How Sparse Attention Approximates Exact Attention?Your Attention is Naturally n^{C} -Sparse

Anonymous authors

006

008 009 010

011

013

014

015

016

017

018

019

021

024

025

026

027

028

029

031 032 Paper under double-blind review

Abstract

Sparse Attention is a technique that approximates standard attention computation with sub-quadratic complexity. This is achieved by selectively ignoring smaller entries in the attention matrix during the softmax function computation. Variations of this technique, such as pruning KV cache, sparsity-based fast attention, and Sparse Transformer, have been extensively utilized for efficient Large Language Models (LLMs) deployment. Despite its widespread use, a theoretical understanding of the conditions under which sparse attention performs on par with traditional attention remains elusive. This work aims to **bridge this gap by examining the inherent sparsity of standard attention processes.** Our theoretical framework reveals several brand-new key insights:

- Attention is n^C -sparse, implying that considering only the largest $\Omega(n^C)$ entries out of all n entries is sufficient for sparse attention to approximate the exact attention matrix with decreasing loss. Here, n represents the input length and $C \in (0, 1)$ is a constant.
- Stable $o(\log(n))$ -sparse attention, which approximates attention computation with $\log(n)$ or fewer entries, may not be feasible since the error will persist at a minimum of O(1).
- An adaptive strategy (α · n^C, α ∈ ℝ) for the window size of efficient attention methods rather than a fixed one is guaranteed to perform more accurately and efficiently in a task for inference on flexible context lengths.
- 033 1 INTRODUCTION

Large Language Models (LLMs) Vaswani et al. (2017); Radford et al. (2018); Devlin et al. (2018);
Radford et al. (2019); Brown et al. (2020); Chowdhery et al. (2022); Zhang et al. (2022); ChatGPT (2022) have emerged as a cornerstone of contemporary artificial intelligence, exhibiting remarkable
capabilities across a plethora of AI domains. Their prowess is grounded in their ability to comprehend and generate human language with a level of sophistication that is unprecedented. This has catalyzed transformative applications in natural language processing, including machine translation He et al. (2021), content creation ChatGPT (2022); OpenAI (2023), and beyond, underscoring the profound impact of LLMs on the field of AI.

However, the architectural backbone of these models, particularly those built on the transformer framework Vaswani et al. (2017), presents a significant challenge: computational efficiency Tay et al. (2022). The essence of the transformer architecture, the Attention mechanism, necessitates a computational and memory complexity of $O(n^2)$, where *n* represents the sequence length. This quadratic dependency limits the scalability of LLMs, especially as we venture into processing longer sequences or expanding model capacities.

In an effort to mitigate this bottleneck, the AI research community has pivoted towards innovative
 solutions, one of which is sparse attention Child et al. (2019); Correia et al. (2019). Sparse attention
 mechanisms aim to approximate the results of the full attention computation by selectively focusing
 on a subset of the input data points. This is typically achieved by omitting certain interactions in
 the Query and Key multiplications within the attention mechanism, thereby inducing sparsity in the
 attention matrix. In order to arrive at the goal of preserving the model's performance while alleviating

the computational and memory demands, prior works, including pruning KV cache, sparsity-based fast attention, and sparse transformer modeling, demonstrate outstanding efficiencies with $O(n^{1+o(1)})$ (sub-quadratic) complexity and sub-linear memory cache with competitive performance compared with standard attention across various tasks Liu et al. (2023b); Zhang et al. (2024); Kacham et al. (2023); Addanki et al. (2023); Lee et al. (2024); Xiong et al. (2021); Zandieh et al. (2023); Alman & Song (2023; 2024b;a); Han et al. (2023).

Despite these advancements, the theoretical underpinnings of sparse attention mechanisms and their
 implications on model performance and behavior remain an area of active inquiry. In detail, it's not
 clear when and when not sparse attention can approximate standard attention with a stable error.
 Also, the sparsity that attention naturally processes, which we call attention sparsity, lacks a strict
 confirmation of its existence and measurement. Especially, we would like to ask:

How Sparse Attention Approximates Exact Attention?

Our Contributions. In this work, we explore the theory of the sparse attention computation problem.
 Particularly, we first provide a analysis framework that first theoretically confirms the sparsity appears in standard attention. In detailed, our analysis describes the relationships between attention sparsity and input boundary, weights of attention networks and context length. Therefore, we derive several incremental insights based on this framework.

072 073

074

079

081

082

083

084

065

066

2 PRELIMINARY

Assumption. In this work, we consider one-layer self-attention computation both in standard form and sparsity-based approximate form. To begin with, we give the assumption of the input matrix of attention computation, denoted as $X \in \mathbb{R}^{n \times d}$ where *n* is the context length and *d* stands the dimension, as follows (refer to Definition C.1 for the formal and detailed version of assumption):

- Independent Entries. For any two entries X_{i1,j1} and X_{i2,j2} in matrix X, ∀i1, i2 ∈ [n] and j1, j2 ∈ [d], they are independent.
- **Bounded Entries.** For failure probability $\delta \in (0, 0.1)$. With a probability 1δ , the entry $X_{i,j}$ in matrix $X, \forall i \in [n]$ and $j \in [d]$, we have $|X_{i,j}| \leq B$ for some positive constant B > 0.

Attention Computation. Hence, we are about to introduce the standard attention computation, which occupies the main time and space complexity $O(n^2)$ in LLMs inference. First, we denote the weights of query, key and value projection as $W_Q, W_K, W_V \in \mathbb{R}^{d \times d}$. Thus, we let query, key and value state matrices be computed by $Q := XW_Q, K := XW_K, V := XW_V \in \mathbb{R}^{n \times d}$. We state the following definition:

Definition 2.1 (Attention computation). Given Query, Key and Value states matrices $Q, K, V \in \mathbb{R}^{n \times d}$. We then define $A := \exp(QK^{\top}/\sqrt{d}), D := \operatorname{diag}(A\mathbf{1}_n) \in \mathbb{R}^{n \times n}$. The attention computation Attn $(Q, K, V) \in \mathbb{R}^{n \times d}$ is given by: Attn $(Q, K, V) := D^{-1}AV \in \mathbb{R}^{n \times d}$. Specially, we denote $D^{-1} = \operatorname{diag}(1/(A\mathbf{1}_n)) \in \mathbb{R}^{n \times n}$.

Attention Sparsity. In Definition 2.1, $D^{-1}A \in \mathbb{R}^{n \times n}$ represents the attention matrix, indicating how 095 much the model focuses on each vector. In much of the sparse attention literature, $D^{-1}A$ is assumed 096 to be sufficiently sparse, allowing sparsity-based efficient attention methods to disregard some zero 097 entries in order to achieve a balance between accuracy and efficiency. In this paper, we introduce 098 a threshold, denoted as ϵ , and define attention sparsity as the number of entries in each row of 099 $D^{-1}A \in \mathbb{R}^n$ that are smaller than ϵ . Specifically, for a softmax vector $u \in \mathbb{R}^n$, if there are at least 100 n-k entries in u that are not greater than ϵ for all integers $k \in [n]$, we say that u is (ϵ, k) -sparse. 101 Since ϵ is intended to be a very small value, we will simply refer to u as being k-sparse. The formal 102 definition is provided below:

Definition 2.2 ((ϵ , k)-sparsity). For a vector $u \in \mathbb{R}^n$ and error $\epsilon > 0$, we define sparse set $S_{\epsilon}(u)$ as: $S_{\epsilon}(u) := \{i \in [n] \mid |u_i| \le \epsilon\}$. Hence, we say u is at least (ϵ , k)-sparse when it holds that $|S_{\epsilon}(u)| \ge n - k$.

Problem Definition I: Estimating ϵ **.** In practical implementations, establishing a clear relationship between ϵ and the sparsity k proves to be challenging. Therefore, we first analyze how to estimate a

boundary for ϵ based on a given sparsity integer $k \in [n]$. Naively, given k, we would like to find a guaranteed value for ϵ that satisfies $|S_{\epsilon}(u)| \ge n - k$. Address this problem will enable us to assess the loss associated with approximating standard attention using sparse attention, ultimately guiding us in finding the optimal trade-off between ϵ (where lower values yield greater accuracy) and k (where lower values lead to higher efficiency).

Sparse Attention and Approximation. Here we state an ideal mathematical definition for the sparse attention in this paper. Initially, we define a set, $\mathcal{T}_k(u)$, to filter out the greatest k entries in a vector $u \in \mathbb{R}^n$. The integer $k \in [1, n]$ is also called window size in some sparse attention works.

Definition 2.3. For a vector $u \in \mathbb{R}^n$, given a sparsity integer k, we denote a top-k set $\mathcal{T}_k(u) := \{i \in [n] \mid S_{u_i}(u) \ge n-k\}$, then we define vector $\mathsf{topk}(u) := [u_1 \cdot \mathbf{1}_{1 \in \mathcal{T}_k(u)}, \cdots, u_n \cdot \mathbf{1}_{n \in \mathcal{T}_k(u)}]^\top \in \mathbb{R}^n$.

119 Note that $\mathbf{1}_{i \in \mathcal{T}_k(u)}$ is an indicator where when $i \in \mathcal{T}_k(u)$, it equals 1, otherwise, 0. We utilize topk(u)120 to compute a sparsity-based approximating version of $A = \exp(QK^{\top})$ in Definition 2.1, we denote 121 it A_{spar} . Accordingly, we provide a universal version for all sparsity-based attention as follows:

Definition 2.4 (Sparse attention). Given Query, Key and Value state matrices $Q, K, V \in \mathbb{R}^{n \times d}$. **R**^{n \times d}. We then define $A := \exp(QK^{\top}/\sqrt{d}) \in \mathbb{R}^{n \times n}$. Especially, we define $A_{\text{spar}} := [\operatorname{topk}(A_{1,*}), \cdots, \operatorname{topk}(A_{n,*})]^{\top}$, $D_{\text{spar}} := \operatorname{diag}(A_{\text{spar}}\mathbf{1}_n) \in \mathbb{R}^{n \times n}$. The sparse attention computation SparseAttn $(Q, K, V) \in \mathbb{R}^{n \times d}$ is given by:

127

 $\mathsf{SparseAttn}(Q, K, V) := D_{\mathrm{spar}}^{-1} A_{\mathrm{spar}} V \in \mathbb{R}^{n \times d}$

128 129 Specially, we denote $D_{\text{spar}}^{-1} = \text{diag}(1/(A_{\text{spar}}\mathbf{1}_n)) \in \mathbb{R}^{n \times n}$.

130 It should be noted that directly accessing top k entries in the attention matrix without any extra 131 computational cost is overly ideal for efficient LLMs in real-world cases. Prior works usually utilize 132 some additional approximate algorithm to meet this condition, e.g. Locality-Sensitive Hashing (LSH) 133 for retrieving larger query-key pairs, but this also brings more approximating errors. We only focus on 134 the part of approximating attention computation in this study and leave the part of pre-approximating 135 top-k entries in $D^{-1}A$ as a future direction.

Problem Definition II: Sparse Attention Approximation. The variations of sparse attention, including pruning KV Cache Liu et al. (2021); Xiao et al. (2023) and sparsity-based attention Kitaev et al. (2020); Zandieh et al. (2023); Han et al. (2023), focus on solving the approximation of the attention matrix, where we call it sparse attention approximation. In particular, we emphasize the importance of *stable sparse attention approximation*, which directly affects the extensibility of sparse attention under long context scenes. We denote $f : \mathbb{N}^+ \to \mathbb{N}^+$ as the strategy to choose a suitable window size due to different input lengths. Hence, we give:

Definition 2.5 (Stable sparse attention approximation SSAA(f)). For some strategy $f : \mathbb{N}^+ \to \mathbb{N}^+$ to choose the sparsity k = f(n) in sparse attention (Definition 2.4), the problem of stable sparse attention approximation SSAA(f) is to solve: $L(f,n) = \|D_{\text{spar}}^{-1}A_{\text{spar}} - D^{-1}A\|_p$, where $\|\cdot\|_p$ denotes some norm. We say this sparse attention approximation is **stable** iff:

• L(f, n) is monotonically decreasing with growing n.

• $\lim_{n \to +\infty} L(f, n) = 0.$

149 150 151

152 153

154

147

148

3 INSIGHTS OVERVIEW

Definition 3.1. Denote $W := W_Q W_K^\top / \sqrt{d} \in \mathbb{R}^{d \times d}$. We define $R := B^2 \cdot ||W||_F$.

We estimate the lower bound on the requirement for (ϵ, k) -sparse softmax vector, proving the vanilla attention computation is naturally sparse.

Theorem 3.2. Let $R \ge 0$ be defined as Definition 3.1. Given sparsity integer $k \le n$. Denote $T := \exp(\sqrt{\log(n(n-k)d/\delta)})$. Let S_{ϵ} be defined as Definition 2.2. $\delta \in (0,0.1)$. If we choose $\epsilon \ge \frac{T^{O(R)}}{n}$, then with a probability at least $1 - \delta$, for all $i \in [n]$, we have $\left|S_{\epsilon}(D_{i,i}^{-1}A_{i,*})\right| \ge n - k$.

Proof sketch of Theorem 3.2. The complete proof is provided in Appendix D and Theorem D.1.

We introduce the concept of *attention collapse*, which demonstrates the number of effective entries in attention matrix provably decrease to 1 or some constant inevitably.

Theorem 3.3. Consider a fixed ϵ with a very small value, $\delta \in (0, 0.1)$. Then with a probability at least $1 - \delta$, there is:

167 168

169 170 171

172

180

181 182

183

184

185

186

187 188 • Part 1. If $R = o(\sqrt{\log(n)})$, then we have $\lim_{n \to +\infty} |\mathcal{S}_{\epsilon}(u)| \ge n - 1$.

• Part 2. If $R = O(\sqrt{\log(n)})$, then we have $\lim_{n \to +\infty} |\mathcal{S}_{\epsilon}(u)| \ge O(1)$.

Proof sketch of Theorem 3.3. Refer to Theorem D.2 for the detailed proof.

We give the sufficient lower bound on the window size of stable sparse attention approximating exact attention computation, $\Omega(n^C)$ for constant $C \in (0, 1)$. This further indicates that sparse attention can recover attention outputs from limited $\Omega(n^C)$ entries while achieving a decreasing error.

Theorem 3.4. $\delta \in (0, 0.1)$. For a constant $C \in (0, 1)$, we then denote $f(n) := \Omega(n^C)$, therefore, with a probability at least $1 - \delta$, window size strategy k = f(n) is sufficient to solve SSAA(f) in Definition 2.5.

Proof sketch of Theorem 3.4. Please see Theorem E.1.

Meanwhile, we also confirm sparse attention approximation from $o(\log(n))$ entries is not enough for stability and extensibility since the lower bound on error will grow with increasing input length.

Theorem 3.5. $\delta \in (0, 0.1)$. For a constant $C \in (0, 1)$, we then denote $f(n) := o(\log(n))$, therefore, with a probability at least $1 - \delta$, window size strategy k = f(n) cannot solve SSAA(f) in Definition 2.5.

Proof sketch of Theorem 3.5. Please refer to Theorem E.2 for the formal version and corresponding detailed proofs.

190 191

198 199

200

201 202

203

204

205 206 207

208 209

210

189

Therefore, we suggest to use adaptive strategy $k = \alpha \cdot n^C$, $\alpha > 0$, $C \in (0, 1)$ for the window size of sparse attention rather than the strategy that fixes the window size for any input. The former is proved more efficient within higher approximation performance. We consider a dataset $\mathcal{D} := \{X^i\}_{i=1}^N$ with dataset size N. For all $i \in [N]$, we use n_i to denote the context length, such that $X^i \in \mathbb{R}^{n_i \times d}$. Hence, the difference of the computational complexities of fixed Top-k strategy and a dynamic (especially $O(n^C)$) strategy could be easily obtained in the Claim below.

Claim 3.6. We have:

• **Part 1.** Choosing the constant window size strategy k = p for some constant integer p > 0. The computational complexity of a one-layer p-sparse attention to inference $\mathcal{D} = \{X^i\}_i^N$ is $\Theta(p\sum_{i=1}^N n_i)$.

• **Part 2.** Choosing the constant window size strategy $k = \alpha \cdot n^C$ for some constant $\alpha, C > 0$. The computational complexity of a one-layer *p*-sparse attention to inference $\mathcal{D} = \{X^i\}_i^N$ is $\Theta(\alpha \sum_{i=1}^N n_i^{1+C})$.

Proof. Recall the computational complexity of Definition 2.4 is O(nk) for input $X \in \mathbb{R}^{n \times d}$ and k window size. We then take the summation for each $X^i \in \mathbb{R}^{n_i \times d}$ in $\mathcal{D} := \{X^i\}_i^N$ to obtain the results of **Part 1** and **Part 2**.

Proposition 3.7. For any window size strategy k = p for some constant integer p > 0, there exist a context-length adaptive strategy $k = \alpha \cdot n^C$ for some constant $\alpha, C > 0$ that performs lower approximating error.

214

215 *Proof.* Following Theorem 3.4 and Theorem 3.5, the conclusion of this proposition can be trivially proved by plugging suitable choices of α and C.

216	References
217	

- Addanki, R., Li, C., Song, Z., and Yang, C. One pass streaming algorithm for super long token attention approximation in sublinear space. *arXiv preprint arXiv:2311.14652*, 2023.
- Alman, J. and Song, Z. Fast attention requires bounded entries. Advances in Neural Information Processing Systems (NeurIPS), 36, 2023.
- Alman, J. and Song, Z. The fine-grained complexity of gradient computation for training large
 language models. *arXiv preprint arXiv:2402.04497*, 2024a.
- Alman, J. and Song, Z. How to capture higher-order correlations? generalizing matrix softmax attention to kronecker computation. In *The Twelfth International Conference on Learning Representations (ICLR)*, 2024b. URL https://openreview.net/forum?id=v0zNCwwkaV.
 - Beltagy, I., Peters, M. E., and Cohan, A. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*, 2020.
- Bernstein, S. On a modification of chebyshev's inequality and of the error formula of laplace. *Ann. Sci. Inst. Sav. Ukraine, Sect. Math*, 1(4):38–49, 1924.
- Brand, J. v. d., Song, Z., and Zhou, T. Algorithm and hardness for dynamic attention maintenance in
 large language models. *arXiv preprint arXiv:2304.02207*, 2023.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam,
 P., Sastry, G., Askell, A., et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Cai, H., Lou, Y., Mckenzie, D., and Yin, W. A zeroth-order block coordinate descent algorithm for
 huge-scale black-box optimization. *arXiv preprint arXiv:2102.10707*, 2021.
- ChatGPT. Optimizing language models for dialogue. OpenAI Blog, November 2022. URL https: //openai.com/blog/chatgpt/.
- Chen, B., Liu, Z., Peng, B., Xu, Z., Li, J. L., Dao, T., Song, Z., Shrivastava, A., and Re, C. Mongoose:
 A learnable lsh framework for efficient neural network training. In *International Conference on Learning Representations*, 2020.
- Chernoff, H. A measure of asymptotic efficiency for tests of a hypothesis based on the sum of observations. *The Annals of Mathematical Statistics*, pp. 493–507, 1952.
- Child, R., Gray, S., Radford, A., and Sutskever, I. Generating long sequences with sparse transformers.
 arXiv preprint arXiv:1904.10509, 2019.
- Choromanski, K., Likhosherstov, V., Dohan, D., Song, X., Gane, A., Sarlos, T., Hawkins, P., Davis, J., Mohiuddin, A., Kaiser, L., et al. Rethinking attention with performers. *arXiv preprint arXiv:2009.14794*, 2020.
- Chowdhery, A., Narang, S., Devlin, J., Bosma, M., Mishra, G., Roberts, A., Barham, P., Chung,
 H. W., Sutton, C., Gehrmann, S., et al. Palm: Scaling language modeling with pathways. *arXiv* preprint arXiv:2204.02311, 2022.
- Chu, T., Song, Z., and Yang, C. How to protect copyright data in optimization of large language models? In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 17871–17879, 2024.
- 264 Correia, G. M., Niculae, V., and Martins, A. F. Adaptively sparse transformers. *arXiv preprint* 265 *arXiv:1909.00015*, 2019.
- Deng, Y., Li, Z., Mahadevan, S., and Song, Z. Zero-th order algorithm for softmax attention optimization. *arXiv preprint arXiv:2307.08352*, 2023a.
- 269 Deng, Y., Mahadevan, S., and Song, Z. Randomized and deterministic attention sparsification algorithms for over-parameterized feature dimension. *arXiv preprint arXiv:2304.04397*, 2023b.

286

287

295

299

300

301

308

- Deng, Y., Song, Z., Xie, S., and Yang, C. Unmasking transformers: A theoretical approach to data recovery via attention weights. *arXiv preprint arXiv:2310.12462*, 2023c.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- Foss, S., Korshunov, D., Zachary, S., et al. An introduction to heavy-tailed and subexponential distributions, volume 6. Springer, 2011.
- Gao, Y., Song, Z., Wang, W., and Yin, J. A fast optimization view: Reformulating single layer attention in llm based on tensor and svm trick, and solving it in matrix multiplication time. *arXiv* preprint arXiv:2309.07418, 2023a.
- Gao, Y., Song, Z., and Xie, S. In-context learning for attention scheme: from single softmax regression to multiple softmax regression via a tensor trick. *arXiv preprint arXiv:2307.02419*, 2023b.
 - Gao, Y., Song, Z., and Yin, J. Gradientcoin: A peer-to-peer decentralized large language models. *arXiv preprint arXiv:2308.10502*, 2023c.
- Haagerup, U. The best constants in the khintchine inequality. *Studia Mathematica*, 70(3):231–283, 1981.
- Han, I., Jayaram, R., Karbasi, A., Mirrokni, V., Woodruff, D. P., and Zandieh, A. Hyperattention:
 Long-context attention in near-linear time. *arXiv preprint arXiv:2310.05869*, 2023.
- Hanson, D. L. and Wright, F. T. A bound on tail probabilities for quadratic forms in independent
 random variables. *The Annals of Mathematical Statistics*, 42(3):1079–1083, 1971.
- He, W., Wu, Y., and Li, X. Attention mechanism for neural machine translation: A survey. In 2021
 IEEE 5th Information Technology, Networking, Electronic and Automation Control Conference (*ITNEC*), volume 5, pp. 1485–1489. IEEE, 2021.
 - Hoeffding, W. Probability inequalities for sums of bounded random variables. The collected works of Wassily Hoeffding, pp. 409–426, 1994.
- Kacham, P., Mirrokni, V., and Zhong, P. Polysketchformer: Fast transformers via sketches for
 polynomial kernels. *arXiv preprint arXiv:2310.01655*, 2023.
- Khintchine, A. Über dyadische brüche. *Mathematische Zeitschrift*, 18(1):109–116, 1923.
- Kitaev, N., Kaiser, Ł., and Levskaya, A. Reformer: The efficient transformer. *arXiv preprint arXiv:2001.04451*, 2020.
- Lample, G., Sablayrolles, A., Ranzato, M., Denoyer, L., and Jégou, H. Large memory layers with
 product keys. *Advances in Neural Information Processing Systems*, 32, 2019.
- Laurent, B. and Massart, P. Adaptive estimation of a quadratic functional by model selection. *Annals of statistics*, pp. 1302–1338, 2000.
- Lee, S., Lee, H., and Shin, D. Proxyformer: Nyström-based linear transformer with trainable
 proxy tokens. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 13418–13426, 2024.
- Li, C., Song, Z., Wang, W., and Yang, C. A theoretical insight into attack and defense of gradient leakage in transformer. *arXiv preprint arXiv:2311.13624*, 2023a.
- Li, Z., Song, Z., and Zhou, T. Solving regularized exp, cosh and sinh regression problems. *arXiv* preprint, 2303.15725, 2023b.
- 323 Liu, H., Li, Z., Hall, D., Liang, P., and Ma, T. Sophia: A scalable stochastic second-order optimizer for language model pre-training. *arXiv preprint arXiv:2305.14342*, 2023a.

324 325 326	Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., and Guo, B. Swin transformer: Hierar- chical vision transformer using shifted windows. In <i>Proceedings of the IEEE/CVF international</i> <i>conference on computer vision</i> , pp. 10012–10022, 2021.
327 328 329 330	Liu, Z., Wang, J., Dao, T., Zhou, T., Yuan, B., Song, Z., Shrivastava, A., Zhang, C., Tian, Y., Re, C., et al. Deja vu: Contextual sparsity for efficient llms at inference time. In <i>International Conference on Machine Learning</i> , pp. 22137–22176. PMLR, 2023b.
331 332	Lu, Y., Dhillon, P., Foster, D. P., and Ungar, L. Faster ridge regression via the subsampled randomized hadamard transform. <i>Advances in neural information processing systems</i> , 26, 2013.
334 335	Malladi, S., Gao, T., Nichani, E., Damian, A., Lee, J. D., Chen, D., and Arora, S. Fine-tuning language models with just forward passes. <i>arXiv preprint arXiv:2305.17333</i> , 2023a.
336 337 338	Malladi, S., Wettig, A., Yu, D., Chen, D., and Arora, S. A kernel-based view of language model fine-tuning. In <i>International Conference on Machine Learning</i> , pp. 23610–23641. PMLR, 2023b.
339	OpenAI. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023.
340 341 342	Panigrahi, A., Saunshi, N., Zhao, H., and Arora, S. Task-specific skill localization in fine-tuned language models. <i>arXiv preprint arXiv:2302.06600</i> , 2023.
343 344 345	Radford, A., Narasimhan, K., Salimans, T., Sutskever, I., et al. Improving language understanding by generative pre-training, 2018.
345 346 347	Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I., et al. Language models are unsupervised multitask learners. <i>OpenAI blog</i> , 1(8):9, 2019.
348 349 350	Rafailov, R., Sharma, A., Mitchell, E., Ermon, S., D.Manning, C., and Finn, C. Direct preference optimization: Your language model is secretly a reward model. <i>arXiv preprint arXiv:2305.18290</i> , 2023.
352	Rudelson, M. and Vershynin, R. Hanson-wright inequality and sub-gaussian concentration. 2013.
353 354 355	Song, Z., Wang, W., and Yin, J. A unified scheme of resnet and softmax. <i>arXiv preprint arXiv:2309.13482</i> , 2023a.
356 357	Song, Z., Yin, J., and Zhang, L. Solving attention kernel regression problem via pre-conditioner. <i>arXiv preprint arXiv:2308.14304</i> , 2023b.
358 359 360	Sun, Z., Yang, Y., and Yoo, S. Sparse attention with learning to hash. In <i>International Conference on Learning Representations</i> , 2021.
361 362	Tay, Y., Dehghani, M., Bahri, D., and Metzler, D. Efficient transformers: A survey. ACM Computing Surveys, 55(6):1–28, 2022.
363 364 365	Tropp, J. A. Improved analysis of the subsampled randomized hadamard transform. <i>Advances in Adaptive Data Analysis</i> , 3(01n02):115–126, 2011.
366 367 368 369	Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention is all you need. <i>Advances in neural information processing systems</i> , 30, 2017.
370 371	Xiao, G., Tian, Y., Chen, B., Han, S., and Lewis, M. Efficient streaming language models with attention sinks. <i>arXiv preprint arXiv:2309.17453</i> , 2023.
372 373 374 375	Xiong, Y., Zeng, Z., Chakraborty, R., Tan, M., Fung, G., Li, Y., and Singh, V. Nyströmformer: A nyström-based algorithm for approximating self-attention. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 35, pp. 14138–14148, 2021.
376 377	Zaheer, M., Guruganesh, G., Dubey, K. A., Ainslie, J., Alberti, C., Ontanon, S., Pham, P., Ravula, A., Wang, Q., Yang, L., et al. Big bird: Transformers for longer sequences. <i>Advances in neural information processing systems</i> , 33:17283–17297, 2020.

378	Zandieh, A., Han, I., Daliri, M., and Karbasi, A. Kdeformer: Accelerating transformers via kernel
379	density estimation. In International Conference on Machine Learning, pp. 40605–40623. PMLR,
380	2023.
381	

- Zhang, S., Roller, S., Goyal, N., Artetxe, M., Chen, M., Chen, S., Dewan, C., Diab, M., Li, X., Lin, X. V., et al. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*, 2022.
- Zhang, Z., Sheng, Y., Zhou, T., Chen, T., Zheng, L., Cai, R., Song, Z., Tian, Y., Ré, C., Barrett,
 C., et al. H20: Heavy-hitter oracle for efficient generative inference of large language models.
 Advances in Neural Information Processing Systems, 36, 2024.

390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429

		Appendix	
C	литі	ENTS	
C	JNII		
1	Intr	oduction	1
2	Prel	iminary	2
3	Inci	ahts Overview	3
5	111.51		5
A	Rela	uted Work	10
B	Prel	iminary	10
	B .1	Notations	10
	B.2	Basic Fact for Softmax	11
	B.3	Basic Facts for Calculation	11
	B.4	Probability Tools	11
С	Prol	blem Definitions	13
	C.1	Input Assumption	13
	C.2	Attention Computation	13
	C.3	ϵ -Approximated k-Sparse Softmax Vector	13
	C.4	Sparse Attention	13
	C.5	Helpful Definitions	14
)	Atte	ention Sparsity	14
	D.1	Main Result 1: Attention Sparsity with Upper Bound on Error	14
	D.2	Main Result 2: Attention Collapse	15
	D.3	Bounding D^{-1}	16
	D.4	Concentration on QK^{\top}	17
E	Sna	rse Attention Approximation	18
	E 1	Main Result 3: Upper Bound on Error	18
	E.1	Lower Bound on Error	19
	E.2	Approximating Softmax Function	20
	E.J	Helpful Bound Toolkit	20
	с.4		22

486 A RELATED WORK

487 488

Sparse and Efficient Transformer. In the landscape of attention mechanisms, Vaswani et al. 489 introduced the transformative transformer model, revolutionizing NLP with its comprehensive self-490 attention mechanism Vaswani et al. (2017). Innovations Child et al. (2019); Lample et al. (2019) in 491 sparse attention presented methods to reduce complexity, maintaining essential contextual information 492 while improving computational efficiency. The Reformer Kitaev et al. (2020) utilized Locality 493 Sensitive Hashing to significantly cut down computational demands, enabling the processing of lengthy sequences. Mongoose Chen et al. (2020) adapted sparsity patterns dynamically, optimizing 494 computation without losing robustness. Sun et al. (2021) introduced a learning-to-hash strategy 495 to generate sparse attention patterns, enhancing data-driven efficiency. HyperAttention Han et al. 496 (2023) refined attention approximation, balancing computational savings with accuracy. Longformer 497 Beltagy et al. (2020) extended transformer capabilities to longer texts through a mix of global and 498 local attention mechanisms. The Performe Choromanski et al. (2020) offered a novel approximation 499 of softmax attention, reducing memory usage for long sequences. Big Bird Zaheer et al. (2020) 500 combined global, local, and random attention strategies to surmount traditional transformer limitations 501 regarding sequence length. 502

Theoretical Approaches to Understanding LLMs. There have been notable advancements in 504 the field of regression models, particularly with the exploration of diverse activation functions, 505 aiding in the comprehension and optimization of these models. The study of over-parameterized 506 neural networks, focusing on exponential and hyperbolic activation functions, has shed light on their 507 convergence traits and computational benefits Brand et al. (2023); Song et al. (2023a); Gao et al. (2023c); Deng et al. (2023b); Gao et al. (2023b); Song et al. (2023b); Zandieh et al. (2023); Alman & 508 Song (2023; 2024b;a); Gao et al. (2023a); Deng et al. (2023c); Li et al. (2023a); Chu et al. (2024). 509 Enhancements in this area include the addition of regularization components and the innovation 510 of algorithms like the convergent approximation Newton method to improve performance Li et al. 511 (2023b). Additionally, employing tensor methods to simplify regression models has facilitated 512 in-depth analyses concerning Lipschitz constants and time complexity Gao et al. (2023b); Deng et al. 513 (2023a). Concurrently, there's a burgeoning interest in optimization algorithms specifically crafted 514 for LLMs, with block gradient estimators being utilized for vast optimization challenges, significantly 515 reducing computational load Cai et al. (2021). Novel methods such as Direct Preference Optimization 516 are revolutionizing the tuning of LLMs by using human preference data, circumventing the need for 517 traditional reward models Rafailov et al. (2023). Progress in second-order optimizers is also notable, 518 offering more leniency in convergence proofs by relaxing the usual Lipschitz Hessian assumptions Liu et al. (2023a). Moreover, a series of studies focus on the intricacies of fine-tuning Malladi et al. 519 (2023a;b); Panigrahi et al. (2023). These theoretical developments collectively push the boundaries 520 of our understanding and optimization of LLMs, introducing new solutions to tackle challenges like 521 the non-strict Hessian Lipschitz conditions. 522

523 524

527 528

529

530 531

532

534

538

B PRELIMINARY

B.1 NOTATIONS

In this work, we use the following notations and definitions:

- For integer n, we use [n] to denote the set $\{1, \ldots, n\}$.
- We use $\mathbf{1}_n$ to denote all-1 vector in \mathbb{R}^n .
- The ℓ_p norm of a vector x is denoted as $||x||_p$, for examples, $||x||_1 := \sum_{i=1}^n |x_i|$, $||x||_2 := (\sum_{i=1}^n x_i^2)^{1/2}$ and $||x||_{\infty} := \max_{i \in [n]} |x_i|$.
- For a vector $x \in \mathbb{R}^n$, $\exp(x) \in \mathbb{R}^n$ denotes a vector where whose i-th entry is $\exp(x_i)$ for all $i \in [n]$.
- For two vectors $x, y \in \mathbb{R}^n$, we denote $\langle x, y \rangle = \sum_{i=1}^n x_i y_i$ for $i \in [n]$.
- Given two vectors $x, y \in \mathbb{R}^n$, we denote $x \circ y$ as a vector whose i-th entry is $x_i y_i$ for all $i \in [n]$.

540 541	• For a vector $x \in \mathbb{R}^n$, diag $(x) \in \mathbb{R}^{n \times n}$ is defined as a diagonal matrix with its diagonal
542	entries given by $\operatorname{diag}(x)_{i,i} = x_i$ for $i = 1,, n$, and all off-diagonal entries are 0.
543	• We use $\operatorname{erf}(x)$ to denote the error function $\operatorname{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x \exp(-t^2) dt$, and erf^{-1} is
544	denoted as the inverse function of $\operatorname{erf}(x)$.
545	• For any matrix $A \in \mathbb{R}^{m \times n}$, we use A^{\top} to denote its transpose, we use $ A _F$ to
546	denote the Frobenius norm and $ A _{\infty}$ to denote its infinity norm, i.e., $ A _F :=$
547	$(\sum_{i \in [m]} \sum_{i \in [n]} A_{i,i}^2)^{1/2}$ and $ A _{\infty} = \max_{i \in [m], i \in [n]} A_{i,i} .$
548	• For $\mu \in \mathbb{R}$ we use $\mathcal{N}(\mu \in \mathbb{R}^2)$ to denote Gaussian distribution with expectation of μ and
549	variance of σ^2 .
550	 For a mean vector μ ∈ ℝ^d and a covariance matrix Σ ∈ ℝ^{d×d}, we use N(μ, Σ²) to denote the vector Gaussian distribution.
552	• We use $\mathbb{E}[\cdot]$ to denote the expectation and $\operatorname{Var}[\cdot]$ to denote the variance.
554	• We use $\Gamma(x)$ to denote the gamma function where $\Gamma(x) = \int_{-\infty}^{\infty} t^{z-1} \exp(-t) dt$
555	We use $f(x)$ to denote the gamma reflection where $f(x) = \int_0^\infty t^{-1} \exp(-t) dt$.
556	• For an integer $k > 0$, we use χ_k^2 to denote the Chi-squared distribution with k degrees of freedom
557	
558	• Usually, we use $C \ge 1$ to denote a sufficient large constant.
559	
560	D.2 DASIC FACT FOR SOFTMAX
561	Fact B.1. For a vector $x \in \mathbb{R}^d$ and a scalar $b \in \mathbb{R}$, we have:
562	$\operatorname{softmax}(r) = \operatorname{softmax}(r + h \cdot 1_{t})$
563	$\operatorname{Solution}(w) = \operatorname{Solution}(w + o - 1_a)$
564 565	B.3 BASIC FACTS FOR CALCULATION
566	Fact B.2. For $a, b \ge 1$ and there exist a constant $C \ge 0$ such that
568	$\sqrt{a} + \sqrt{b} \le C\sqrt{a+b}$
569	Fact B 3 For a sufficient large $r \in \mathbb{R}$ $(r > 55)$ we have
570	$\frac{1}{1-\sqrt{1-\sqrt{1-\sqrt{1-\sqrt{1-\sqrt{1-\sqrt{1-\sqrt{1-\sqrt{1-\sqrt{1$
571	$\exp(\sqrt{\log(x)}) \le \sqrt{x}$
572 573	B.4 PROBABILITY TOOLS
574	
575	Here, we state a probability toolkit in the following, including several helpful lemmas we'd like to
576	use. Firstly, we provide the feminia about Chernon bound in Chernon (1952) below. L = D \mathbf{A} (Cl = \mathbb{S}^{1} = 1 Cl = \mathbb{S} (1052)) \mathbf{A} + \mathbf{X} = $\sum_{n=1}^{n}$ X = \mathbf{A} - \mathbf{X} = 1 · ··· \mathbf{A} = \mathbf{A} · ··· \mathbf{A}
577	Lemma B.4 (Chernoff bound, Chernoff (1952)). Let $X = \sum_{i=1} X_i$, where $X_i = 1$ with probability
578	p_i and $X_i = 0$ with probability $1 - p_i$, and all X_i are independent. Let $\mu = \mathbb{E}[X] = \sum_{i=1} p_i$. Then
579	• $\Pr[X > (1+\delta)\mu] < \exp(-\delta^2 \mu/3), \forall \delta > 0;$
580	$ \sum_{i=1}^{n} \left[\frac{1}{2} + \frac{1}{2} \right] = \sum_{i=1}^{n} \left(\frac{1}{2} + \frac{1}{2} \right) \left(\frac{1}{2} + \frac{1}{2} \right) $
581	• $\Pr[X \le (1-\delta)\mu] \le \exp(-\delta^2\mu/1), \forall 0 < \delta < 1.$
582	Next, we offer the lemma about Hoeffding bound as in Hoeffding (1994).
583	Lemma B.5 (Hoeffding bound, Hoeffding (1994)). Let X_1, \dots, X_n denote n independent bounded
585	variables in $[a_i, b_i]$ for $a_i, b_i \in \mathbb{R}$. Let $X := \sum_{i=1}^{n} X_i$, then we have
586	$2t^2$
587	$\Pr[X - \mathbb{E}[X] \ge t] \le 2\exp(-\frac{2\pi}{\sum^n (b - a_1)^2})$
588	$\angle_i = 1(o_i - u_i)$
589	We show the lemma of Bernstein inequality as Bernstein (1924).

Lemma B.6 (Bernstein inequality, Bernstein (1924)). Let X_1, \dots, X_n denote *n* independent zeromean random variables. Suppose $|X_i| \leq M$ almost surely for all *i*. Then, for all positive *t*,

$$\Pr[\sum_{i=1}^{n} X_i \ge t] \le \exp(-\frac{t^2/2}{\sum_{j=1}^{n} \mathbb{E}[X_j^2] + Mt/3})$$

592

Then, we give the Khintchine's inequality in Khintchine (1923); Haagerup (1981) as follows:

Lemma B.7 (Khintchine's inequality, Khintchine (1923); Haagerup (1981)). Let $\sigma_1, \dots, \sigma_n$ be *i.i.d* sign random variables, and let $z_1 \cdots, z_n$ be real numbers. Then there are constants C > 0 so that for all t > 0

$$\Pr[|\sum_{i=1}^{n} z_i \sigma_i| \ge t ||z||_2] \le \exp(-Ct^2)$$

We give Hason-wright inequality from Hanson & Wright (1971); Rudelson & Vershynin (2013) below.

Lemma B.8 (Hason-wright inequality, Hanson & Wright (1971); Rudelson & Vershynin (2013)). Let $x \in \mathbb{R}^n$ denote a random vector with independent entries x_i with $\mathbb{E}[x_i] = 0$ and $|x_i| \leq K$ Let A be an $n \times n$ matrix. Then, for every $t \ge 0$

$$\Pr[|x^{\top}Ax - \mathbb{E}[x^{\top}Ax]| > t] \le 2\exp(-c\min\{t^2/(K^4 ||A||_F^2), t/(K^2 ||A||)\})$$

We state Lemma 1 on page 1325 of Laurent and Massart Laurent & Massart (2000).

Lemma B.9 (Lemma 1 on page 1325 of Laurent and Massart, Laurent & Massart (2000)). Let $X \sim \mathcal{X}_k^2$ be a chi-squared distributed random variable with k degrees of freedom. Each one has zero mean and σ^2 variance. Then

$$\Pr[X - k\sigma^2 \ge (2\sqrt{kt} + 2t)\sigma^2] \le \exp(-t)$$
$$\Pr[X - k\sigma^2 \ge 2\sqrt{kt}\sigma^2] \le \exp(-t)$$

Here, we provide a tail bound for sub-exponential distribution Foss et al. (2011).

Lemma B.10 (Tail bound for sub-exponential distribution, Foss et al. (2011)). We say $X \in SE(\sigma^2, \alpha)$ with parameters $\sigma > 0$, $\alpha > 0$, if

$$\mathbb{E}[e^{\lambda X}] \le \exp(\lambda^2 \sigma^2/2), \forall |\lambda| < 1/\alpha$$

Let $X \in SE(\sigma^2, \alpha)$ and $\mathbb{E}[X] = \mu$, then:

$$\Pr[|X - \mu| \ge t] \le \exp(-0.5\min\{t^2/\sigma^2, t/\alpha\})$$

In the following, we show the helpful lemma of matrix Chernoff bound as in Tropp (2011); Lu et al. (2013).

Lemma B.11 (Matrix Chernoff bound, Tropp (2011); Lu et al. (2013)). Let X be a finite set of positive-semidefinite matrices with dimension $d \times d$, and suppose that

$$\max_{X \in \mathcal{X}} \lambda_{\max}(X) \le B.$$

Sample $\{X_1, \dots, X_n\}$ uniformly at random from \mathcal{X} without replacement. We define μ_{\min} and μ_{\max} as follows:

$$\mu_{\min} := n \cdot \lambda_{\min}(\underset{X \in \mathcal{X}}{\mathbb{E}}(X))$$
$$\mu_{\max} := n \cdot \lambda_{\max}(\underset{X \in \mathcal{X}}{\mathbb{E}}(X)).$$

Then

$$\Pr[\lambda_{\min}(\sum_{i=1}^{n} X_i) \le (1-\delta)\mu_{\min}] \le d \cdot \exp(-\delta^2 \mu_{\min}/B) \text{ for } \delta \in (0,1],$$

$$\Pr[\lambda_{\max}(\sum_{i=1}^{n} X_i) \ge (1+\delta)\mu_{\max}] \le d \cdot \exp(-\delta^2 \mu_{\max}/(4B)) \text{ for } \delta \ge 0.$$

648 **PROBLEM DEFINITIONS** С 649 650 C.1 INPUT ASSUMPTION 651 **Definition C.1.** We consider for any input matrix to an attention network $X \in \mathbb{R}^{n \times d}$ where integer 652 *n* denotes the input length and *d* denotes the dimension. We assume: 653 654 • Independent Entries. For any two entries X_{i_1,j_1} and X_{i_2,j_2} in matrix $X, \forall i_1, i_2 \in [n]$ and 655 $j_1, j_2 \in [d]$, they are independent. 656 657 • **Bounded Entries.** For failure probability $\delta \in (0, 0.1)$. With a probability $1 - \delta$, the entry 658 $X_{i,j}$ in matrix $X, \forall i \in [n]$ and $j \in [d]$, we have $|X_{i,j}| \leq B$ for some positive constant B > 0.659 660 661 C.2 ATTENTION COMPUTATION 662 **Definition C.2** (Attention computation). If the following conditions hold: 663 664 • Let $W_Q, W_K, W_V \in \mathbb{R}^{d \times d}$ be denoted as Query, Key and Value projection matrices of 665 attention. 666 • Given an input $X \in \mathbb{R}^{n \times d}$ that holds properties in Definition C.1. 667 668 • Define Query, Key and Value states matrices $Q := XW_Q, K := XW_K, V := XW_V \in$ 669 $\mathbb{R}^{n \times d}$. 670 • $A := \exp(QK^{\top}/\sqrt{d}) \in \mathbb{R}^{n \times n}.$ 671 672 • $D := \operatorname{diag}(A\mathbf{1}_n) \in \mathbb{R}^{n \times n}$. 673 674 Then we have attention computation $Attn(Q, K, V) \in \mathbb{R}^{n \times d}$ as follows: 675 $\operatorname{Attn}(Q, K, V) := D^{-1}AV$ 676 677 C.3 ϵ -Approximated k-Sparse Softmax Vector 678 679 In order to describe the sparsity of the softmax, we define the following notation. 680 **Definition C.3.** For a vector $u \in \mathbb{R}^n$ and $\epsilon \ge 0$, we define sparse set $S_{\epsilon}(u)$ as follows: 681 682 $\mathcal{S}_{\epsilon}(u) := \left\{ i \in [n] \mid |u_i| \le \epsilon \right\}$ 683 684 **Definition C.4** ((ϵ , k)-sparsity). For a vector $u \in \mathbb{R}^n$, we say u is (ϵ , k)-sparse if for a constant 685 $\epsilon \in (0, 1)$, it holds that 686 $|\mathcal{S}_{\epsilon}(u)| > n - k.$ 687 688 C.4 SPARSE ATTENTION 689 690 **Definition C.5.** For a vector $u \in \mathbb{R}^n$. Given a sparsity integer k. Let S_{ϵ} be defined as Definition C.3 691 for some error $\epsilon > 0$. We define the top-k set $\mathcal{T}_k(u) := \{i \in [n] \mid S_{u_i}(u) \ge n - k\}$. 692 **Definition C.6.** If the following conditions hold: 693 • For a vector $u \in \mathbb{R}^n$. 694 • *Given a sparsity integer k.* 696 697 • Let S_{ϵ} be defined as Definition C.3 for some error $\epsilon > 0$. 698 • Denote a top-k set $\mathcal{T}_k(u) := \{i \in [n] \mid S_{u_i}(u) \ge n-k\}$ as Definition C.5. 699 700 Then we define 701 n

$$\mathsf{topk}(u) := [u_i \cdot \mathbf{1}_{i \in \mathcal{T}_k(u)}]_{i \in [n]} \in \mathbb{R}$$

702 703	We define the sparse attention computation as follows:		
704	Definition C.7. If the following conditions hold:		
705 706	• Let $W_Q, W_K, W_V \in \mathbb{R}^{d \times d}$ be denoted as Query, Key and Value projection matrices of attention.		
707 708	• Given an input $X \in \mathbb{R}^{n \times d}$ that holds properties in Definition C.1.		
709 710 711	• Define Query, Key and Value states matrices $Q := XZW_Q, K := XW_K, V := XW_V \in \mathbb{R}^{n \times d}$.		
712	• $A := \exp(QK^{\top}/\sqrt{d}) \in \mathbb{R}^{n \times n}.$		
713 714	• Let topk be defined as Definition C.6.		
715 716	• Define $A_{\text{spar}} := [\text{topk}(A_1), \text{ topk}(A_2), \cdots, \text{ topk}(A_n)]^{\top} \in \mathbb{R}^{n \times n}.$		
717	• $D_{\text{spar}} := \text{diag}(A_{\text{spar}} 1_n) \in \mathbb{R}^{n \times n}.$		
718 719	• $\delta \in (0, 0.1).$		
720 721	• Let $R \ge 0$ be defined as Definition C.9.		
722	The sparse attention computation $\text{SparseAttn}(Q, K, V) \in \mathbb{R}^{n \times d}$ is given by:		
723 724	$SparseAttn(Q,K,V) := D_{\mathrm{spar}}^{-1} A_{\mathrm{spar}} V \in \mathbb{R}^{n \times d}$		
725 726	C.5 HELPFUL DEFINITIONS		
727 728 729	We introduce the following algebraic lemmas to be used later. Definition C.8. Let Query and Key projection matrices $W_Q, W_K \in \mathbb{R}^{d \times d}$ be defined as Definition C.2. We define		
731	$W := W_Q W_K^\top / \sqrt{d}.$		
732 733	Definition C.9. If the following conditions hold:		
734 735	• Let $W \in \mathbb{R}^{d \times d}$ be define as Definition C.8.		
736 737	Then for any $i \in [n]$, we define:		
738 739	$R := B^2 \ W\ _F.$		
740 741	D ATTENTION SPARSITY		
742 743	D.1 MAIN RESULT 1: ATTENTION SPARSITY WITH UPPER BOUND ON ERROR		
744	Theorem D.1. If the following conditions hold:		
745 746 747	• Let $W_Q, W_K, W_V \in \mathbb{R}^{d \times d}$ be denoted as Query, Key and Value projection matrices of attention.		
748	• Given an input $X \in \mathbb{R}^{n \times d}$ that holds properties in Definition C.1.		
749 750 751	• Define Query, Key and Value states matrices $Q := XW_Q, K := XW_K, V := XW_V \in \mathbb{R}^{n \times d}$.		
752 753	• $A := \exp(QK^{\top}/\sqrt{d}) \in \mathbb{R}^{n \times n}.$		
754	• $D := \operatorname{diag}(A1_n) \in \mathbb{R}^{n \times n}$.		
755	• Denote $\beta_i := Bd \max_{j_1 \in [d]} \mathbb{E}[\sum_{j_2=1}^d W_{j_1, j_2} \cdot X_{i, j_2}] .$		

756	• Define $\Gamma := [\beta_1 \cdot 1_n, \beta_2 \cdot 1_n, \cdots, \beta_n \cdot 1_n]^\top \in \mathbb{R}^{n \times n}$
758	• $\widetilde{A} := \exp(QK^{\top}/\sqrt{d} + \Gamma) \in \mathbb{R}^{n \times n}$.
759	
760	• $D := \operatorname{diag}(A1_n) \in \mathbb{R}^{n \times n}.$
761	• $\delta \in (0, 0.1).$
763	• Let $R > 0$ be defined as Definition C.9.
764	
765	• Given sparsity integer $k \leq n$.
766	• Denote $T := \exp(\sqrt{\log(n(n-k)d/\delta)}).$
767	• Let S_{ϵ} be defined as Definition C.3.
769 770	If we choose $\epsilon \geq \frac{T^{O(R)}}{r}$), then with a probability at least $1 - \delta$, we have
771	
772	$\left S_{\epsilon}(D_{i_{1},i_{1}}^{-1}A_{i_{1}})\right = \left S_{\epsilon}(D_{i_{1},i_{1}}^{-1}A_{i_{1}})\right \ge n - k$
773	
775	<i>Proof.</i> Remark. We re-denote $S_{\epsilon} = S_{\epsilon}(D_{i_1,i_1}^{-1}A_{i_1,*})$ in the statement, and $i_2 \in S_{\epsilon}$.
776	Following Part 1 of Lemma D.3, with a probability at least $1 - \delta_1$, we have
777	$\widetilde{A}_{i_1,i_2} \le \exp(Q(R) \cdot \sqrt{\log((n-k)d/\delta_1)}) \tag{1}$
778	$c_1, c_2 = c_1 \left(c_1 \left(c_1 \right) \right) \left(c_1 \left(c_1 \right) \left(c_1 \left(c_1 \right) \left(c_1 \left(c_1 \right) \right) \left(c_1 \left(c_1 \right) \right) \left(c_1 \left(c_1 \right) \left(c_1 \left(c_1 \right) \right) \left(c_1 \left(c_1 \right) \left(c_1 \left(c_1 \right) \right) \left(c_1 \left(c_1 \right) \left(c_1 \left(c_1 \right) \right) \left(c$
779	Following Part 3 of Lemma D.3, with a probability at least $1 - \delta_2$, we have
781	$\widetilde{D}_{i,-i}^{-1} < \exp(O(R) \cdot \sqrt{\log(nd/\delta_2)})/n \tag{2}$
782	$a_{1},a_{1} = 1 \cdot (1 \cdot 1 \cdot$
783	Now we combine Eq. (1) and Eq. (2), with a probability at least $1 - \delta_1 - \delta_2$, we have
784 785	$D_{i_1,i_1}^{-1}A_{i_1,i_2} = \widetilde{D}_{i_1,i_1}^{-1}\widetilde{A}_{i_1,i_2}$
786	$\leq \exp(O(R) \cdot \sqrt{\log((n-k)d/\delta_1)} + O(R) \cdot \sqrt{\log(nd/\delta_2)})/n$
787	$\leq \exp(O(R) \cdot \sqrt{\log((n-k)d/\delta_1) + \log(nd/\delta_2)})/n$
788	$\leq \exp(O(R) \cdot \sqrt{\log((n-\kappa)a/b_1) + \log(na/b_2))/n}$
789	$\leq \exp(O(R) \cdot \sqrt{\log(n(n-k)d^2/(\delta_1\delta_2))})/n$
791	$\leq \exp(O(R) \cdot \sqrt{\log(n(n-k)d/\delta)})/n$
792	$\leq \exp(\sqrt{\log(n(n-k)/\delta)})^{O(R)}/n$
793	$\leq T^{O(R)} \cdot n^{-1}$
794	where the first step follows from Fact B 1, the second step follows from Eq. (1) and Eq. (2) , the third
795	step follows from Fact B.2, the fourth step follows from simple algebras, the fifth step follows from
790	choosing $\delta_1 = \delta_2 = \delta/2$, the sixth step follows from simple algebras, the last step follows from the
798	definition of T .
799	D 2 MAIN DESULT 2. ATTENTION COLLARSE
800	D.2 MAIN RESULT 2: ATTENTION COLLAPSE
801	Theorem D.2. If the following conditions hold:

- Let $W_Q, W_K, W_V \in \mathbb{R}^{d \times d}$ be denoted as Query, Key and Value projection matrices of attention.
- Given an input $X \in \mathbb{R}^{n \times d}$ that holds properties in Definition C.1.
- Define Query, Key and Value states matrices $Q := XW_Q, K := XW_K, V := XW_V \in$ $\mathbb{R}^{n \times d}$.
- $A := \exp(QK^{\top}/\sqrt{d}) \in \mathbb{R}^{n \times n}.$

803 804 805

806

807

810 811	• $D := \operatorname{diag}(A1_n) \in \mathbb{R}^{n \times n}$.
812	• Denote $\beta_i := Bd \max_{i \in [d]} \mathbb{E}[\sum_{i=1}^d W_{i}, i \in X_i, j] _{i \in [d]}$
813	Denote $p_i := Da \max_{j_1 \in [a]} \mathbb{E}[\angle_{j_2=1} \cup j_1, j_2 \cup j_1, j_2] $
814	• Define $\Gamma := [\beta_1 \cdot 1_n, \beta_2 \cdot 1_n, \cdots, \beta_n \cdot 1_n]^\top \in \mathbb{R}^{n \times n}$
815 816	• $\widetilde{A} := \exp(QK^{\top}/\sqrt{d} + \Gamma) \in \mathbb{R}^{n \times n}.$
817	• $\widetilde{D} := \operatorname{diag}(\widetilde{A}1_{+}) \in \mathbb{R}^{n \times n}$
818	$D := \operatorname{diag}(\operatorname{rig}_n) \subset \operatorname{rig}_n$
819	• $\delta \in (0, 0.1).$
820 821	• Let $R \ge 0$ be defined as Definition C.9.
822	• Assuming $R = o(\sqrt{\log(n)})$.
823 824	• Given sparsity integer $k \leq n$.
825	• Denote $T := \exp(\sqrt{\log(n(n-k)d/\delta)})$.
020 827	
828	• Let S_{ϵ} be defined as Definition C.3.
829	For any $\epsilon > 0$, with probability at least $1 - \delta$, we have:
830	$\lim_{n \to \infty} S(D^{-1}A) = n - 1$
831	$\lim_{n \to +\infty} \mathcal{C}_{\epsilon}(\mathcal{D}_{i,i} \Omega_i) = n - 1$
832	Dur of the anders to share a light the state of a second state of a second to second the second state of t
833	<i>Proof.</i> In order to choose κ that meets the ϵ -approximated sparsity, we have:
834	$\epsilon \ge \exp\left(O(R) \cdot \sqrt{\log(n \cdot (n-k)d/\delta)}\right)$
836	
837	where this step follows from Theorem D.1.
838	We obtain:
839	$(\cos^2(\epsilon \cdot n)) = \delta$
840 841	$k \le n - \exp\left(O(\frac{-1}{R^2})\right) \cdot \frac{1}{nd} \tag{3}$
842	Hence, we have:
843	$\lim_{i \to \infty} \mathcal{S}_{\epsilon}(D_{i,i}^{-1}A_i) \ge \lim_{i \to \infty} (n-k)$
844	$n \rightarrow +\infty$ $i < i, i < j = n \rightarrow +\infty$
846	$\geq \lim_{n \to +\infty} \exp\left(O(\frac{\log^2(\epsilon \cdot n)}{R^2})\right) \cdot \frac{\delta}{nd}$
847	= n - 1
848	where the first star follows from Definition C_2 the second star follows from $F_2(2)$ the last star
849 850	follows from $R = o(\sqrt{\log(n)})$ and $\langle D_{i,i}^{-1}A_i, 1_n \rangle = 1$, then
851	$ C(D^{-1}A) = 1$
852	$\max \mathcal{S}_{\epsilon}(D_{i,i}A_i) = n - 1.$
853	
854	
855 856	D.3 BOUNDING D^{-1}
857	Lemma D.3. If the following conditions hold:
858	• Lat W W C DdXd be denoted as Quam. Know d Weber moderation of
859 860	• Let $W_Q, W_K, W_V \in \mathbb{R}^{n \times n}$ be denoted as Query, Key and Value projection matrices of attention.
861	• Given an input $X \in \mathbb{R}^{n \times d}$ that holds properties in Definition C.1.
862 863	• Define Query, Key and Value states matrices $Q := XW_Q, K := XW_K, V := XW_V \in \mathbb{R}^{n \times d}$.

864 865	• Denote $\beta_i := Bd \max_{j_1 \in [d]} \mathbb{E}[\sum_{j_2=1}^d W_{j_1, j_2} \cdot X_{i, j_2}] .$
866	• Define $\Gamma := [\beta_1 \cdot 1_n, \beta_2 \cdot 1_n, \cdots, \beta_n \cdot 1_n]^\top \in \mathbb{R}^{n \times n}$
867	$\widetilde{\mathcal{L}} = \left[$
868	• $A := \exp(QK^+/\sqrt{d} + \Gamma) \in \mathbb{R}^{n \times n}.$
869 870	• $\widetilde{D} := \operatorname{diag}(\widetilde{A}1_n) \in \mathbb{R}^{n \times n}.$
871	• $\delta \in (0, 0, 1)$
872	
873	• Let $R \ge 0$ be defined as Definition C.9.
874 875	Then with a probability at least $1 - \delta$, we have
876	• Part 1. For $i_1, i_2 \in [n]$
877	$\frac{1}{2} \left(\frac{1}{2} \left(\frac{1}{2} \right) \right) = \left(\frac{1}{2} \left(\frac{1}{2} \right) \right) = \left(\frac{1}{2} \left(\frac{1}{2} \right) \right)$
878 879	$\exp(-O(R) \cdot \sqrt{\log(d/\delta)}) \le A_{i_1,i_2} \le \exp(O(R) \cdot \sqrt{\log(d/\delta)})$
880	• Part 2. For $i_1 \in [n]$
881	$n \cdot \exp(-O(R) \cdot \sqrt{\log(nd/\delta)}) \le \widetilde{D} \cdot \cdot \le n \cdot \exp(O(R) \cdot \sqrt{\log(nd/\delta)})$
882	$n \exp(-O(n) + \sqrt{\log(na/b)}) \le D_{i_1,i_1} \le n \exp(O(n) + \sqrt{\log(na/b)})$
884	• Part 2. For $i_1 \in [n]$
885	$\widetilde{D}^{-1} \leq \exp(O(R) \cdot \sqrt{\log(nd/\delta)})/n$
886	$D_{i_1,i_1} \leq \exp(O(n) + \sqrt{\log(nn/\delta)})/n$
887	Proof. Proof of Part 1. We have
888	$ \widetilde{A}, \ldots = \exp((OK^{T}), \ldots)$
890	$ A_{i_1,i_2} = \exp((Q \Lambda)_{i_1,i_2})$
891	$\leq \exp(O(R) \cdot \sqrt{\log(d/\delta)})$
892	where the first step follows from the definition of A, the second step follows from Lemma D.4.
893 894	Proof of Part 2. This proof follows from the union bound of Part 1 of this Lemma and the Definition of D
895	Proof of Part 3 This proof follows from the lower bound on D_{1} , and simple algebras \Box
896	The off are 5. This proof follows from the lower bound on D_{i_1,i_1} and simple algebras.
898	D.4 Concentration on QK^{\top}
899	\mathbf{L} amma $\mathbf{D}\mathbf{A}$. If the following conditions hold:
900	Lemma D.4. If the following conditions nota.
901	• Let $W_Q, W_K, W_V \in \mathbb{R}^{d \times d}$ be denoted as Query, Key and Value projection matrices of
903	attention.
904	• Given an input $X \in \mathbb{R}^{n \times d}$ that holds properties in Definition C.1.
905	• Define Ouery, Key and Value states matrices $Q := XW_Q, K := XW_K, V := XW_V \in$
906	$\mathbb{R}^{n \times d}$.
907	• Let $W \in \mathbb{R}^{d \times d}$ be define as Definition C.8.
909	• Denote $C > 1$ as a sufficient lange constant
910	• Denote $C \ge 1$ as a sufficient targe constant.
911 912	• $\delta \in (0, 0.1).$
913	• Let $R \ge 0$ be defined as Definition C.9.
914	• For $i_1, i_2 \in [n]$
915 916	Then with a probability at least $1 - \delta$ we have
917	Then with a probability at least 1 – 0, we have
	$ (QK^{\top})_{i_1,i_2} \le O(R) \cdot \sqrt{\log(d/\delta)} + \beta_{i_2}$

Proof. We have:

$$|(QK^{\top})_{i_{1},i_{2}}| = |(XW_{Q}W_{K}^{\top}X^{\top}/\sqrt{d})_{i_{1},i_{2}}|$$

$$= |(XWX^{\top})_{i_{1},i_{2}}|$$

$$= |X_{i_{1}}^{\top}WX_{i_{2}}|$$

$$= |\sum_{j_{1}=1}^{d}\sum_{j_{2}=1}^{d}W_{j_{1},j_{2}} \cdot X_{i_{1},j_{1}}X_{i_{2},j_{2}}|$$

$$\leq B|\sum_{j_{1}=1}^{d}\sum_{j_{2}=1}^{d}W_{j_{1},j_{2}} \cdot X_{i_{2},j_{2}}|$$

where the first step follows from $Q := XW_Q, K := XW_K$, the second step follows from Definition C.9, the third and fourth steps follow from simple algebras, the fifth step follows from Definition C.1.

We then apply Hoeffding inequality (Lemma B.5) to each $W_{j_1,j_2} \cdot X_{i_2,j_2}$. Hence, with a probability at least $1 - \delta$, we have:

$$\left|\sum_{j_{2}=1}^{d} W_{j_{1},j_{2}} \cdot X_{i_{2},j_{2}} - \mathbb{E}\left[\sum_{j_{2}=1}^{d} W_{j_{1},j_{2}} \cdot X_{i_{2},j_{2}}\right]\right| \le O(B) \cdot \|W_{j_{1}}\|_{2} \cdot \sqrt{\log(d/\delta)}$$

since $|X_{i_2,j_2}| \leq B$ for any $j_2 \in [d]$.

By triangle inequality, we have:

$$\left|\sum_{j_{2}=1}^{d} W_{j_{1},j_{2}} \cdot X_{i_{2},j_{2}}\right| \le O(B) \cdot \|W_{j_{1}}\|_{2} \cdot \sqrt{\log(d/\delta)} + \left|\mathbb{E}\left[\sum_{j_{2}=1}^{d} W_{j_{1},j_{2}} \cdot X_{i_{2},j_{2}}\right]\right|$$
(4)

We obtain:

$$B|\sum_{j_1=1}^d \sum_{j_2=1}^d W_{j_1,j_2} \cdot X_{i_2,j_2}| \le O(B^2) \cdot \|W\|_F \cdot \sqrt{\log(d/\delta)} + Bd \max_{j_1 \in [d]} |\mathbb{E}[\sum_{j_2=1}^d W_{j_1,j_2} \cdot X_{i_2,j_2}]|$$
$$= O(R) \cdot \sqrt{\log(d/\delta)} + \beta_{i_2}$$

where the first step follows from Eq 4 and Fact B.2 ($||W||_F = \sum_{j_1=1}^d ||W_{j_1}||_2^2$), the second step follows from Definition C.9 and define

$$\beta_{i_2} := Bd \max_{j_1 \in [d]} |\mathbb{E}[\sum_{j_2=1}^d W_{j_1, j_2} \cdot X_{i_2, j_2}]|$$

Remark D.5. The formal results of Lemma D.4 in the appendix have slight differences with the informal forms, in which we omit the additional terms of each upper bound since such terms are trivially some constants. Fact B.1 shows any constant bias term added before the softmax function will not change the output. We thus simplify the equations for tighter boundaries and more convenient notation.

E SPARSE ATTENTION APPROXIMATION

E.1 MAIN RESULT 3: UPPER BOUND ON ERROR

Theorem E.1. If the following conditions hold:

- Let Query, Key and Value states matrices $Q, K, V \in \mathbb{R}^{n \times d}$ be defined as Definition C.2.
 - $A := \exp(QK^{\top}/\sqrt{d}) \in \mathbb{R}^{n \times n}.$

• $D := \operatorname{diag}(A\mathbf{1}_n) \in \mathbb{R}^{n \times n}$ • Let \mathcal{T}_k be defined as Definition C.5. • Let topk be defined as Definition C.6. • Let S_{ϵ} be defined as Definition C.3, we omit $S_{\epsilon}(D_{i,i}^{-1}A_{i,*})$ for $i \in [n]$ to $S_{\epsilon,i}$. • Define $A_{\text{spar}} := [\text{topk}(A_{1,*}) \quad \text{topk}(A_{2,*}) \quad \cdots \quad \text{topk}(A_{n,*})]^{\top}$. • $D_{\text{spar}} := \text{diag}(A_{\text{spar}} \mathbf{1}_n) \in \mathbb{R}^{n \times n}.$ • $\delta \in (0, 0.1)$. • Let R be defined as Definition C.9. Then with a probability at least $1 - \delta$, we have • Part 1. Choosing $k = \Omega(n^C)$ for $C \in (0, 1)$, we have $C_{\text{error}} \in (0, C)$: $\|D_{\text{spar}}^{-1}A_{\text{spar}}V - D^{-1}AV\|_{\infty} \le o(n^{-C_{\text{error}}})$ • Part 2. Choosing $k = o(\log(n))$ for $C \in (0, 1)$, we have $C_{\text{error}} \in (0, C)$: $\|D_{\text{spar}}^{-1}A_{\text{spar}}V - D^{-1}AV\|_{\infty} \le \Omega(n^{C_{\text{error}}})$ Proof. We have $\|D_{\text{spar}}^{-1}A_{\text{spar}} - D^{-1}A\|_{\infty} = \|D_{\text{spar}}^{-1}A_{\text{spar}} - D^{-1}A_{\text{spar}} + D^{-1}A_{\text{spar}} - D^{-1}A\|_{\infty}$ $\leq \|D_{\text{spar}}^{-1}A_{\text{spar}} - D^{-1}A_{\text{spar}}\|_{\infty} + \|D^{-1}A_{\text{spar}} - D^{-1}A\|_{\infty}$ $\leq \left(\frac{n-k}{nk} + \frac{1}{n}\right) \cdot \exp(O(R) \cdot \sqrt{\log(nd/\delta)})$ $\leq \frac{1}{k} \cdot \exp(O(R) \cdot \sqrt{\log(nd/\delta)})$ where the first step follows from simple algebras, the second step follows from triangle inequality, the third step follows from Part 1 and Part 4 of Lemma E.3, the fourth step follows from simple algebras. **Part 1.** Choosing $k = \Omega(n^C)$ for $C \in (0, 1)$, we have: $\|D_{\operatorname{spar}}^{-1}A_{\operatorname{spar}} - D^{-1}A\|_{\infty} \le \frac{1}{k} \cdot \exp(O(R) \cdot \sqrt{\log(nd/\delta)})$ $< o(n^{-C_{\rm error}})$ where the first step follows from Eq. (5), the second step follows from $0 < C_{\text{error}} < C$. **Part 2.** Choosing $k = o(\log(n))$ for $C \in (0, 1)$, we have: $\|D_{\text{spar}}^{-1}A_{\text{spar}} - D^{-1}A\|_{\infty} \le \frac{1}{k} \cdot \exp(O(R) \cdot \sqrt{\log(nd/\delta)})$ $\leq \Omega(n^{C_{\text{error}}})$ where the first step follows from Eq. (5), the second step follows from $0 < C_{\text{error}}$. E.2 LOWER BOUND ON ERROR **Theorem E.2.** If the following conditions hold: • Let Query, Key and Value states matrices $Q, K, V \in \mathbb{R}^{n \times d}$ be defined as Definition C.2. • $A := \exp(QK^{\top}/\sqrt{d}) \in \mathbb{R}^{n \times n}$. • $D := \operatorname{diag}(A\mathbf{1}_n) \in \mathbb{R}^{n \times n}$

(5)

1026 • Let \mathcal{T}_k be defined as Definition C.5. 1027 • Let topk be defined as Definition C.6. 1028 1029 • Let S_{ϵ} be defined as Definition C.3, we omit $S_{\epsilon}(D_{i,i}^{-1}A_{i,*})$ for $i \in [n]$ to $S_{\epsilon,i}$. 1030 1031 • Define $A_{\text{spar}} := [\operatorname{topk}(A_{1,*}) \quad \operatorname{topk}(A_{2,*}) \quad \cdots \quad \operatorname{topk}(A_{n,*})]^\top$. 1032 • $D_{\text{spar}} := \text{diag}(A_{\text{spar}} \mathbf{1}_n) \in \mathbb{R}^{n \times n}$. 1034 • $\delta \in (0, 0.1).$ 1035 1036 • Let R be defined as Definition C.9. 1037 • Choosing $k = o(\log(n))$ 1038 1039 Then with a probability at least $1 - \delta$, we have 1040 $\|D_{\rm spar}^{-1}A_{\rm spar} - D^{-1}A\|_F > O(1)$ 1041 1042 Proof. We have: 1043 1044 $\|D_{\text{spar}}^{-1}A_{\text{spar}} - D^{-1}A\|_{F}^{2} = \sum_{i=-1}^{n} \sum_{i=-1}^{n} (D_{\text{spar},i_{1},i_{1}}^{-1}A_{\text{spar},i_{1},i_{2}} - D_{i_{1},i_{1}}^{-1}A_{i_{1},i_{2}})^{2}$ 1045 1046 1047 $\geq \sum_{k=1}^{n} \sum_{k \in \mathcal{T}} \left(\frac{1}{k} \exp(-O(R) \cdot \sqrt{\log(\frac{nd}{\delta})}) - \frac{1}{n} \exp(O(R) \cdot \sqrt{\log(\frac{nd}{\delta})})\right)^2$ 1048 1049 $+\sum_{i_{n}=1}^{n}\sum_{i_{n}\in[n]/\mathcal{T}}\left(\frac{1}{n}\exp(-O(R)\cdot\sqrt{\log(\frac{nd}{\delta})}\right)\right)^{2}$ 1051 1052 1053 $\geq \sum_{i_1=1}^{n} \sum_{i_2 \in \mathcal{T}_{k}} (O(\frac{1}{\sqrt{n \cdot o(\log(n))}}))^2$ 1054 1055 > O(1)1056 1057 where the first step follows from the definition of Frobenius norm ℓ_F , the second step follows from 1058 plugging $k = o(\sqrt{\log(n)})$ and we have: 1059 $\frac{\mathrm{d}}{\mathrm{d}n} \Big(\frac{1}{k} \exp(-O(R) \cdot \sqrt{\log(\frac{nd}{\delta})}) - \frac{1}{n} \exp(O(R) \cdot \sqrt{\log(\frac{nd}{\delta})}) \Big) \ge \frac{\mathrm{d}}{\mathrm{d}n} \frac{1}{\sqrt{n \cdot o(\log(n))}},$ 1060 1061 1062 and the last step follows from simple algebras. 1063 1064 E.3 **APPROXIMATING SOFTMAX FUNCTION** Lemma E.3. If the following conditions hold: 1068 • Let Query, Key and Value states matrices $Q, K, V \in \mathbb{R}^{n \times d}$ be defined as Definition C.2. 1069 • $A := \exp(QK^{\top}/\sqrt{d}) \in \mathbb{R}^{n \times n}$. 1070 1071 • $D := \operatorname{diag}(A\mathbf{1}_n) \in \mathbb{R}^{n \times n}$ • Let \mathcal{T}_k be defined as Definition C.5. 1074 • Let topk be defined as Definition C.6. 1075 1076 • Let S_{ϵ} be defined as Definition C.3, we omit $S_{\epsilon}(D_{i,i}^{-1}A_{i,*})$ for $i \in [n]$ to $S_{\epsilon,i}$. 1077 • Define $A_{\text{spar}} := [\text{topk}(A_{1,*}) \quad \text{topk}(A_{2,*}) \quad \cdots \quad \text{topk}(A_{n,*})]^\top$. 1078 1079 • $D_{\text{spar}} := \text{diag}(A_{\text{spar}} \mathbf{1}_n) \in \mathbb{R}^{n \times n}$.

• $\delta \in (0, 0.1)$. 1081 • Let $R \ge 0$ be defined as Definition C.9. 1082 1083 Then with a probability at least $1 - \delta$, we have 1084 • Part 1. $\|D^{-1}A_{\text{spar}} - D^{-1}A\|_{\infty} \le \frac{1}{n} \cdot \exp(O(R) \cdot \sqrt{\log(nd/\delta)})$ 1087 1088 1089 • Part 2. 1090 $||D_{\text{spar}} - D||_{\infty} \le (n-k) \cdot \exp(O(R) \cdot \sqrt{\log(nd/\delta)})$ 1091 • Part 3. 1093 $\|D_{\operatorname{spar}}^{-1} - D^{-1}\|_{\infty} \le \frac{n-k}{nk} \cdot \exp(O(R) \cdot \sqrt{\log(nd/\delta)})$ 1094 1095 • Part 4. $\|D_{\text{spar}}^{-1}A_{\text{spar}} - D^{-1}A_{\text{spar}}\|_{\infty} \le \frac{n-k}{n^k} \cdot \exp(O(R) \cdot \sqrt{\log(nd/\delta)})$ 1099 1100 *Proof.* Before we begin the proof, we construct a toolkit as follows: 1101 1102 For $x_1 < \exp(\sqrt{\log(a/\delta_1)})$ and $x_2 < \exp(\sqrt{\log(a/\delta_2)})$, we have 1103 $x_1 x_2 \leq \exp(\sqrt{\log(a/\delta_1)}) \cdot \exp(\sqrt{\log(b/\delta_2)})$ 1104 $< \exp(\sqrt{\log(a/\delta_1)} + \sqrt{\log(b/\delta_2)})$ 1105 1106 $< \exp(C\sqrt{\log(a/\delta_1) + \log(b/\delta_2)})$ 1107 $< \exp(C_{\sqrt{\log(ab/\delta)}})$ 1108 (6)1109 where these steps follow from simple algebras, Fact B.3 and choose $\delta_1 = \delta_2 = \delta/2$. 1110 **Proof of Part 1.** This proof follows from Theorem D.1 and $n - k \le n$. 1111 1112 Proof of Part 2. We have 1113 $\|D_{\text{spar}} - D\|_{\infty} = \|A_{\text{spar}}\mathbf{1}_n - A\mathbf{1}_n\|_{\infty}$ 1114 $= \|D \circ (D^{-1}A_{\operatorname{spar}}\mathbf{1}_n - D^{-1}A\mathbf{1}_n)\|_{\infty}$ 1115 1116 $\leq \|D\|_{\infty} \cdot \|D^{-1}A_{\operatorname{spar}}\mathbf{1}_n - D^{-1}A\mathbf{1}_n\|_{\infty}$ 1117 $\leq \|D\|_{\infty} \cdot \frac{n-k}{n} \cdot \exp(O(R) \cdot \sqrt{\log(nd/\delta)})$ 1118 1119 $\leq (n-k) \cdot \exp(O(R) \cdot \sqrt{\log(nd/\delta)}) \cdot \exp(O(R) \cdot \sqrt{\log(nd/\delta)})$ 1120 1121 $\leq (n-k) \cdot \exp(O(R) \cdot \sqrt{\log(nd/\delta)})$ 1122 where the first step follows from the definitions of D_{spar} and D, the second step follows from simple 1123 algebras, the third step follows from Cauchy-Schwarz inequality, the fourth step follows from Part 1 1124 of this Lemma and the definition of A_{spar} , the fifth step follows from Part 4 of Lemma E.4, the sixth 1125 step follows from Eq. (6). 1126 Proof of Part 3. We have 1127 $||D_{\text{spar}}^{-1} - D^{-1}||_{\infty} = ||D_{\text{spar}}^{-1}||_{\infty} \cdot ||D^{-1}||_{\infty} \cdot ||D_{\text{spar}} - D||_{\infty}$ 1128 1129 $\leq \|D_{\text{spar}}^{-1}\|_{\infty} \cdot \|D^{-1}\|_{\infty} \cdot (n-k) \cdot \exp(O(R) \cdot \sqrt{\log(nd/\delta)})$ 1130 1131 $\leq \frac{n-k}{nk} \cdot \exp(O(R) \cdot \sqrt{\log(nd/\delta)}) \cdot \exp(O(R) \cdot \sqrt{\log(nd/\delta)})$ 1132 $\leq \frac{n-k}{nk} \cdot \exp(O(R) \cdot \sqrt{\log(nd/\delta)})$ 1133

1134 where the first step follows from simple algebras, the second step follows from Part 2 of this Lemma, 1135 the third step follows from Part 5 and Part 6 of Lemma E.4, the last step follows from Eq. (6). 1136 Proof of Part 4. This proof follows from combining Part 2 of Lemma E.4 and Part 3 of this 1137 Lemma. \square 1138 1139 E.4 HELPFUL BOUND TOOLKIT 1140 1141 Lemma E.4. If the following conditions hold: 1142 • Let Query, Key and Value states matrices $Q, K, V \in \mathbb{R}^{n \times d}$ be defined as Definition C.2. 1143 1144 • $A := \exp(QK^{\top}/\sqrt{d}) \in \mathbb{R}^{n \times n}$. 1145 1146 • $D := \operatorname{diag}(A\mathbf{1}_n) \in \mathbb{R}^{n \times n}$ 1147 • Let \mathcal{T}_k be defined as Definition C.5. 1148 1149 • Let topk be defined as Definition C.6. 1150 • Let S_{ϵ} be defined as Definition C.3, we omit $S_{\epsilon}(D_{i,i}^{-1}A_{i,*})$ for $i \in [n]$ to $S_{\epsilon,i}$. 1151 1152 • Define $A_{\text{spar}} := [\text{topk}(A_{1,*}) \quad \text{topk}(A_{2,*}) \quad \cdots \quad \text{topk}(A_{n,*})]^\top$. 1153 1154 • $D_{\text{spar}} := \text{diag}(A_{\text{spar}} \mathbf{1}_n) \in \mathbb{R}^{n \times n}$. 1155 1156 • $\delta \in (0, 0.1)$. 1157 • Let $R \ge 0$ be defined as Definition C.9. 1158 1159 Then with a probability at least $1 - \delta$, we have 1160 1161 • Part 1. $\exp(-O(R) \cdot \sqrt{\log(\frac{n-k}{\delta}d)}) \le A_{i_1,i_2} \le \exp(O(R) \cdot \sqrt{\log(\frac{n-k}{\delta}d)}), \forall i_1 \in [n], i_2 \in \mathbb{N}$ 1162 Se 1163 • Part 2. $\exp(-O(R) \cdot \sqrt{\log(\frac{n}{\delta}d)}) \le A_{i_1,i_2} \le \exp(O(R) \cdot \sqrt{\log(\frac{n}{\delta}d)}), \forall i_1, i_2 \in [n]$ 1164 1165 • Part 3. $k \exp(-O(R) \cdot \sqrt{\log(\frac{n}{\delta}d)}) \leq \sum_{i_2 \in \mathcal{T}_h} A_{i_1,i_2} \leq k \exp(O(R) \cdot \sqrt{\log(\frac{n}{\delta}d)}), \forall i_1 \in [n]$ 1166 1167 • Part 4. $n \exp(-O(R) \cdot \sqrt{\log(\frac{n}{\delta}d)}) \le D_{i_1,i_1} \le n \exp(O(R) \cdot \sqrt{\log(\frac{n}{\delta}d)}), \forall i_1 \in [n]$ 1168 1169 • Part 5. $\frac{1}{k}\exp(-O(R)\cdot\sqrt{\log(\frac{n}{\delta}d)}) \leq \left(\sum_{i_2\in\mathcal{T}_k}A_{i_1,i_2}\right)^{-1} \leq \frac{1}{k}\exp(O(R)\cdot\sqrt{\log(\frac{n}{\delta}d)}),$ 1170 1171 $\forall i_1 \in [n]$ 1172 • Part 6. $\frac{1}{n} \exp(-O(R) \cdot \sqrt{\log(\frac{n}{\delta}d)}) \le D_{i_1,i_1}^{-1} \le \frac{1}{n} \exp(O(R) \cdot \sqrt{\log(\frac{n}{\delta}d)}), \forall i_1 \in [n]$ 1173 1174 Proof. Proof of Part 1. 1175 1176 This proof follows from Eq. (1). 1177 **Proof of Part 2.** 1178 1179 This proof follows from Part 1 and Part 2 of Lemma D.3. 1180 **Proof of Part 3.** 1181 1182 This proof follows from Part 1 of this Lemma. 1183 **Proof of Part 4**. 1184 This proof follows from Part 2 of Lemma D.3. 1185 1186 **Proof of Part 5.** 1187 This proof follows from Part 3 of this Lemma.

1188	Proof of Part 6 . This proof follows from Part 1 of this Lemma	
1189	11001 01 1 at 0. This proof follows from 1 at 4 of this Echinia.	
1190		
1191		
1192		
1193		
1194		
1195		
1196		
1197		
1198		
1199		
1200		
1201		
1202		
1202		
1203		
1204		
1205		
1200		
1207		
1200		
1209		
1210		
1211		
1212		
1213		
1214		
1210		
1210		
1217		
1218		
1219		
1220		
1221		
1222		
1223		
1224		
1225		
1226		
1227		
1228		
1229		
1230		
1231		
1232		
1233		
1234		
1235		
1236		
1237		
1238		
1239		
1240		
1241		