, Se Jung Kwon, and Dongsoo Lee. 2024b. No token left behind: Reliable kv cache compression via importanceaware mixed precision quantization. *Preprint*, arXiv:2402.18096.
- Shuo Yang, Ying Sheng, Joseph E. Gonzalez, Ion Stoica, and Lianmin Zheng. 2024c. Post-training sparse attention with double sparsity. *Preprint*, arXiv:2408.07092.
- Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, and Amr Ahmed. 2020. Big bird: Transformers for longer sequences. In *Advances in Neural Information Processing Systems*, volume 33, pages 17283–17297. Curran Associates, Inc.
- Xinrong Zhang, Yingfa Chen, Shengding Hu, Zihang Xu, Junhao Chen, Moo Hao, Xu Han, Zhen Thai, Shuo Wang, Zhiyuan Liu, and Maosong Sun. 2024a. ∞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.

937 941

934

935

- 942 943
- 946
- 947 949 951
- 953 954
- 960 961

962 963

957

944

ference of large language models. Advances in Neural Information Processing Systems, 36.

58840-58850. PMLR.

Liang Zhao, Xiaocheng Feng, Xiachong Feng, Dongliang Xu, Qing Yang, Hongtao Liu, Bing Qin, and Ting Liu. 2024. Length extrapolation of transformers: A survey from the perspective of positional encoding. Preprint, arXiv:2312.17044.

Yuxin Zhang, Yuxuan Du, Gen Luo, Yunshan Zhong,

Zhenyu Zhang, Shiwei Liu, and Rongrong Ji. 2024b.

CaM: Cache merging for memory-efficient LLMs

inference. In Proceedings of the 41st International

Conference on Machine Learning, volume 235 of

Proceedings of Machine Learning Research, pages

Zhenyu Zhang, Ying Sheng, Tianyi Zhou, Tianlong

Chen, Lianmin Zheng, Ruisi Cai, Zhao Song, Yuan-

dong Tian, Christopher Ré, Clark Barrett, et al. 2024c.

H2o: Heavy-hitter oracle for efficient generative in-

- Lianmin Zheng, Liangsheng Yin, Zhiqiang Xie, Chuyue Sun, Jeff Huang, Cody Hao Yu, Shiyi Cao, Christos Kozyrakis, Ion Stoica, Joseph E. Gonzalez, Clark Barrett, and Ying Sheng. 2024. Sglang: Efficient execution of structured language model programs. Preprint, arXiv:2312.07104.
- Zixuan Zhou, Xuefei Ning, Ke Hong, Tianyu Fu, Jiaming Xu, Shiyao Li, Yuming Lou, Luning Wang, Zhihang Yuan, Xiuhong Li, Shengen Yan, Guohao Dai, Xiao-Ping Zhang, Yuhan Dong, and Yu Wang. 2024. A survey on efficient inference for large language models. Preprint, arXiv:2404.14294.

#### A **Related Works**

Long-context LLMs Due to computational complexity constraints, current LLMs based on Transformers often utilize limited context lengths during pre-training (Touvron et al., 2023; Dubey et al., 2024; Jiang et al., 2023; Yang et al., 2024a; GLM et al., 2024; AI et al., 2024). To extend the long-context capabilities of LLMs, current methods can be broadly categorized into three approaches (Huang et al., 2024; Zhou et al., 2024; Zhao et al., 2024): 1) Modifying positional encodings: A widely adopted method is positional interpolation (Chen et al., 2023). Chen et al. first proposed linear scaling of RoPE (Su et al., 2024) to map longer positional ranges within the original training window. Subsequent works (bloc97, 2023; emozilla, 2023) further improved this method using Neural Tangent Kernel (NTK) theory (Jacot et al., 2018), achieving longer context windows while maintaining model performance. Methods like YaRN (Peng et al., 2024) and Giraffe (Pal et al., 2023) optimize interpolation effects by adjusting frequency components or introducing temperature parameters. 2) Long-context post-training: This approach extends the model's context length through additional training steps on longer documents after pre-training (Yang et al., 2024c; Tian et al., 2024). It has been widely adopted by leading LLMs (Team et al., 2024; Yang et al., 2024a; GLM et al., 2024) with the support of sequence parallelism techniques (Shoeybi et al., 2020; Jacobs et al., 2023; Liu et al., 2024b). 3) Incorporating additional memory modules: Notable examples include Transformer-XL (Dai\* et al., 2019), Compressive Transformer (Rae et al., 2020), RMT (Bulatov et al., 2022) and Infini-attention (Munkhdalai et al., 2024). Although these methods have ex-1000 panded the context length of LLMs, long-context 1001 inference still faces the challenge of high computa-1002 tional costs. 1003

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1004

1005

1006

1007

1008

1009

1010

1011

1012

1013

1014

Efficient Long-context Inference In state-ofthe-art LLMs serving systems (Kwon et al., 2023; Huggingface, 2024; NVIDIA, 2024; Zheng et al., 2024), technologies such as Flash Attention (Dao et al., 2022; Dao, 2024) and Paged Attention (Kwon et al., 2023) have greatly optimized LLMs inference efficiency by improving GPU I/O bottlenecks. However, in long-context inference scenarios, the quadratic computational complexity of attention with respect to sequence length poses new challenges for LLMs inference. Numerous



Figure 9: The illustration of *TokenSelect*, which involves calculating per-head criticality using the Paged Dot Product Kernel, performing head soft vote to get selected indices, and executing selective sparse attention via the Paged Attention Kernel.

studies focus on the sparsity of attention, select-1015 ing partial KV Cache for attention calculations to 1016 improve long-context inference efficiency. Sliding window (Wang et al., 2019; Zaheer et al., 2020) is one of the most widely used sparse patterns, 1019 reducing complexity to linear by executing atten-1020 tion computations within localized windows. Re-1021 cent works like StreamingLLM (Xiao et al., 2024b) 1022 and LM-infinite (Han et al., 2024) retain the ini-1023 tial tokens of the sequence in addition to sliding 1024 windows, effectively maintaining LLMs' perfor-1025 mance when processing long sequences. While 1026 these approaches are simple to implement, they 1027 cannot retain information from long contexts. An-1028 other approach focuses on dynamic KV Cache selection during inference. Methods like H2O (Zhang et al., 2024c), TOVA (Oren et al., 2024), Fast-Gen (Ge et al., 2024), Scissorhands (Liu et al., 1032 2023), and SnapKV (Li et al., 2024) evaluate to-1033 ken criticality based on historical attention scores, 1034 selecting tokens within a limited budget. However, these methods permanently discard parts of 1036 the KV Cache, causing information loss from long 1037 contexts. To address this, InfLLM (Xiao et al., 1038 2024a) introduces Block Memory Units for KV Cache management, retrieving information from 1040 long contexts and offloading less-used blocks to 1041 CPU. Similarly, QUEST (Tang et al., 2024b) pro-1042

poses query-aware sparsity at page granularity, 1043 while MInference (Jiang et al., 2024) optimizes 1044 long-context inference using three sparse patterns. 1045 Apart from considering all attention heads, some other works (Ribar et al., 2024; Lee et al., 2024; 1047 Tang et al., 2024a) attempt to focus on only a sub-1048 set of attention heads. Beyond selection, some other research focuses on KV Cache quantiza-1050 tion (Liu et al., 2024c; Yang et al., 2024b; He et al., 1051 2024; Kang et al., 2024) and merging (Liu et al., 1052 2024a; Wan et al., 2024; Zhang et al., 2024b; Wang 1053 et al., 2024). While existing methods have shown 1054 progress, opportunities for further improvement remain in achieving optimal accuracy and compu-1056 tational efficiency for real-world deployment.

#### **B** The Illustration of *TokenSelect*

The workflow of *TokenSelect* are illustrated in *Fig.* 9.

1059

#### C The Selection Cache Algorithm

Algorithm I Selection Cache Algor	unm
-----------------------------------	-----

<b>Require:</b> Q: current Query, k: number of selected tokens,
$\mathbf{C}_Q$ : Query cache, $\mathbf{C}_{\mathcal{I}}$ : selection cache,
$\theta$ : similarity threshold,
S: selection function (Eq.(10)),
<b>f</b> : first query flags (default True)
<b>Ensure:</b> $\mathcal{I}$ : selected indices
1: if f or $\cos(\mathbf{Q}, \mathbf{C}_Q) < \theta$ then
2: $\mathcal{I} \leftarrow \mathcal{S}(\mathbf{Q}, k)$
3: $\mathbf{C}_{\mathcal{I}} \leftarrow \mathcal{I}$
4: $\mathbf{C}_Q \leftarrow \mathbf{Q}$
5: $\mathbf{f} \leftarrow \text{False}$
6: else
7: $\mathcal{I} \leftarrow C_{\mathcal{I}}$
8: end if
9. return $\mathcal{T}$

# D Formal Statement and Proof of Lemma

Lemma 1 (Invariant Top-k Key Selection under Cosine Similarity Threshold, Formal). Assumptions:

- 1. Let  $\mathbf{q}_1, \mathbf{q}_2 \in \mathbb{R}^d$  be two query vectors.
- 2. Let  $\{\mathbf{k}_i\}_{i=1}^N \subset \mathbb{R}^d$  be a finite set of key vectors.
- 3. Let k be a positive integer such that  $1 \le k \le N$ .
- 4. Define the cosine similarity between vectors  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^d$  as:

$$\cos(\mathbf{a}, \mathbf{b}) = rac{\mathbf{a}\mathbf{b}}{\|\mathbf{a}\|_2 \|\mathbf{b}\|_2},$$

where  $\|\cdot\|_2$  denotes the Euclidean norm.

- 5. Define the top-k selection function based on dot product similarity as:  $\mathcal{I}(\mathbf{q}) = \arg \max_{S \subseteq \{1,2,...,N\}, |S|=k} \sum_{i \in S} \mathbf{q} \cdot \mathbf{k}_i$ . Assume that for any query vectors  $\mathbf{q}$ , the top-k set  $\mathcal{I}(\mathbf{q})$  is uniquely determined.
- 6. Let  $\epsilon \in (0, 1]$  be a predefined threshold.

1081Lemma Statement: If the cosine similarity be-1082tween the two query vectors  $\mathbf{q}_1$  and  $\mathbf{q}_2$  satisfies

$$\cos(\mathbf{q}_1,\mathbf{q}_2) > \epsilon,$$

1084then the indices of the top-k keys selected by  $q_1$ 1085and  $q_2$  are identical, i.e.,

1086 
$$\mathcal{I}(\mathbf{q}_1) = \mathcal{I}(\mathbf{q}_2)$$

**Proof:** We start with the given condition:

$$\min_{1 \le i \le k} \mathbf{q}_1 \mathbf{k}_i - \max_{j > k} \mathbf{q}_1 \mathbf{k}_j > \eta,$$
 1080

1087

1089

1094

1095

1097

1104

1108

1112

which we aim to use to demonstrate that:

$$\min_{1 \le i \le k} \mathbf{q}_2 \mathbf{k}_i - \max_{j > k} \mathbf{q}_2 \mathbf{k}_j > 0.$$

To facilitate our analysis, we introduce the follow-1091ing notations:1092

$$\hat{\eta} = \frac{\eta}{\|\mathbf{q}_1\|}, \quad \hat{\mathbf{q}}_1 = \frac{\mathbf{q}_1}{\|\mathbf{q}_1\|}, \quad \hat{\mathbf{q}}_2 = \frac{\mathbf{q}_2}{\|\mathbf{q}_2\|}.$$
 1093

With these definitions, the original condition becomes:

$$\min_{1 \le i \le k} \hat{\mathbf{q}}_1 \mathbf{k}_i - \max_{j > k} \hat{\mathbf{q}}_1 \mathbf{k}_j > \hat{\eta},$$
 1096

and our goal transforms to showing:

$$\min_{1 \le i \le k} \hat{\mathbf{q}}_2 \mathbf{k}_i - \max_{j > k} \hat{\mathbf{q}}_2 \mathbf{k}_j > 0.$$
 1090

Next, let  $\theta$  denote the angle between  $\mathbf{q}_1$  and  $\mathbf{q}_2$ , 1099  $\cos \theta = \hat{\mathbf{q}}_1 \cdot \hat{\mathbf{q}}_2$ . We can further define: 1100

$$\mathbf{p}_1 = \mathbf{q}_2 - \mathbf{q}_1 \cos \theta, \quad \hat{\mathbf{p}}_1 = \frac{\mathbf{p}_1}{\|\mathbf{p}_1\|},$$
 110

then 
$$\sin \theta = \hat{\mathbf{p}}_1 \cdot \hat{\mathbf{q}}_2$$
, and 110

$$\hat{\mathbf{q}}_2 = \hat{\mathbf{q}}_1 \cos \theta + \hat{\mathbf{p}}_1 \sin \theta.$$
 1103

Then we have:

$$\geq \hat{\mathbf{q}}_1 \mathbf{k}_k \cos \theta - \|\mathbf{k}\|_{\max} \sin \theta, \qquad 110$$

and

$$\max_{j>k} \hat{\mathbf{q}}_2 \mathbf{k}_j = \max_{j>k} \left( \hat{\mathbf{q}}_1 \cos \theta + \hat{\mathbf{p}}_1 \sin \theta \right) \mathbf{k}_j$$
 1109

$$\leq \max_{j>k} \dot{\mathbf{q}}_1 \mathbf{k}_i \cos \theta + \max_{j>k} \dot{\mathbf{p}}_1 \mathbf{k}_i \sin \theta, \qquad 1110$$

$$\leq \hat{\mathbf{q}}_1 \mathbf{k}_{p+1} \cos \theta + \|\mathbf{k}\|_{\max} \sin \theta.$$
 111

Therefore,

 $\min_{1 \le i \le k} \hat{\mathbf{q}}_2 \mathbf{k}_i - \max_{j > k} \hat{\mathbf{q}}_2 \mathbf{k}_j \ge \hat{\mathbf{q}}_1 \mathbf{k}_p \cos \theta - \|\mathbf{k}\|_{\max} \sin \theta$  1113

$$-\left(\hat{\mathbf{q}}_{1}\mathbf{k}_{p+1}\cos\theta + \|\mathbf{k}\|_{\max}\sin\theta\right) \qquad 1114$$

$$= (\hat{\mathbf{q}}_1 \mathbf{k}_p \cos \theta - \hat{\mathbf{q}}_1 \mathbf{k}_{p+1} \cos \theta)$$
 1115

$$-2\|\mathbf{k}\|_{\max}\sin\theta$$
 1116

$$\geq \hat{\eta}\cos\theta - 2\|\mathbf{k}\|_{\max}\sin\theta. \tag{11}$$

1062

1063

1064

1066

1067

1068

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

1080

In order to have Eqn. (11) > 0, we require 1118

1119 
$$\hat{\eta}\cos\theta > 2\|\mathbf{k}\|_{\max}\sin\theta$$

1121

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

1137

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

$$\Rightarrow \frac{\sin \theta}{\cos \theta} < \frac{\hat{\eta}}{2 \|\mathbf{k}\|_{\max}}, \\ \Rightarrow \frac{1 - \cos^2 \theta}{\cos^2 \theta} < \left(\frac{\hat{\eta}}{2 \|\mathbf{k}\|_{\max}}\right)^2,$$

$$\Rightarrow \cos\theta \ge \frac{1}{\sqrt{1 + \left(\frac{\hat{\eta}}{2\|\mathbf{k}\|_{\max}}\right)^2}}.$$

This final inequality establishes a sufficient con-1123 dition for the original statement to hold, thereby 1124 completing the proof. 1125

#### E **Detailed Descriptions on Baselines** 1126

In this paper, we use the following baselines:

- NTK-Aware Scaled RoPE (bloc97, 2023): A nonlinear RoPE interpolation method.
- SelfExtend: A RoPE interpolation method that reuses the position ids of neighboring tokens.
- StreamingLLM (Xiao et al., 2024b): The state-of-the-ar method for long-context inference with predefined sparse patterns. Similar approaches include LM-Infinite (Han et al., 2024).
- InfLLM (Xiao et al., 2024a): The state-ofthe-art method for long-context inference and length extrapolation using a block-level selective sparse attention method.
  - MInference (Jiang et al., 2024): The state-ofthe-ar method for long-context prefilling acceleration, utilizing three sparse patterns including block-level sparse attention.

#### More Information on Datasets F

In this paper, we use the following datasets:

- InfiniteBench (Zhang et al., 2024a): The main-٠ stream long-context benchmark consisting of multi-tasks. The average length of it exceeds 200K tokens.
- RULER (Hsieh et al., 2024): A challenging 1151 long-context benchmark containing 13 differ-1152 ent tasks, with subsets of varying lengths up to 1153 128K tokens. 1154

LongBench (Bai et al., 2024): Another main-1155 stream long-context benchmark comprising 6 1156 types of tasks. The 95% percentile for its 1157 lengths is 31K tokens. 1158

For InfiniteBench (Zhang et al., 2024a), we use 1159 longbook\_sum\_eng (En.Sum), longbook\_qa\_eng 1160 (En.QA), longbook\_choice\_eng (En.MC), longdi-1161 alogue qa eng (En.Dia), code debug (Code.D), 1162 math find (Math.F), passkey (R.PK), num-1163 ber\_string (R.Num) and kv\_retrieval (R.KV) as 1164 evaluation datasets. The corresponding evaluation 1165 metrics are shown in Table 7. RULER (Hsieh et al., 1166 2024) consists of various evaluation tasks: Single 1167 NIAH (needle in a haystack), Multi-keys NIAH, 1168 Multi-values NIAH, Multi-values NIAH, Multi-1169 queries NIAH, Variable Tracking, Common Words 1170 Extraction, Frequent Words Extraction and Ques-1171 tion Answering. The evaluation metric is match 1172 rate. For LongBench, we use all English tasks with 1173 evaluation metrics in Table 8. 1174

#### **Effiency Comparison with MInference** G

We note that Minference (Jiang et al., 2024) has gained widespread adoption in real-world longcontext inference applications due to its novel design of attention sparse patterns and efficient implementation based on vLLM. In the main text, we demonstrated TokenSelect's performance advantages. To further prove its efficiency readiness for real-world applications, we followed Minference's approach by comparing the end-to-end prefill latency under paged KV Cache management for different input token lengths on Llama-3-8B using a single A100, with results shown in Table 6. The results indicate that TokenSelect demonstrates significant advantages with shorter input token lengths, while maintaining efficiency comparable to MInference as input token lengths increase.

Length	FlashAttention-2 (vLLM)	MInference (vLLM)	TokenSelect
1K	0.081	3.017	0.092
10K	0.832	2.762	1.290
50K	7.717	7.540	5.712
100K	21.731	14.081	12.088
128K	32.863	18.827	15.920
200K	OOM	OOM	26.500
300K	OOM	OOM	43.406

Table 6: Comparison of end-to-end prefill latency (s).

1175

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

1188

1189

### 1192 H Experimental Results on LongBench

Compared to InfiniteBench and RULER, Long-1193 Bench has much shorter text lengths. The 95% per-1194 centile for its lengths is 31K tokens. Considering 1195 that recent LLMs after SFT generally have context 1196 lengths of up to 32K tokens (Yang et al., 2024a), 1197 LongBench is less suitable for evaluating state-of-1198 the-art long-context inference methods. Neverthe-1199 less, as shown in Table 9, our TokenSelect still 1200 demonstrates superior overall performance com-1201 pared to most baseline methods. It's worth noting 1202 that Yi-1.5-6B did not yield effective results on 1203 the SAMSum task because it failed to correctly 1204 follow instructions. 1205

Datasets	En.Sum	En.QA	En.MC	En.Dia	Code.D	Math.F	R.PK	R.Num	R.KV
Metrics	Rouge-L-Sum	QA F1 Score	Accuracy						

Table 7: Evaluation metrics of different datasets on InfiniteBench.

Datasets	NQA	Qasper	MFQA	HQA	2WikiMQA	Musique	GovReport	QMSum
Metrics	QA F1 Score	Rouge-L	Rouge-L					
Datasets	MultiNews	TREC	TQA	SAMSum	PsgCount	PsgRetrieval	LCC	RepoBench-P
Metrics	Rouge-L	Accuracy	QA F1 Score	Rouge-L	Accuracy	Accuracy	Code Sim Score	Code Sim Score

Table 8: Evaluation metrics of different datasets on LongBench.

Methods	NQA	Qasper	MFQA	HQA	2WikiMQA	Musique	GovReport	QMSum	MultiNews
Qwen2-7B	24.24	45.42	47.79	42.76	44.38	24.16	33.80	23.78	26.17
NTK	26.25	45.94	50.76	53.20	50.31	30.83	32.75	23.21	25.94
SelfExtend	7.15	20.37	24.06	14.91	13.73	4.75	16.92	16.53	18.74
StreamLLM	19.49	42.56	39.63	42.43	44.67	15.22	31.51	20.57	26.00
InfLLM	27.47	41.44	46.99	47.47	49.29	25.62	32.68	23.10	26.77
TokenSelect	24.18	42.29	45.77	48.62	49.08	27.85	33.69	23.03	26.35
Llama-3-8B	19.85	42.36	41.03	47.38	39.20	22.96	29.94	21.45	27.51
NTK	9.90	45.35	49.41	48.86	29.22	24.56	34.31	23.82	27.27
SelfExtend	1.72	8.90	20.80	8.65	6.97	3.27	13.99	15.36	17.66
StreamLLM	20.05	42.46	39.54	43.69	37.89	19.68	29.17	21.33	27.56
InfLLM	22.64	43.70	49.03	49.04	35.61	26.06	30.76	22.70	27.57
TokenSelect	22.44	40.74	47.73	50.33	31.38	24.53	32.56	23.50	27.92
Yi-1.5-6B	17.18	32.56	39.06	36.26	39.25	16.32	30.53	20.21	26.20
NTK	0.80	35.06	29.05	7.47	24.38	0.73	13.66	6.25	25.43
SelfExtend	3.29	19.03	26.00	17.11	11.88	7.73	20.38	17.46	21.79
StreamLLM	15.05	33.27	38.31	34.91	36.92	16.33	29.38	20.02	26.14
InfLLM	17.65	36.25	45.40	41.25	35.89	16.94	30.22	20.85	26.04
TokenSelect	19.36	33.98	48.14	45.05	40.13	22.98	31.59	21.51	26.48
-									
Methods	TREC	TQA	SAMSum	PsgCount	PsgRetrieval	LCC	RepoBench-P	Ave	rage
Methods Qwen2-7B	TREC 78.50	TQA 88.77	SAMSum 46.33	PsgCount 5.50	PsgRetrieval 70.00	LCC 62.40	RepoBench-P 61.95	<b>Ave</b> : 45	<b>rage</b> .37
Methods Qwen2-7B NTK	TREC 78.50 79.50	TQA 88.77 89.51	SAMSum 46.33 46.03	PsgCount 5.50 5.50	PsgRetrieval 70.00 60.00	LCC 62.40 59.36	RepoBench-P 61.95 59.69	Ave: 45 46	rage .37 .17
Methods Qwen2-7B NTK SelfExtend	TREC 78.50 79.50 16.50	TQA 88.77 89.51 27.54	SAMSum 46.33 46.03 29.42	PsgCount 5.50 5.50 4.50	PsgRetrieval 70.00 60.00 0.00	LCC 62.40 59.36 41.42	RepoBench-P 61.95 59.69 41.89	Ave 45 46 18	<b>rage</b> .37 .17 .65
Methods Qwen2-7B NTK SelfExtend StreamLLM	TREC 78.50 79.50 16.50 75.50	TQA 88.77 89.51 27.54 87.19	SAMSum 46.33 46.03 29.42 46.27	PsgCount 5.50 5.50 4.50 3.50	PsgRetrieval 70.00 60.00 0.00 27.50	LCC 62.40 59.36 41.42 61.18	RepoBench-P 61.95 59.69 41.89 61.12	Ave 45 46 18 40	37 .17 .65 .27
Methods Qwen2-7B NTK SelfExtend StreamLLM InfLLM	TREC 78.50 79.50 16.50 75.50 70.50	TQA 88.77 89.51 27.54 87.19 87.51	SAMSum 46.33 46.03 29.42 46.27 44.53	PsgCount 5.50 5.50 4.50 3.50 4.00	PsgRetrieval 70.00 60.00 0.00 27.50 46.50	LCC 62.40 59.36 41.42 61.18 55.08	RepoBench-P 61.95 59.69 41.89 61.12 57.53	Ave 45 46 18 40 42	.37 .17 .65 .27 .90
Methods Qwen2-7B NTK SelfExtend StreamLLM InfLLM TokenSelect	TREC 78.50 79.50 16.50 75.50 70.50 74.00	TQA 88.77 89.51 27.54 87.19 87.51 89.26	SAMSum 46.33 46.03 29.42 46.27 44.53 45.94	PsgCount 5.50 5.50 4.50 3.50 4.00 5.00	PsgRetrieval 70.00 60.00 0.00 27.50 46.50 42.50	LCC 62.40 59.36 41.42 61.18 55.08 61.48	RepoBench-P 61.95 59.69 41.89 61.12 57.53 59.33	Ave 45 46 18 40 42 43	rage .37 .17 .65 .27 .90 .64
Methods Qwen2-7B NTK SelfExtend StreamLLM InfLLM TokenSelect Llama-3-8B	TREC           78.50           79.50           16.50           75.50           70.50           74.00	TQA 88.77 89.51 27.54 87.19 87.51 89.26 90.50	SAMSum           46.33           46.03           29.42           46.27           44.53           45.94           42.30	PsgCount 5.50 5.50 4.50 3.50 4.00 5.00 8.50	PsgRetrieval 70.00 60.00 0.00 27.50 46.50 42.50 62.50	LCC 62.40 59.36 41.42 61.18 55.08 61.48 60.83	RepoBench-P 61.95 59.69 41.89 61.12 57.53 59.33 49.14	Ave 45 46 18 40 42 43 42	rage .37 .17 .65 .27 .90 .64 .46
Methods Qwen2-7B NTK SelfExtend StreamLLM InfLLM TokenSelect Llama-3-8B NTK	TREC           78.50           79.50           16.50           75.50           70.50           74.00           73.00	TQA 88.77 89.51 27.54 87.19 87.51 89.26 90.50 88.74	SAMSum           46.33           46.03           29.42           46.27           44.53           45.94           42.30           42.51	PsgCount 5.50 5.50 4.50 3.50 4.00 5.00 8.50 8.87	PsgRetrieval 70.00 60.00 0.00 27.50 46.50 42.50 62.50 99.50	LCC 62.40 59.36 41.42 61.18 55.08 61.48 60.83 33.62	RepoBench-P 61.95 59.69 41.89 61.12 57.53 59.33 49.14 35.04	Ave 45 46 18 40 42 43 43 42 42 42	rage .37 .17 .65 .27 .90 .64 .12
Methods Qwen2-7B NTK SelfExtend StreamLLM InfLLM TokenSelect Llama-3-8B NTK SelfExtend	TREC           78.50           79.50           16.50           75.50           70.50           74.00           73.00           20.50	TQA 88.77 89.51 27.54 87.19 87.51 89.26 90.50 88.74 16.82	SAMSum           46.33           46.03           29.42           46.27           44.53           45.94           42.30           42.51           25.39	PsgCount 5.50 5.50 4.50 3.50 4.00 5.00 8.50 8.50 8.87 5.75	PsgRetrieval 70.00 60.00 27.50 46.50 42.50 62.50 99.50 7.50	LCC 62.40 59.36 41.42 61.18 55.08 61.48 60.83 33.62 26.24	RepoBench-P 61.95 59.69 41.89 61.12 57.53 59.33 49.14 35.04 31.22	Ave 45 46 18 40 42 43 42 42 42 42	rage .37 .17 .65 .27 .90 .64 .64 .12 .42
Methods Qwen2-7B NTK SelfExtend StreamLLM InfLLM TokenSelect Llama-3-8B NTK SelfExtend StreamLLM	TREC           78.50           79.50           16.50           75.50           70.50           74.00           74.00           73.00           20.50           73.50	TQA 88.77 89.51 27.54 87.19 87.51 89.26 90.50 88.74 16.82 90.08	SAMSum           46.33           46.03           29.42           46.27           44.53           45.94           42.30           42.51           25.39           41.55	PsgCount 5.50 5.50 4.50 3.50 4.00 5.00 8.50 8.50 8.87 5.75 5.00	PsgRetrieval 70.00 60.00 0.00 27.50 46.50 42.50 62.50 99.50 7.50 49.00	LCC 62.40 59.36 41.42 61.18 55.08 61.48 60.83 33.62 26.24 60.35	RepoBench-P           61.95           59.69           41.89           61.12           57.53           59.33           49.14           35.04           31.22           48.95	Ave 45 46 18 40 42 43 42 42 42 42 42 42 42 42 42 42 42 42 42	rage .37 .17 .65 .27 .90 .64 .64 .12 .42 .61
Methods Qwen2-7B NTK SelfExtend StreamLLM InfLLM TokenSelect Llama-3-8B NTK SelfExtend StreamLLM InfLLM	TREC           78.50           79.50           16.50           75.50           70.50           74.00           74.00           73.00           20.50           73.50           73.50	TQA 88.77 89.51 27.54 87.19 87.51 89.26 90.50 88.74 16.82 90.08 90.91	SAMSum           46.33           46.03           29.42           46.27           44.53           45.94           42.30           42.51           25.39           41.55           42.43	PsgCount 5.50 5.50 4.50 3.50 4.00 5.00 8.50 8.87 5.75 5.00 7.17	PsgRetrieval 70.00 60.00 27.50 46.50 42.50 62.50 99.50 7.50 49.00 84.00	LCC 62.40 59.36 41.42 61.18 55.08 61.48 60.83 33.62 26.24 60.35 59.88	RepoBench-P           61.95           59.69           41.89           61.12           57.53           59.33           49.14           31.22           48.95           46.48	Ave 45 46 18 40 42 43 42 42 42 42 42 42 40 40 44	rage .37 .17 .65 .27 .90 .64 .46 .12 .42 .61 .46
Methods         Qwen2-7B         NTK         SelfExtend         StreamLLM         InfLLM         TokenSelect         Llama-3-8B         NTK         SelfExtend         StreamLLM         InfLLM         TokenSelect	TREC           78.50           79.50           16.50           75.50           70.50           74.00           74.00           73.00           20.50           73.50           73.50           67.50	TQA 88.77 89.51 27.54 87.19 87.51 89.26 90.50 88.74 16.82 90.08 90.91 92.22	SAMSum           46.33           46.03           29.42           46.27           44.53           45.94           42.30           42.51           25.39           41.55           42.43	PsgCount 5.50 5.50 4.50 3.50 4.00 5.00 8.50 8.50 8.87 5.75 5.00 7.17 4.54	PsgRetrieval 70.00 60.00 27.50 46.50 42.50 62.50 99.50 7.50 49.00 84.00 87.00	LCC 62.40 59.36 41.42 61.18 55.08 61.48 60.83 33.62 26.24 60.35 59.88 58.86	RepoBench-P           61.95           59.69           41.89           61.12           57.53           59.33           49.14           35.04           31.22           48.95           46.48           51.24	Ave 45 46 18 40 42 43 42 42 42 42 42 42 42 42 44 40 44	rage .37 .17 .65 .27 .90 .64 .64 .46 .12 .42 .61 .46 .04
Methods Qwen2-7B NTK SelfExtend StreamLLM InfLLM TokenSelect Llama-3-8B NTK SelfExtend StreamLLM InfLLM TokenSelect Yi-1.5-6B	TREC           78.50           79.50           16.50           75.50           70.50           74.00           74.00           73.00           20.50           73.50           73.50           67.50           71.50	TQA 88.77 89.51 27.54 87.19 87.51 89.26 90.50 88.74 16.82 90.08 90.91 92.22 48.79	SAMSum           46.33           46.03           29.42           46.27           44.53           45.94           42.30           42.51           25.39           41.55           42.43           42.16	PsgCount 5.50 5.50 4.50 3.50 4.00 5.00 8.50 8.87 5.75 5.00 7.17 4.54 3.00	PsgRetrieval 70.00 60.00 27.50 46.50 42.50 62.50 99.50 7.50 49.00 84.00 87.00 28.50	LCC 62.40 59.36 41.42 61.18 55.08 61.48 60.83 33.62 26.24 60.35 59.88 58.86 57.10	RepoBench-P           61.95           59.69           41.89           61.12           57.53           59.33           49.14           35.04           31.22           48.95           46.48           51.24           52.53	Ave 45 46 18 40 42 43 42 42 42 42 42 42 42 42 42 42 42 42 42	rage
Methods           Qwen2-7B           NTK           SelfExtend           StreamLLM           InfLLM           TokenSelect           Llama-3-8B           NTK           SelfExtend           StreamLLM           InfLLM           TokenSelect           JinfLLM           TokenSelect           Yi-1.5-6B           NTK	TREC           78.50           79.50           16.50           75.50           70.50           74.00           74.00           73.00           20.50           73.50           73.50           73.50           71.50           40.00	TQA 88.77 89.51 27.54 87.19 87.51 89.26 90.50 88.74 16.82 90.08 90.91 92.22 48.79 12.71	SAMSum           46.33           46.03           29.42           46.27           44.53           45.94           42.30           42.51           25.39           41.55           42.43           42.16           0.79           1.34	PsgCount           5.50           5.50           4.50           3.50           4.00           5.00           8.50           8.87           5.75           5.00           7.17           4.54           3.00           0.50	PsgRetrieval           70.00           60.00           0.00           27.50           46.50           42.50           62.50           99.50           7.50           49.00           84.00           87.00           28.50           3.35	LCC 62.40 59.36 41.42 61.18 55.08 61.48 60.83 33.62 26.24 60.35 59.88 58.86 57.10 54.55	RepoBench-P           61.95           59.69           41.89           61.12           57.53           59.33           49.14           35.04           31.22           48.95           46.48           51.24           52.53           37.24	Ave 45 46 18 40 42 43 42 42 42 42 42 42 42 42 42 42 42 42 14 40 44 44 5 8 2 18	rage
Methods Qwen2-7B NTK SelfExtend StreamLLM InfLLM TokenSelect Llama-3-8B NTK SelfExtend StreamLLM InfLLM TokenSelect Yi-1.5-6B NTK SelfExtend	TREC           78.50           79.50           16.50           75.50           70.50           74.00           74.00           73.00           20.50           73.50           67.50           71.50           40.00           23.75	TQA 88.77 89.51 27.54 87.19 87.51 89.26 90.50 88.74 16.82 90.08 90.91 92.22 48.79 12.71 30.61	SAMSum           46.33           46.03           29.42           46.27           44.53           45.94           42.30           42.51           25.39           41.55           42.43           42.16           0.79           1.34           2.58	PsgCount 5.50 5.50 4.50 3.50 4.00 5.00 8.50 8.87 5.75 5.00 7.17 4.54 3.00 0.50 2.75	PsgRetrieval           70.00           60.00           0.00           27.50           46.50           42.50           62.50           99.50           7.50           49.00           84.00           87.00           28.50           3.35           13.50	LCC 62.40 59.36 41.42 61.18 55.08 61.48 60.83 33.62 26.24 60.35 59.88 58.86 57.10 54.55 43.17	RepoBench-P           61.95           59.69           41.89           61.12           57.53           59.33           49.14           35.04           31.22           48.95           46.48           51.24           52.53           37.24           35.45	Ave 45 46 18 40 42 43 42 42 42 42 42 42 42 42 44 44 44 44 8 18	rage
Methods Qwen2-7B NTK SelfExtend StreamLLM InfLLM TokenSelect Llama-3-8B NTK SelfExtend StreamLLM InfLLM TokenSelect Yi-1.5-6B NTK SelfExtend StreamLLM	TREC           78.50           79.50           16.50           75.50           70.50           74.00           74.00           73.00           20.50           73.50           67.50           71.50           40.00           23.75           69.00	TQA 88.77 89.51 27.54 87.19 87.51 89.26 90.50 88.74 16.82 90.08 90.91 92.22 48.79 12.71 30.61 73.36	SAMSum           46.33           46.03           29.42           46.27           44.53           45.94           42.30           42.51           25.39           41.55           42.43           42.16           0.79           1.34           2.58           0.82	PsgCount 5.50 5.50 4.50 3.50 4.00 5.00 8.50 8.87 5.75 5.00 7.17 4.54 3.00 0.50 2.75 2.50	PsgRetrieval 70.00 60.00 0.00 27.50 46.50 42.50 62.50 99.50 7.50 49.00 84.00 87.00 28.50 3.35 13.50 18.50	LCC 62.40 59.36 41.42 61.18 55.08 61.48 60.83 33.62 26.24 60.35 59.88 58.86 57.10 54.55 43.17 56.37	RepoBench-P           61.95           59.69           41.89           61.12           57.53           59.33           49.14           35.04           31.22           48.95           46.48           51.24           52.53           37.24           35.45           49.05	Ave 45 46 18 40 42 43 42 42 42 42 42 42 14 40 44 44 44 8 18 8 32	rage .37 .17 .65 .27 .90 .64 .46 .12 .42 .61 .46 .04 .48 .28 .53 .49
Methods Qwen2-7B NTK SelfExtend StreamLLM InfLLM TokenSelect Llama-3-8B NTK SelfExtend StreamLLM InfLLM StreamLLM InfLLM	TREC           78.50           79.50           16.50           75.50           70.50           74.00           74.00           73.00           20.50           73.50           67.50           71.50           40.00           23.75           69.00           71.50	TQA 88.77 89.51 27.54 87.19 87.51 89.26 90.50 88.74 16.82 90.08 90.91 92.22 48.79 12.71 30.61 73.36 71.49	SAMSum           46.33           46.03           29.42           46.27           44.53           45.94           42.30           42.51           25.39           41.55           42.43           42.16           0.79           1.34           2.58           0.82           1.01	PsgCount           5.50           5.50           4.50           3.50           4.00           5.00           8.50           8.87           5.75           5.00           7.17           4.54           3.00           0.50           2.75           2.50           4.00	PsgRetrieval 70.00 60.00 0.00 27.50 46.50 42.50 62.50 99.50 7.50 49.00 84.00 87.00 28.50 3.35 13.50 18.50 10.50	LCC 62.40 59.36 41.42 61.18 55.08 61.48 60.83 33.62 26.24 60.35 59.88 58.86 57.10 54.55 43.17 56.37 56.88	RepoBench-P           61.95           59.69           41.89           61.12           57.53           59.33           49.14           35.04           31.22           48.95           46.48           51.24           52.53           37.24           35.45           49.05           46.28	Ave 45 46 18 40 42 43 42 42 42 14 40 44 44 44 32 18 18 8 23 33	rage

Table 9: Comparison of different methods with different origin models on LongBench.