HSR-Enhanced Sparse Attention Acceleration

Bo Chen¹, Yingyu Liang^{2,3}, Zhizhou Sha⁴, Zhenmei Shi³, Zhao Song⁵ ¹Middle Tennessee State University ²The University of Hong Kong ³University of Wisconsin-Madison ⁴Tsinghua University ⁵The Simons Institute for the Theory of Computing at the University of California, Berkeley, bc7b@mtmail.mtsu.edu, yingyul@hku.hk, shazz20@mails.tsinghua.edu.cn, zhmeishi@cs.wisc.edu, magic.linuxkde@gmail.com

Large Language Models (LLMs) have demonstrated remarkable capabilities across various applications, but their performance on long-context tasks is often limited by the computational complexity of attention mechanisms. We introduce a novel approach to accelerate attention computation in LLMs, particularly for longcontext scenarios. We leverage the inherent sparsity within attention mechanisms, both in conventional Softmax attention and ReLU attention (with ReLU^{α} activation, $\alpha \in \mathbb{N}_+$), to significantly reduce the running time complexity. Our method employs a Half-Space Reporting (HSR) data structure to identify non-zero or "massively activated" entries in the attention matrix. We present theoretical analyses for two key scenarios: generation decoding and prompt prefilling. Our approach achieves a running time of $O(mn^{4/5})$ significantly faster than the naive approach O(mn) for generation decoding, where n is the context length, m is the query length, and d is the hidden dimension. We can also reduce the running time for prompt prefilling from O(mn) to $O(mn^{1-1/\lfloor d/2 \rfloor} + mn^{4/5})$. Our method introduces only provably negligible error for Softmax attention. This work represents a significant step towards enabling efficient long-context processing in LLMs.

1. Introduction

Large Language Models (LLMs) have showcased remarkable capabilities across various applications, including context-aware question answering, content generation, summarization, and dialogue systems, among others [1–4]. Long-context tasks of LLMs have gained more and more attention. Several LLMs extend their context length to 128K tokens, such as Yarn [5], GPT-4 [6], Claude 3.5 [7], Llama 3.1 [8], Phi-3.5 [9], Mistral Nemo [10], etc. A bottleneck for long-context tasks is the computational cost of the attention mechanism in LLMs. The key to LLM success is the transformer architecture [11], wildly used in various practical scenarios [12–16], whose critical component is the attention mechanism. Let n be the data length, m be the length of query tokens, and d be the feature dimension¹. The conventional attention uses Softmax activation and is defined as follows:

Definition 1.1 (Softmax attention). Let $Q \in \mathbb{R}^{m \times d}$ and $K, V \in \mathbb{R}^{n \times d}$ denote the query, key, and value *matrix*. The Softmax attention is:

$$\mathsf{Attn}_s(Q, K, V) := \mathsf{Softmax}(QK^\top / \sqrt{d}) V = D^{-1} A_s V \in \mathbb{R}^{m \times d},$$

where (1) $A_s := \exp(QK^{\top}/\sqrt{d}) \in \mathbb{R}^{m \times n}$ and exp is applied element-wise , (2) $D := \operatorname{diag}(A_s \cdot \mathbf{1}_n) \in \mathbb{R}^{m \times m}$ denotes the normalization matrix, (3) $D^{-1}A_s \in \mathbb{R}^{m \times n}$ denotes the attention matrix.

In practical LLM applications, there are two scenarios for attention computation depending on the context length n and query length m. The first case, $m = \Theta(1)$, represents the generation decoding based on the pre-computed Key Value Cache (KV), which stores the intermediate attention key and value matrices. The second case, $m = \Theta(n)$, represents the prompt prefilling before text generation

¹As *d* is always fixed in practice, there is no need to scale up *d* in analysis. Thus, in this work, we always assume *d* is a small constant.

or the cross-attention computation. However, in both cases, when the context window n becomes larger, the running time will increase correspondingly, i.e., it will be linear and quadratic in n for $m = \Theta(1)$ and $m = \Theta(n)$, respectively. Thus, reducing the running time of attention computations with long context input becomes essential to minimize response latency and increase throughput for LLMs.

In this work, we introduce novel methods to reduce the running time complexity for both cases, i.e., $m = \Theta(1)$ and $m = \Theta(n)$. We are inspired by the inherent sparsity within attention mechanisms. Numerous prior studies have highlighted the significant sparsity in the attention matrix [17–21]. This manifestation of sparsity in Softmax attention is that a large number of attention scores, i.e., QK^{\top} , concentrate on a small number of entries, which is known as "massive activation". Due to this nature, Softmax attention can be accelerated by only calculating the entries that contain large attention scores, introducing negligible approximation errors [22, 23].

When talking about ReLU activation, one can easily accelerate the computation process by only calculating the entries activated by ReLU (since other non-activated entries will eventually be set to zero by ReLU). ReLU attention is another attention mechanism widely used, substituting the conventional Softmax activation function with ReLU. ReLU attention has demonstrated performance comparable to Softmax attention in various downstream tasks [24, 25] (see Section 2 for more details). We present the formal definition of ReLU attention as follows.

Definition 1.2 (ReLU attention). Let $Q \in \mathbb{R}^{m \times d}$ and $K, V \in \mathbb{R}^{n \times d}$ denote the query, key, and value *matrix*. Let $\alpha \in \mathbb{N}_+$. The ReLU attention is:

$$\operatorname{Attn}_r(Q, K, V) := D^{-1}A_r V \in \mathbb{R}^{m \times d},$$

where (1) $A_r := \text{ReLU}^{\alpha}(QK^{\top}/\sqrt{d} - b) \in \mathbb{R}^{m \times n}$ and ReLU^{α} denotes the α -th power of ReLU activation for any $\alpha \in \mathbb{N}_+$, (2) $D := \text{diag}(A_r \cdot \mathbf{1}_n) \in \mathbb{R}^{m \times m}$ denotes the normalization matrix, (3) $b \in \mathbb{R}$ denotes position bias, (4) $D^{-1}A_r \in \mathbb{R}^{m \times n}$ denotes the attention matrix.

To expedite the computation, the critical task is to identify the large/nonzero entries for Softmax/ReLU attention, respectively. In this work, We utilize the half-space reporting (HSR) data structure to tackle this problem. HSR was first introduced in [26] to address the half-space range reporting problem (More details can be found in Section 3.2).

In our framework, we define the halfspace as the region where the attention scores (the inner products of key and query vectors) exceed some threshold. We leverage the HSR data structure's ability to quickly answer range reporting to expedite the identification of non-zero entries within the ReLU attention matrix and large entries in Softmax attention. Consequently, we accelerate the computa-



Figure 1: The trending of the Softmax activation (exp) and the ReLU activation with different powers. Here, we choose b = 1.5 as the threshold for the ReLU activation.

tion for ReLU attention and expedite the computation of Softmax attention with negligible approximation error, resulting in a substantial reduction in computation time.

Then, we state our results under the generation decoding scenario $(m = \Theta(1))$ and prompt prefilling scenario $(m = \Theta(n))$. When $m = \Theta(1)$, we accelerate ReLU and Softmax attention computation time over the naive approach from O(mn) to $O(mn^{4/5})$ with pre-processed KV cache (Algorithm 1).

When $m = \Theta(n)$, we accelerate ReLU and Softmax attention computation time over the naive approach from O(mn) to $O(mn^{1-1/\lfloor d/2 \rfloor} + mn^{4/5})$ (Algorithm 2). Furthermore, Section 7 shows that the approximation error associated with Softmax attention utilizing "massive activated" entries only is small in practice, which is consistent with our theoretical analysis.

Our contributions:

- To the best of our knowledge, this is the first work incorporating the HSR data structure with attention computation to reduce the running time complexity with the help of the sparsity within the attention mechanisms.
- We provide rigorous theoretical proofs for reducing the computational time (1) for ReLU attention generation decoding from O(mn) to $O(mn^{4/5})$ (Algorithm 1); (2) for ReLU attention prompt prefilling from O(mn) to $O(mn^{1-1/\lfloor d/2 \rfloor} + mn^{4/5})$ (Algorithm 2).
- We achieve the same running time speed up for the conventional Softmax attention (Theorem 4.2, 5.2), and we give rigorous theoretical proofs to ensure that the approximation error remains negligible (Theorem 4.3). And we provide an empirical evaluation to prove our theoretical error analysis (Section 7).
- We conduct empirical experiments on prominent LLMs to verify the approximation error associated with Softmax attention utilizing "massive activated" entries only. The results show that the error using a few top entries is already insignificant, consistent with our theoretical analysis.

Comparison with Previous Works. [27] uses an approximated nearest neighbors search, which requires that the data are well-conditioned, e.g., uniformly separated, while our algorithm supports an exact nearest neighbor search. Their nearest neighbors search may introduce large approximation errors, while the approximation error of our algorithm can be small. On the other hand, [28] uses low-rank approximation methods to accelerate the attention computation while they require bounded entries assumptions, which is not required in this work. The bounded entry assumption may not always be held in practical scenarios. Our work is based on practical observation that the attention matrix is sparse.

Roadmap. Section 2 presents related work. Section 3 introduces essential concepts. Section 4 presents our main results, i.e., guarantees on run time reduction and approximation error. Section 5 introduces the extension of our method on prompt prefilling scenarios. Section 6 provides a summary of the techniques used in our proof. Section 7 provides our empirical evaluation for the approximation error on Softmax attention. Section 8 discusses the potential of extending our method. Section 9 concludes our algorithm and contributions.

2. Related Work

2.1. Attention Acceleration for Long Context Input

A long context window is essential for transformer-based LLMs in many downstream tasks. However, due to the quadratic time complexity associated with attention mechanisms, transformers are usually hard to run inference efficiently. Numerous methods have been proposed to enhance the inference efficiency. One approach involves using alternative architectures as proxies for attention to support faster inference, such as Mamba [29, 30], PolySketchFormer [31], Hopfield Models [32–38] and Linearizing Transformers [39, 40]. Another line of research focuses on approximating attention matrix computation [28, 41–54]. Nevertheless, these methods often rely on assumptions that may not be practical. For instance, some approaches use polynomial methods to approximate the exponential function, which requires all entries to be bounded by a small constant. However, our HSR-enhanced attention framework is designed based on practical observation and validated by empirical support. These advancements not only improve general model performance but also play a crucial role in enhancing in-context learning capabilities, where models leverage information from the immediate context to perform tasks without fine-tuning. We refer the readers to some other related works [45, 55–64, 64–83].

2.2. ReLU Attention

ReLU attention employs the ReLU activation function in place of the traditional Softmax function for attention computation. Previous studies have highlighted the promising potential of ReLU attention in various domains. From the empirical side, [24] has demonstrated that incorporating ReLU as the activation function in vision transformers enhances performance on downstream tasks. [84] has shown that transformers equipped with ReLU attention outperform those with Softmax attention, particularly when dealing with large key-value memory in machine translation tasks. From the theoretical side, the scale-invariant property of ReLU attention [85] facilitates the scalability of transformer networks. Furthermore, [86–88] have shown that the inherent properties of ReLU attention contribute positively to the learning process of transformer models. Another key advantage is that the ReLU function effectively sets all negative values to zero, allowing us to bypass these non-contributory elements during attention computation and thereby reducing its running time. Omitting these zero and negative entries does not introduce any error into the final output of the ReLU attention mechanism.

2.3. Half-Space Reporting (HSR) Data Structure

The HSR data structure, initially proposed by [26], was developed to address the half-space range reporting problem. The expedited range query capability inherent to HSR has been demonstrated to significantly enhance computational efficiency across a variety of tasks. Studies such as [89] and [90] have applied HSR to facilitate solving general linear programming (LP) problems. Another line of research has highlighted HSR's potential in expediting the training process of contemporary neural networks [91, 92]. There is also a collection of research that concentrates on leveraging HSR for the advancement of solutions to geometric and graphical challenges [93–95].

3. Preliminary

3.1. Notations

We first introduce basic notations used in this paper. For any positive integer n, we use [n] to denote set $\{1, 2, \dots, n\}$. We use Var[] to denote the variance. For two vectors $x \in \mathbb{R}^n$ and $y \in \mathbb{R}^n$, we use $\langle x, y \rangle$ to denote the inner product between x, y. We use $\mathbf{1}_n$ to denote a length-n vector where all the entries are ones. We use $X_{i,j}$ to denote the *i*-row, *j*-th column of $X \in \mathbb{R}^{m \times n}$. We use $||A||_{\infty}$ to denote the ℓ_{∞} norm of a matrix $A \in \mathbb{R}^{n \times d}$, i.e. $||A||_{\infty} := \max_{i \in [n], j \in [d]} |A_{i,j}|$.

3.2. Half-Space Reporting (HSR) Data Structure

Due to the limitation of space, we provide only the core corollary of the HSR data structure here. We refer the readers to Appendix B.4 for more details.

In [26], the author introduces a data structure named HSR to solve the half-space range reporting problem. The interface of the HSR data structure can be summarized as in Algorithm 3. Intuitively, the HSR data structure recursively partitions the set *S* and organizes the points in a tree data structure. Then for a given query (a, b), all k points of *S* with $sgn(\langle a, x \rangle - b) \ge 0$ are reported quickly. Note that the query (a, b) here defines the half-space *H* in Definition B.10. We summarize the running time complexity of the HSR data structure as follows:

Corollary 3.1 (HSR data-structure time complexity [26], informal version of Corollary B.12). Let \mathcal{T}_{init} denote the pre-processing time to build the data structure, \mathcal{T}_{query} denote the time per query, and \mathcal{T}_{update} time per update. Given a set of *n* points in \mathbb{R}^d , the half-space range reporting problem can be solved with the following performances:

- Part 1. $\mathcal{T}_{init}(n,d) = O_d(n\log n), \mathcal{T}_{query}(n,d,k) = O(dn^{1-1/\lfloor d/2 \rfloor} + dk).$
- Part 2. $\mathcal{T}_{init}(n,d) = O(n^{\lfloor d/2 \rfloor}), \mathcal{T}_{query}(n,d,k) = O(d\log(n) + dk).$

4. Main Results on Generation Decoding

In this section, we present our key findings regarding generation decoding, $m = \Theta(1)$, for both ReLU and Softmax attention mechanisms. We reduce the time complexity from a naive O(mn) to $O(mn^{4/5})$. Our method only introduces a negligible approximation error for Softmax attention.



Figure 2: An outline of our principal algorithms. **Top:** Algorithm 1 for generation decoding is depicted, with the key matrix K is fixed. During each inference step, the input query Q interacts with the HSR data structure to get the activated indices set $\tilde{S}_{i,j}$. Then, we can calculate the attention matrix according to $\tilde{S}_{i,j}$. **Bottom:** Algorithm 2 for prompt prefilling is shown, where both the key matrix K and the query matrix Q are variable across iterations. Consequently, the HSR data structure must first be initialized with K, followed by querying it using Q. Finally, according to the activated entries set $\tilde{S}_{i,j}$ reported by the HSR data structure, the attention matrix can be calculated. For more information, please refer to Remark 6.4.

We begin with introducing our result on ReLU attention generation decoding as follows:

Theorem 4.1 (Running time of ReLU attention generation decoding, informal version of Theorem D.2). Let ReLU attention be defined as Definition 1.2. Assume each entry of K conforms Gaussian $\mathcal{N}(0, \sigma_k^2)$, and each entry of Q conforms Gaussian $\mathcal{N}(0, \sigma_q^2)$. Let $\delta \in (0, 1)$ denote the failure probability. Let $\sigma_a = 4 \cdot (1 + d^{-1} \log(m/\delta))^{1/2} \cdot \sigma_q \sigma_k$. Let $b = \sigma_a \cdot \sqrt{0.4 \log n}$. Suppose we have KV Cache $K, V \in \mathbb{R}^{n \times d}$. We want to generate a m length answer, where $n \gg m$. Then, our inference function in Algorithm 1, with probability at least $1 - \delta$, takes $O(mn^{4/5})$ time to generate the answer.

Theorem 4.1 shows that our Algorithm 1 accelerates the running time of ReLU attention generation decoding from naive O(mn) to $O(mn^{4/5})$, which is a significant speed up when the KV Cache is large. We provide an example for the sparsity with different n in Table 1 in the Appendix. The at least $1 - \delta$ success probability originates from the sparsity analysis of ReLU attention (Lemma 6.1), where with probability at least $1 - \delta$, we have the number of non-zero entries of each row of the attention matrix no larger than $n^{4/5}$.

Algorithm 1 Generation decoding

1: data structure GenerationDecoding ⊳ Lemma 6.2 2: members HALFSPACEREPORT HSR ▷ Part 2 of Corollary 3.1 3: 4: $\{K_i\}_{i\in[n]}$ ⊳ Key matrix $V \in \mathbb{R}^{n \times d}$ 5: ▷ Value matrix $b \in \mathbb{R}$ > Threshold of ReLU activation 6: 7: end members procedure INIT $({K_i}_{i \in [n]}, V, n, d)$ 8: ▷ Store necessary matrices 9: $\{K_i\}_{i \in [n]}, V \leftarrow \{K_i\}_{i \in [n]}, V$ 10: $b \leftarrow \sigma_a \cdot \sqrt{0.4 \log n}$ > Init essential parameters and data structure. Lemma 6.1 \triangleright It takes $O(n^{\lfloor d/2 \rfloor})$ time HSR.INIT $({K_i}_{i \in [n]}, n, d)$ 11: 12: end procedure 13: **procedure** INFERENCE $(Q \in \mathbb{R}^{m \times d}, m)$ $A \leftarrow \mathbf{0}_{m \times n}$ 14: for $i = 1 \rightarrow m$ do \triangleright Loop for *m* query vectors 15: $\widetilde{S}_{i,\text{fire}} \leftarrow \text{hsr.Query}(Q_i, b)$ \triangleright It takes $O(n^{4/5})$ time 16: for $j \in \widetilde{S}_{i,\text{fire}}$ do \triangleright Calculate the ReLU attention output according to $\tilde{S}_{i,\text{fire}}$ 17: $A_{i,j} \leftarrow \mathsf{ReLU}^{\alpha}(\langle Q_i, K_j \rangle / \sqrt{d} - b) \text{ or } A_{i,j} \leftarrow \mathsf{Softmax}(\langle Q_i, K_j \rangle / \sqrt{d})$ 18: 19: end for 20: end for return $D^{-1}AV$ 21: 22: end procedure 23: end data structure

Then, we move on to presenting our result on Softmax attention generation decoding. Our results consist of two parts: the improved running time and the approximation error analysis. Firstly, we introduce our result about the improved running time of Softmax attention generation decoding as follows:

Theorem 4.2 (Running time of Softmax attention generation decoding, informal version of Theorem F.1). Let $Q \in \mathbb{R}^{m \times d}$, $K, V \in \mathbb{R}^{n \times d}$ and the Softmax attention Attn_s be defined in Definition 1.1. Let $\operatorname{NN}(r, q, K) \subseteq [n]$ and the Softmax attention with index set Attn_s be defined as Definition B.2. We choose the threshold $b \in \mathbb{R}$ in Algorithm 1 such that $R = \operatorname{NN}(n^{4/5}, q, K)$. Then, we can show that the Softmax attention with index set Attn_s achieves outstanding running time under the Softmax attention generation scenario: Suppose we have KV Cache $K, V \in \mathbb{R}^{n \times d}$. We want to generate a m length answer, where $n \gg m$. Our inference function in Algorithm 1 (replacing ReLU attention with Softmax attention) takes $O(mn^{4/5})$ time to generate the answer.

Theorem 4.2 demonstrates that if we choose the threshold *b* satisfying $R = NN(n^{4/5}, q, K)$, we can achieve a significant running time improvement of the Softmax attention generation.

Indeed, this method introduces an approximation error due to the exclusion of certain entries. Nevertheless, under mild assumptions about the distribution of the attention scores, we demonstrate that this approximation error is indeed negligible. The proof's intuitive explanation lies in the fact that the majority of attention scores are focused on the small subset of entries that we retain. We organize our results as follows:

Theorem 4.3 (Error analysis of Softmax attention with index set, informal version of Theorem G.2). Let $Q \in \mathbb{R}^{m \times d}$, $K, V \in \mathbb{R}^{n \times d}$ and the Softmax attention Attn_s be defined in Definition 1.1. Let $q \in \mathbb{R}^d$ denote a single row of $Q \in \mathbb{R}^{m \times d}$. Let $\gamma \in [0, 1]$, $\beta_1 \ge \beta_2 \ge 0$. Let the index set R and the Softmax attention with index set Attn_s be defined as Definition B.2. Let NN $(r, q, K) \subseteq [n]$ denote the indices of top-r entries of qK. Let $R = NN(n^{\gamma}, q, K) \subseteq [n]$, where $|R| = n^{\gamma}$. Assume the query q and key cache K have $(\gamma, \beta_1, \beta_2)$ massive activation property (Definition B.3). Then, we have

$$\|\widehat{\mathsf{Attn}}_s(q,K,V) - \mathsf{Attn}_s(q,K,V)\|_{\infty} \le \frac{2\|V\|_{\infty}}{n^{\gamma + (\beta_1 - \beta_2) \cdot \|q\|_2 - 1}}.$$

Theorem 4.3 presents the error of Softmax attention with index set is relatively small. Consequently, omitting the remaining less significant entries is a justifiable compromise.

Remark 4.4. With mild assumptions on V, we can have more precious results from Theorem 4.3. For example, if the entries in V conform to subgaussian distribution with constant variance, we have $||V||_{\infty} = O(\log(n))$ with high probability.

5. Extension on Prompt Prefilling

In this section, we extend our results to prompt prefilling scenarios where the number of queries and keys is proportional, i.e., $m = \Theta(n)$. For ReLU attention, we leverage Part 1 result of Corollary 3.1 to accelerate the identification of non-zero entries (activated entries). We introduce our result on ReLU attention as follows:

Algorithm 2 Prompt Prefilling

```
1: data structure PromptPrefilling
                                                                                                                                    ⊳ Lemma 6.3
 2: members
 3:
          HALFSPACEREPORT HSR
                                                                                                                   ▷ Part 1 of Corollary 3.1
 4: end members
 5:
 6: procedure INFERENCE(\{K_i\}_{i \in [n]}, \{Q_r\}_{r \in [m]}, V, n, m, d)
           b \leftarrow \sigma_a \cdot \sqrt{0.4 \log n}.
                                                                                  ▷ Threshold of ReLU activation (Lemma 6.1)
 7:
 8:
           HSR.INIT({K_i}_{i \in [n]}, n, d)
                                                                                                                  \triangleright It takes O(n \log n) time
           A \leftarrow \mathbf{0}_{m \times n}
 9:
           for i = 1 \rightarrow m do
                                                                                                              \triangleright Loop for m query vectors
10:
                                                                                                    \triangleright It takes O(n^{1-1/\lfloor d/2 \rfloor} + \widetilde{k}_i) time.
                S_{i,\text{fire}} \leftarrow \text{hsr.Query}(Q_i, b)
11:
                                                               \triangleright Calculate the ReLU attention output according to \widetilde{S}_{i,\mathrm{fire}}
                for j \in \widetilde{S}_{i, \text{fire}} do
12:
                     A_{i,j} \leftarrow \mathsf{ReLU}^{\alpha}(\langle Q_i, K_j \rangle / \sqrt{d} - b) \text{ or } A_{i,j} \leftarrow \mathsf{Softmax}(\langle Q_i, K_j \rangle / \sqrt{d})
13:
                end for
14:
           end for
15:
           return D^{-1}AV
16:
17: end procedure
18: end data structure
```

Theorem 5.1 (Running time of ReLU attention prompt prefilling, informal version of Theorem C.2). Let ReLU attention be defined as Definition 1.2. Assume each entry of K is from Gaussian $\mathcal{N}(0, \sigma_k^2)$, and each entry of Q is from Gaussian $\mathcal{N}(0, \sigma_q^2)$. Let $\delta \in (0, 1)$ denote the failure probability. Let $\sigma_a = 4 \cdot (1 + d^{-1}\log(m/\delta))^{1/2} \cdot \sigma_q \sigma_k$. Let $b = \sigma_a \cdot \sqrt{0.4 \log n}$. Suppose we have $Q, K, V \in \mathbb{R}^{n \times d}$. There exist an algorithm (Algorithm 2), with probability at least $1 - \delta$, takes $O(n^{2-1/\lfloor d/2 \rfloor} + n^{1+4/5})$ time to compute the full ReLU attention of Q, K, V.

In Theorem 5.1, we improve the running time of ReLU attention prompt prefilling from $O(n^2)$ to $O(n^{2-1/\lfloor d/2 \rfloor} + n^{1+4/5})$, which is a notable uplift of the running time when n is extremely large.

Then, we present our result on Softmax attention. Intuitively, we use the Part 1 result of Corollary 3.1 to identify those "massive activated" entries (top-r indices) within the attention matrix of Softmax attention and calculate the Softmax attention with top-r indices. We organize our results as follows:

Theorem 5.2 (Running time of Softmax attention prompt prefilling, informal version of Theorem F.2). Let $Q \in \mathbb{R}^{m \times d}$, $K, V \in \mathbb{R}^{n \times d}$ and the Softmax attention Attn_s be defined in Definition 1.1. Let $NN(r, q, K) \subseteq [n]$ and the Softmax attention with index set $Attn_s$ be defined as Definition B.2. We choose the threshold $b \in \mathbb{R}$ in Algorithm 2 such that $R = NN(n^{4/5}, q, K)$. Then, we have the Softmax attention with index set $Attn_s$ achieves outstanding running time under Softmax attention prompt prefilling scenario: Suppose we have $m = \Theta(n)$. Algorithm 2 (replacing ReLU attention with Softmax attention) takes $O(n^{2-1/\lfloor d/2 \rfloor} + n^{1+4/5})$ time to compute the full ReLU attention of Q, K, V. Theorem 5.2 demonstrates our $O(n^{2-1/\lfloor d/2 \rfloor} + n^{1+4/5})$ running time on Softmax attention prompt prefilling, which improves from naive running time $O(n^2)$.

6. Technical Overview

Section 6.1 introduces our analysis of the sparsity in the ReLU attention mechanism. Section 6.2 presents our results of two general attention frameworks. Section 6.3 provides our error analysis of Softmax attention with index set. We have shown that with mild assumption on the distribution of attention scores, the error of Softmax attention with index set is negligible.

6.1. Sparsity Analysis of ReLU Attention

Intuitively, the ReLU activation will deactivate some key and query pairs. We introduce the results of employing the concentration inequality to quantitatively analyze the number of non-zero entries.

Lemma 6.1 (Sparsity analysis, informal version of Lemma E.3). Let the ReLU attention be defined as Definition 1.2. Let $Q \in \mathbb{R}^{m \times d}$ and $K, V \in \mathbb{R}^{n \times d}$ be defined as Definition 1.2. Let $b \in \mathbb{R}$ denote the threshold of ReLU activation, as defined in Definition 1.2. For $i \in [m]$, let \tilde{k}_i denote the number of non-zero entries in *i*-th row of $A \in \mathbb{R}^{m \times n}$. Assume each entry of K is from Gaussian $\mathcal{N}(0, \sigma_k^2)$, and each entry of Q is from Gaussian $\mathcal{N}(0, \sigma_q^2)$. Let $\delta \in (0, 1)$ denote the failure probability. Let $\sigma_a = 4 \cdot (1 + d^{-1} \log(m/\delta))^{1/2} \cdot \sigma_q \sigma_k$. Let $b = \sigma_a \cdot \sqrt{0.4 \log n}$. Then, we have, with probability at least $1 - \delta$, for all $i \in [m]$, the number of non-zero entries of the *i*-th row \tilde{k}_i is at most $2n^{4/5}$.

In Lemma 6.1, we use \tilde{k}_i to denote the number of non-zero entries in *i*-th row of attention matrix $A_r \in \mathbb{R}^{m \times n}$. It indicates that if we choose $b = \sigma_a \sqrt{0.4 \log n}$, with high probability, the number of activated (non-zero) entries can be bounded by $O(n^{4/5})$.

6.2. General Attention Frameworks

First, we introduce our general framework for generation decoding. Here, we use the Part 2 result of the HSR data structure. As this framework is designed for the generation decoding task, the key matrix K is fixed in each inference. Therefore, in the INIT procedure, we initialize the HSR data structure with the key matrix K. Then, in each inference, we use the same HSR data structure to answer the query from each row of the query matrix Q. We provide the result of this general generation decoding framework as follows.

Lemma 6.2 (General generation decoding framework, informal version of Lemma D.1). Let $Q \in \mathbb{R}^{m \times d}$ and $K, V \in \mathbb{R}^{n \times d}$ be defined as Definition 1.2. Assume each entry of K is from Gaussian $\mathcal{N}(0, \sigma_k^2)$, and each entry of Q is from Gaussian $\mathcal{N}(0, \sigma_q^2)$. Let $\delta \in (0, 1)$ denote the failure probability. Let $\sigma_a = 4 \cdot (1 + d^{-1} \log(m/\delta))^{1/2} \cdot \sigma_q \sigma_k$. Let $b = \sigma_a \cdot \sqrt{0.4 \log n}$. Let HSR data structure be defined as Part 2 in Corollary 3.1. There exists an algorithm (Algorithm 1), with at least $1 - \delta$ probability, has the following performance:

- **Part 1.** The INIT procedure runs in $O(n^{\lfloor d/2 \rfloor})$ time.
- **Part 2.** For each query, the INFERENCE procedure runs in $O(mn^{4/5})$ time.

The general framework for prompt prefilling is quite different from the previous one. Namely, we choose the Part 1 result of the HSR data structure. Since in each inference, both the query matrix Q and the key matrix K differ from any other inference, we first initialize the HSR data structure with the key matrix K. Then, for each row of the query matrix Q, we query the HSR data structure to find the activated entries.

Lemma 6.3 (General prompt prefilling framework, informal version of Lemma C.1). Let $Q \in \mathbb{R}^{m \times d}$ and $K, V \in \mathbb{R}^{n \times d}$ be defined as Definition 1.2. Assume each entry of K is from Gaussian $\mathcal{N}(0, \sigma_k^2)$, and each entry of Q is from Gaussian $\mathcal{N}(0, \sigma_q^2)$. Let $\delta \in (0, 1)$ denote the failure probability. Let $\sigma_a = 4 \cdot (1 + d^{-1} \log(m/\delta))^{1/2} \cdot \sigma_q \sigma_k$. Let $b = \sigma_a \cdot \sqrt{0.4 \log n}$. Let HSR data structure be defined as Part 1 in Corollary 3.1. *There exists an algorithm (Algorithm 2), with at least* $1 - \delta$ *probability, computes full attention of* Q, K, V *in* $O(mn^{1-1/\lfloor d/2 \rfloor} + mn^{4/5})$ *time.*

Then, we use the following Remark to demonstrate the different intuitions behind Lemma 6.2 (Algorithm 1) and Lemma 6.3 (Algorithm 2).

Remark 6.4. Algorithm 1 is tailored for generation decoding scenarios, where the key matrix K remains constant throughout each inference. Consequently, our optimization efforts are directed at decreasing the time required for individual inferences, which is achieved by adopting Part 2 of Corollary 3.1. In contrast, Algorithm 2 is intended for prompt prefilling, a context in which the key matrix K varies with each inference. Thus, our objective shifts to minimizing the time complexity associated with initializing the HSR data structure, leading us to select Part 1 of Corollary 3.1. For more details, please refer to Figure 2.

6.3. Error Analysis of Softmax Attention with Top-r Indices

Calculating the Softmax attention on the "massive activated" index set will introduce an approximation error. In the following Lemma, we analyze the quantity of this approximation error. Here, we use α to denote the summation of all entries activated by $\exp(x)$ function, and we use $\overline{\alpha}$ to denote the summation of those entries which are excluded from "massive activated" index set. We provide the general error bound of Softmax attention with an index set as follows:

Lemma 6.5 (General error analysis of Softmax attention with index set, informal version of Lemma G.1). Let $Q \in \mathbb{R}^{m \times d}$, $K, V \in \mathbb{R}^{n \times d}$ and the Softmax attention Attn_s be defined in Definition 1.1. Let $q \in \mathbb{R}^d$ denote a single row of $Q \in \mathbb{R}^{m \times d}$. Let $\alpha, \overline{\alpha}$ and Attn_s be defined as Definition B.2. Then, we have

$$\|\mathsf{Attn}_s(q,K,V) - \widehat{\mathsf{Attn}}_s(q,K,V)\|_{\infty} \leq \frac{2\overline{\alpha}}{\alpha} \cdot \|V\|_{\infty}.$$

Note that Lemma 6.5 only provides a general error analysis of Softmax attention with an index set. Under mild assumption on the distribution of attention scores, we show that this error is actually very small. For more details, please refer to Theorem 4.3.

7. Experiments

In this section, we present our empirical results of evaluating three mainstream LLMs with Softmax attention with top-r indices on different r, showing that the results of the experiments are consistent with our theoretical analysis.

Datasets. To estimate the approximation error of the Softmax attention with "massive activation" entries, we conduct experiments on the PaulGrahamEssays datasets from LLMTest-NeedleInAHaystack [96]. Specifically, for each article in the dataset, we first input $2^{15} = 32768$ tokens to the LLMs, then generate 1024 tokens.

Metric. We evaluate the generation quality by the classical perplexity. Perplexity is defined as the exponentiated average negative log-likelihood of a sequence. If we have a tokenized sequence $X = (x_0, x_1, \dots, x_N)$, then the perplexity of X is: $Perplexity(X) = exp(-\frac{1}{N}\sum_{i=1}^{N} \log p_{\theta}(x_i|x_{<i}))$, where $\log p_{\theta}(x_i|x_{<i})$ is the log-likelihood of the *i*-th token conditioned on the preceding tokens. Intuitively, it can be thought of as an evaluation of the model's ability to predict uniformly among the set of specified tokens in a corpus. Importantly, the tokenization procedure has a direct impact on a model's perplexity, which should be taken into consideration when comparing different models.

Models. To demonstrate the generalization of our approximation error bound, we conducted experiments on three mainstream large models: LLaMA 3.1 8B Instruct² [8], Mistral Nemo 12B Instruct³ [10], and Phi 3.5 Mini 3.8B Instruct⁴ [9].

²https://huggingface.co/meta-llama/Meta-Llama-3.1-8B-Instruct

³https://huggingface.co/mistralai/Mistral-Nemo-Base-2407

⁴https://huggingface.co/microsoft/Phi-3.5-mini-instruct

Results. The experiments are conducted in the setting discussed in previous paragraphs. We evaluated the performance of three mainstream LLMs using Softmax attention with top-r indices. In particular, we chose r from the set $\{2^2, 2^4, 2^6, 2^8, 2^{10}, 2^{12}, 2^{15}\}$. As depicted in Figure 3, a significant increase in the perplexity (drop in performance) of LLMs is observed only when r falls below 2^4 . This suggests that the "massive activated" tokens are predominantly found within the top- 2^4 entries. In comparison to the total of 2^{15} entries, the "massive activated" entries constitute a relatively minor fraction. The observed results align with our theoretical analysis, confirming the approximation error of the Softmax attention with top-r indices is negligible for larger values of r.



Figure 3: We evaluated the perplexity of three mainstream language models: LLaMA 3.1 8B Instruct, Mistral Nemo 12B, and Phi 3.5 Mini 3.8B Instruct, using Softmax attention with top-r indices on the PaulGrahamEssays dataset. The results indicate a significant increase in perplexity only when the number of selected entries, r, falls below 2^4 . This observation aligns with our earlier findings that the proportion of "massive activated" entries is minimal compared to the total number of entries. Furthermore, the approximation error introduced by using top-r indices in Softmax attention remains negligible unless r becomes excessively small.

8. Discussion and Future Work

The sparsity within neural networks arises primarily from the incorporation of non-linear activation functions. These non-linear functions determine the mechanism or circuit of the neural networks [97]. Gaining insight into these non-linear layers not only enhances our understanding of how neural networks work but also paves the way for optimizing training and inference. We believe our analysis may inspire efficient neural network architecture design. This work represents the initial point of this envisioned blueprint. We concentrate on analyzing the combinations of LLMs and fundamental non-linear activation functions, ReLU and Softmax. By analyzing these functions, we aim to demonstrate to the research community that a thorough examination of a model's nonlinear characteristics can significantly enhance neural networks' running-time complexity.

In real-world scenarios, various non-linear activation functions exist beyond ReLU and Softmax, such as SELU(x) = scale \cdot (max(0, x) + min(0, $\alpha \cdot$ (exp(x) - 1))) [98], CELU(x) = max(0, x) + min(0, $\alpha \cdot$ (exp(x/α) - 1)) [99], and PRELU(x) = max(0, x) + weight \cdot min(0, x) [100]. Analyzing these alternative functions poses multiple challenges. We will explore these additional functions in the future.

9. Conclusion

This work investigates the exploitation of the intrinsic sparsity present in both ReLU and Softmax attention mechanisms to decrease the computational complexity of generation decoding and prompt prefilling scenarios. We employ the HSR data structure to accelerate the process of identifying nonzero or "massive activated" entries within ReLU and Softmax attentions. Our approach results in only a negligible approximation error for Softmax attention.

Acknowledgement

Research is partially supported by the National Science Foundation (NSF) Grants 2023239-DMS, CCF-2046710, and Air Force Grant FA9550-18-1-0166.

References

- [1] Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, et al. Lamda: Language models for dialog applications. arXiv preprint arXiv:2201.08239, 2022.
- [2] Andy Coenen, Luke Davis, Daphne Ippolito, Emily Reif, and Ann Yuan. Wordcraft: A human-ai collaborative editor for story writing. *arXiv preprint arXiv:*2107.07430, 2021.
- [3] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682*, 2022.
- [4] Tianyi Zhang, Faisal Ladhak, Esin Durmus, Percy Liang, Kathleen McKeown, and Tatsunori B Hashimoto. Benchmarking large language models for news summarization. *Transactions of the Association for Computational Linguistics*, 12:39–57, 2024.
- [5] Bowen Peng, Jeffrey Quesnelle, Honglu Fan, and Enrico Shippole. Yarn: Efficient context window extension of large language models. *arXiv preprint arXiv:2309.00071*, 2023.
- [6] OpenAI. Gpt-4 turbo, 2023. URL https://openai.com/blog/ new-models-and-developer-products-announced-at-devday.
- [7] Anthropic. Claude 3.5 sonnet, 2024. URL https://www.anthropic.com/news/ claude-3-5-sonnet.
- [8] Meta. Llama 3, 2024. URL https://ai.meta.com/blog/meta-llama-3/.
- [9] Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, et al. Phi-3 technical report: A highly capable language model locally on your phone. arXiv preprint arXiv:2404.14219, 2024.
- [10] MistralAI. Mistral nemo, 2024. URL https://mistral.ai/news/mistral-nemo/.
- [11] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information* processing systems, 30, 2017.
- [12] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [13] Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of naacL-HLT*, volume 1, page 2. Minneapolis, Minnesota, 2019.
- [14] Zirui Wang, Zhizhou Sha, Zheng Ding, Yilin Wang, and Zhuowen Tu. Tokencompose: Grounding diffusion with token-level supervision. *arXiv preprint arXiv:*2312.03626, 2023.
- [15] Yilin Wang, Zeyuan Chen, Liangjun Zhong, Zheng Ding, Zhizhou Sha, and Zhuowen Tu. Dolfin: Diffusion layout transformers without autoencoder. *arXiv preprint arXiv*:2310.16305, 2023.

- [16] Yilin Wang, Haiyang Xu, Xiang Zhang, Zeyuan Chen, Zhizhou Sha, Zirui Wang, and Zhuowen Tu. Omnicontrolnet: Dual-stage integration for conditional image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7436– 7448, 2024.
- [17] Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse transformers. *arXiv preprint arXiv:1904.10509*, 2019.
- [18] Sotiris Anagnostidis, Dario Pavllo, Luca Biggio, Lorenzo Noci, Aurelien Lucchi, and Thomas Hofmann. Dynamic context pruning for efficient and interpretable autoregressive transformers. Advances in Neural Information Processing Systems, 36, 2023.
- [19] Zichang Liu, Jue Wang, Tri Dao, Tianyi Zhou, Binhang Yuan, Zhao Song, Anshumali Shrivastava, Ce Zhang, Yuandong Tian, Christopher Re, et al. Deja vu: Contextual sparsity for efficient llms at inference time. In *International Conference on Machine Learning*, pages 22137– 22176. PMLR, 2023.
- [20] Jiaming Tang, Yilong Zhao, Kan Zhu, Guangxuan Xiao, Baris Kasikci, and Song Han. Quest: Query-aware sparsity for efficient long-context llm inference. arXiv preprint arXiv:2406.10774, 2024.
- [21] Mingjie Sun, Xinlei Chen, J Zico Kolter, and Zhuang Liu. Massive activations in large language models. arXiv preprint arXiv:2402.17762, 2024.
- [22] Zhenyu Zhang, Ying Sheng, Tianyi Zhou, Tianlong Chen, Lianmin Zheng, Ruisi Cai, Zhao Song, Yuandong Tian, Christopher Ré, Clark Barrett, et al. H20: Heavy-hitter oracle for efficient generative inference of large language models. *Advances in Neural Information Processing Systems*, 36, 2023.
- [23] Yuhong Li, Yingbing Huang, Bowen Yang, Bharat Venkitesh, Acyr Locatelli, Hanchen Ye, Tianle Cai, Patrick Lewis, and Deming Chen. Snapkv: Llm knows what you are looking for before generation. arXiv preprint arXiv:2404.14469, 2024.
- [24] Mitchell Wortsman, Jaehoon Lee, Justin Gilmer, and Simon Kornblith. Replacing softmax with relu in vision transformers. *arXiv preprint arXiv:2309.08586*, 2023.
- [25] Weizhe Hua, Zihang Dai, Hanxiao Liu, and Quoc Le. Transformer quality in linear time. In *International conference on machine learning*, pages 9099–9117. PMLR, 2022.
- [26] Pankaj K Agarwal, David Eppstein, and Jirí Matousek. Dynamic half-space reporting, geometric optimization, and minimum spanning trees. In *Annual Symposium on Foundations of Computer Science*, volume 33, pages 80–80. IEEE COMPUTER SOCIETY PRESS, 1992.
- [27] Elias Frantar and Dan Alistarh. Sparsegpt: Massive language models can be accurately pruned in one-shot. In *International Conference on Machine Learning*, pages 10323–10337. PMLR, 2023.
- [28] Josh Alman and Zhao Song. Fast attention requires bounded entries. *Advances in Neural Information Processing Systems*, 36, 2023.
- [29] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*, 2023.
- [30] Tri Dao and Albert Gu. Transformers are ssms: Generalized models and efficient algorithms through structured state space duality. *arXiv preprint arXiv:2405.21060, 2024*.
- [31] Praneeth Kacham, Vahab Mirrokni, and Peilin Zhong. Polysketchformer: Fast transformers via sketches for polynomial kernels. *arXiv preprint arXiv:2310.01655*, 2023.

- [32] Jerry Yao-Chieh Hu, Donglin Yang, Dennis Wu, Chenwei Xu, Bo-Yu Chen, and Han Liu. On sparse modern hopfield model. In *Thirty-seventh Conference on Neural Information Processing Systems (NeurIPS)*, 2023.
- [33] Dennis Wu, Jerry Yao-Chieh Hu, Weijian Li, Bo-Yu Chen, and Han Liu. STanhop: Sparse tandem hopfield model for memory-enhanced time series prediction. In *The Twelfth International Conference on Learning Representations (ICLR)*, 2024.
- [34] Jerry Yao-Chieh Hu, Thomas Lin, Zhao Song, and Han Liu. On computational limits of modern hopfield models: A fine-grained complexity analysis. In *Forty-first International Conference on Machine Learning (ICML)*, 2024.
- [35] Chenwei Xu, Yu-Chao Huang, Jerry Yao-Chieh Hu, Weijian Li, Ammar Gilani, Hsi-Sheng Goan, and Han Liu. Bishop: Bi-directional cellular learning for tabular data with generalized sparse modern hopfield model. In *Forty-first International Conference on Machine Learning* (*ICML*), 2024.
- [36] Dennis Wu, Jerry Yao-Chieh Hu, Teng-Yun Hsiao, and Han Liu. Uniform memory retrieval with larger capacity for modern hopfield models. In *Forty-first International Conference on Machine Learning (ICML)*, 2024.
- [37] Jerry Yao-Chieh Hu, Pei-Hsuan Chang, Haozheng Luo, Hong-Yu Chen, Weijian Li, Wei-Po Wang, and Han Liu. Outlier-efficient hopfield layers for large transformer-based models. In *Forty-first International Conference on Machine Learning (ICML)*, 2024.
- [38] Jerry Yao-Chieh Hu, Bo-Yu Chen, Dennis Wu, Feng Ruan, and Han Liu. Nonparametric modern hopfield models. *arXiv preprint arXiv:*2404.03900, 2024.
- [39] Michael Zhang, Kush Bhatia, Hermann Kumbong, and Christopher Ré. The hedgehog & the porcupine: Expressive linear attentions with softmax mimicry. *arXiv preprint arXiv*:2402.04347, 2024.
- [40] Jean Mercat, Igor Vasiljevic, Sedrick Keh, Kushal Arora, Achal Dave, Adrien Gaidon, and Thomas Kollar. Linearizing large language models. *arXiv preprint arXiv:*2405.06640, 2024.
- [41] Josh Alman and Zhao Song. The fine-grained complexity of gradient computation for training large language models. *arXiv preprint arXiv*:2402.04497, 2024.
- [42] Josh Alman and Zhao Song. How to capture higher-order correlations? generalizing matrix softmax attention to kronecker computation. In *The Twelfth International Conference on Learning Representations*, 2024.
- [43] Insu Han, Rajesh Jayaram, Amin Karbasi, Vahab Mirrokni, David Woodruff, and Amir Zandieh. Hyperattention: Long-context attention in near-linear time. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum? id=Eh00d2BJIM.
- [44] Amir Zandieh, Insu Han, Vahab Mirrokni, and Amin Karbasi. Subgen: Token generation in sublinear time and memory. *arXiv preprint arXiv:*2402.06082, 2024.
- [45] Yingyu Liang, Zhenmei Shi, Zhao Song, and Yufa Zhou. Tensor attention training: Provably efficient learning of higher-order transformers. *arXiv preprint arXiv:*2405.16411, 2024.
- [46] Michael Poli, Stefano Massaroli, Eric Nguyen, Daniel Y Fu, Tri Dao, Stephen Baccus, Yoshua Bengio, Stefano Ermon, and Christopher Ré. Hyena hierarchy: Towards larger convolutional language models. In *International Conference on Machine Learning*, pages 28043–28078. PMLR, 2023.
- [47] Ruisi Cai, Yuandong Tian, Zhangyang Wang, and Beidi Chen. Lococo: Dropping in convolutions for long context compression. *arXiv preprint arXiv:2406.05317*, 2024.

- [48] Yingyu Liang, Zhenmei Shi, Zhao Song, and Chiwun Yang. Toward infinite-long prefix in transformer. *arXiv preprint arXiv:2406.14036*, 2024.
- [49] Yingyu Liang, Heshan Liu, Zhenmei Shi, Zhao Song, and Junze Yin. Conv-basis: A new paradigm for efficient attention inference and gradient computation in transformers. *arXiv preprint arXiv*:2405.05219, 2024.
- [50] Yeqi Gao, Zhao Song, Weixin Wang, and Junze Yin. 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*, 2023.
- [51] Harry Dong, Xinyu Yang, Zhenyu Zhang, Zhangyang Wang, Yuejie Chi, and Beidi Chen. Get more with less: Synthesizing recurrence with kv cache compression for efficient llm inference. arXiv preprint arXiv:2402.09398, 2024.
- [52] Yingyu Liang, Zhizhou Sha, Zhenmei Shi, Zhao Song, and Yufa Zhou. Multi-layer transformers gradient can be approximated in almost linear time. *arXiv preprint arXiv:2408.13233*, 2024.
- [53] Yekun Ke, Xiaoyu Li, Yingyu Liang, Zhizhou Sha, Zhenmei Shi, and Zhao Song. On computational limits and provably efficient criteria of visual autoregressive models: A fine-grained complexity analysis. arXiv preprint arXiv:2501.04377, 2025.
- [54] Yifang Chen, Jiayan Huo, Xiaoyu Li, Yingyu Liang, Zhenmei Shi, and Zhao Song. Fast gradient computation for rope attention in almost linear time. *arXiv preprint arXiv*:2412.17316, 2024.
- [55] Yingyu Liang, Zhizhou Sha, Zhenmei Shi, Zhao Song, and Yufa Zhou. Looped relu mlps may be all you need as practical programmable computers, 2024.
- [56] Xiaoyu Li, Yingyu Liang, Zhenmei Shi, Zhao Song, and Yufa Zhou. Fine-grained attention i/o complexity: Comprehensive analysis for backward passes, 2024.
- [57] Yingyu Liang, Jiangxuan Long, Zhenmei Shi, Zhao Song, and Yufa Zhou. Beyond linear approximations: A novel pruning approach for attention matrix, 2024.
- [58] Yingyu Liang, Zhenmei Shi, Zhao Song, and Yufa Zhou. Differential privacy of crossattention with provable guarantee. *arXiv preprint arXiv:*2407.14717, 2024.
- [59] Yingyu Liang, Zhenmei Shi, Zhao Song, and Yufa Zhou. Unraveling the smoothness properties of diffusion models: A gaussian mixture perspective. *arXiv preprint arXiv*:2405.16418, 2024.
- [60] Zhenmei Shi, Junyi Wei, Zhuoyan Xu, and Yingyu Liang. Why larger language models do in-context learning differently? *arXiv preprint arXiv:2405.19592*, 2024.
- [61] Zhuoyan Xu, Zhenmei Shi, and Yingyu Liang. Do large language models have compositional ability? an investigation into limitations and scalability. In ICLR 2024 Workshop on Mathematical and Empirical Understanding of Foundation Models, 2024.
- [62] Zhenmei Shi, Yifei Ming, Xuan-Phi Nguyen, Yingyu Liang, and Shafiq Joty. Discovering the gems in early layers: Accelerating long-context llms with 1000x input token reduction. *arXiv preprint arXiv*:2409.17422, 2024.
- [63] Xiaoyu Li, Yingyu Liang, Zhenmei Shi, and Zhao Song. A tighter complexity analysis of sparsegpt. arXiv preprint arXiv:2408.12151, 2024.
- [64] Zhao Song, Mingquan Ye, and Lichen Zhang. Streaming semidefinite programs: $O(\sqrt{n})$ passes, small space and fast runtime. *arXiv preprint arXiv:2309.05135*, 2023.

- [65] Lianke Qin, Zhao Song, Lichen Zhang, and Danyang Zhuo. An online and unified algorithm for projection matrix vector multiplication with application to empirical risk minimization. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 101–156. PMLR, 2023.
- [66] Zhao Song, Xin Yang, Yuanyuan Yang, and Lichen Zhang. Sketching meets differential privacy: fast algorithm for dynamic kronecker projection maintenance. In *International Conference on Machine Learning (ICML)*, pages 32418–32462. PMLR, 2023.
- [67] Lichen Zhang. *Speeding up optimizations via data structures: Faster search, sample and maintenance.* PhD thesis, Master's thesis, Carnegie Mellon University, 2022.
- [68] Zhao Song, Lichen Zhang, and Ruizhe Zhang. Training multi-layer over-parametrized neural network in subquadratic time. In *Innovations in Theoretical Computer Science (ITCS)*, pages 93:1–93:15, 2024.
- [69] Yang Cao, Bo Chen, Xiaoyu Li, Yingyu Liang, Zhizhou Sha, Zhenmei Shi, Zhao Song, and Mingda Wan. Force matching with relativistic constraints: A physics-inspired approach to stable and efficient generative modeling. *arXiv preprint arXiv:2502.08150*, 2025.
- [70] Bo Chen, Chengyue Gong, Xiaoyu Li, Yingyu Liang, Zhizhou Sha, Zhenmei Shi, Zhao Song, and Mingda Wan. High-order matching for one-step shortcut diffusion models. *arXiv preprint arXiv*:2502.00688, 2025.
- [71] Bo Chen, Xiaoyu Li, Yingyu Liang, Zhenmei Shi, and Zhao Song. Bypassing the exponential dependency: Looped transformers efficiently learn in-context by multi-step gradient descent. In *International Conference on Artificial Intelligence and Statistics*, 2025.
- [72] Bo Chen, Xiaoyu Li, Yingyu Liang, Jiangxuan Long, Zhenmei Shi, and Zhao Song. Circuit complexity bounds for rope-based transformer architecture. *arXiv preprint arXiv*:2411.07602, 2024.
- [73] Xuan Shen, Zhao Song, Yufa Zhou, Bo Chen, Yanyu Li, Yifan Gong, Kai Zhang, Hao Tan, Jason Kuen, Henghui Ding, et al. Lazydit: Lazy learning for the acceleration of diffusion transformers. *arXiv preprint arXiv:2412.12444*, 2024.
- [74] Xuan Shen, Zhao Song, Yufa Zhou, Bo Chen, Jing Liu, Ruiyi Zhang, Ryan A Rossi, Hao Tan, Tong Yu, Xiang Chen, et al. Numerical pruning for efficient autoregressive models. arXiv preprint arXiv:2412.12441, 2024.
- [75] Bo Chen, Xiaoyu Li, Yingyu Liang, Zhao Song, and Zhizhou Sha. Nrflow: Towards noiserobust generative modeling via second-order flow matching. *manuscript*, 2025.
- [76] Xiaoyu Li, Yingyu Liang, Jiangxuan Long, Zhenmei Shi, Zhao Song, and Zhen Zhuang. Neural algorithmic reasoning for hypergraphs with looped transformers. *arXiv preprint arXiv:2501.10688*, 2025.
- [77] Yekun Ke, Yingyu Liang, Zhenmei Shi, Zhao Song, and Chiwun Yang. Curse of attention: A kernel-based perspective for why transformers fail to generalize on time series forecasting and beyond. In *Conference on Parsimony and Learning*. PMLR, 2025.
- [78] Chenyang Li, Yingyu Liang, Zhenmei Shi, Zhao Song, and Tianyi Zhou. Fourier circuits in neural networks and transformers: A case study of modular arithmetic with multiple inputs. In *International Conference on Artificial Intelligence and Statistics*, 2025.
- [79] Jiayu Wang, Yifei Ming, Zhenmei Shi, Vibhav Vineet, Xin Wang, Yixuan Li, and Neel Joshi. Is a picture worth a thousand words? delving into spatial reasoning for vision language models. *Advances in Neural Information Processing Systems*, 36, 2024.

- [80] Jerry Yao-Chieh Hu, Weimin Wu, Zhuoru Li, Sophia Pi, , Zhao Song, and Han Liu. On statistical rates and provably efficient criteria of latent diffusion transformers (dits). *Advances in Neural Information Processing Systems*, 38, 2024.
- [81] Jerry Yao-Chieh Hu, Weimin Wu, Yi-Chen Lee, Yu-Chao Huang, Minshuo Chen, and Han Liu. On statistical rates of conditional diffusion transformers: Approximation, estimation and minimax optimality. arXiv preprint arXiv:2411.17522, 2024.
- [82] Jerry Yao-Chieh Hu, Wei-Po Wang, Ammar Gilani, Chenyang Li, Zhao Song, and Han Liu. Fundamental limits of prompt tuning transformers: Universality, capacity and efficiency. *arXiv preprint arXiv*:2411.16525, 2024.
- [83] Weimin Wu, Maojiang Su, Jerry Yao-Chieh Hu, Zhao Song, and Han Liu. Transformers are deep optimizers: Provable in-context learning for deep model training. *arXiv preprint arXiv*:2411.16549, 2024.
- [84] Kai Shen, Junliang Guo, Xu Tan, Siliang Tang, Rui Wang, and Jiang Bian. A study on relu and softmax in transformer. *arXiv preprint arXiv:2302.06461*, 2023.
- [85] Zhiyuan Li, Srinadh Bhojanapalli, Manzil Zaheer, Sashank Reddi, and Sanjiv Kumar. Robust training of neural networks using scale invariant architectures. In *International Conference on Machine Learning*, pages 12656–12684. PMLR, 2022.
- [86] Yu Bai, Fan Chen, Huan Wang, Caiming Xiong, and Song Mei. Transformers as statisticians: Provable in-context learning with in-context algorithm selection. *Advances in neural information processing systems*, 36, 2023.
- [87] Hengyu Fu, Tianyu Guo, Yu Bai, and Song Mei. What can a single attention layer learn? a study through the random features lens. *Advances in Neural Information Processing Systems*, 36, 2023.
- [88] Yingyu Liang, Zhizhou Sha, Zhenmei Shi, Zhao Song, and Yufa Zhou. Looped relu mlps may be all you need as practical programmable computers. In *International Conference on Artificial Intelligence and Statistics*, 2025.
- [89] Shunhua Jiang, Zhao Song, Omri Weinstein, and Hengjie Zhang. A faster algorithm for solving general lps. In *Proceedings of the 53rd Annual ACM SIGACT Symposium on Theory of Computing*, pages 823–832, 2021.
- [90] Sayan Bhattacharya, Peter Kiss, and Thatchaphol Saranurak. Dynamic algorithms for packing-covering lps via multiplicative weight updates. In *Proceedings of the 2023 Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 1–47. SIAM, 2023.
- [91] Lianke Qin, Zhao Song, and Yuanyuan Yang. Efficient sgd neural network training via sublinear activated neuron identification. *arXiv preprint arXiv:2307.06565*, 2023.
- [92] Yeqi Gao, Lianke Qin, Zhao Song, and Yitan Wang. A sublinear adversarial training algorithm. *arXiv preprint arXiv:2208.05395*, 2022.
- [93] Danny Z Chen, Michiel Smid, and Bin Xu. Geometric algorithms for density-based data clustering. International Journal of Computational Geometry & Applications, 15(03):239–260, 2005.
- [94] Wenqi Ju, Chenglin Fan, Jun Luo, Binhai Zhu, and Ovidiu Daescu. On some geometric problems of color-spanning sets. *Journal of Combinatorial Optimization*, 26:266–283, 2013.
- [95] David Eppstein, Michael T Goodrich, Doruk Korkmaz, and Nil Mamano. Defining equitable geographic districts in road networks via stable matching. In *Proceedings of the 25th ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems*, pages 1–4, 2017.

- [96] Greg Kamradt. Llmtest-needleinahaystack, 2024. URL https://github.com/gkamradt/ LLMTest_NeedleInAHaystack.
- [97] Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, et al. In-context learning and induction heads. *arXiv preprint arXiv:2209.11895*, 2022.
- [98] Günter Klambauer, Thomas Unterthiner, Andreas Mayr, and Sepp Hochreiter. Selfnormalizing neural networks. *Advances in neural information processing systems*, 30, 2017.
- [99] Jonathan T Barron. Continuously differentiable exponential linear units. *arXiv preprint arXiv*:1704.07483, 2017.
- [100] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034, 2015.
- [101] Beatrice Laurent and Pascal Massart. Adaptive estimation of a quadratic functional by model selection. *Annals of Statistics*, pages 1302–1338, 2000.
- [102] Sergei Bernstein. 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.
- [103] Zhao Song, Shuo Yang, and Ruizhe Zhang. Does preprocessing help training overparameterized neural networks? *Advances in Neural Information Processing Systems*, 34:22890– 22904, 2021.

Appendix

Roadmap. In Section A, we discuss some limitations of our work. In Section B, we introduce more fundamental lemmas and facts. In Section C, we extend the analysis to ReLU attention calculation, demonstrating improved performance over standard attention computation under specific conditions. In Section D, we first introduce and analyze the time complexity of ReLU attention generation using half-space reporting (HSR) data structures. In Section E, we analyze the sparsity of ReLU attention matrices. In Section F, we introduce our results on reducing the running time of Softmax attention. In Section G, we analyze error bounds for Softmax attention with index sets, balancing efficiency and accuracy.

| Sequence length | Activated entries | Sparsity ratio |
|--------------------|----------------------|-------------------|
| 1k | 251 | 0.75 |
| $2\mathbf{k}$ | 437 | 0.78 |
| $4\mathbf{k}$ | 761 | 0.81 |
| 8k | 1325 | 0.83 |
| 16k | 2308 | 0.86 |
| 32k | 4019 | 0.87 |
| 64k | 6997 | 0.89 |
| 128k | 12183 | 0.90 |
| 256k | 21212 | 0.92 |
| 512k | 36933 | 0.93 |
| 1024k | 64304 | 0.94 |

Table 1: An illustration of the sparsity level attained by our algorithm across varying sequence lengths, n. Our approach activates merely $n^{4/5}$ entries per inference, resulting in a computational savings of up to 90% when n = 1024k.

A. Limitations

Our work is fundamentally theoretical in nature. We concentrate on the theoretical analysis of our proposed algorithms (Algorithms 1 and 2). Our study does not include an implementation of the suggested algorithms, a limitation that arises from the absence of an existing implementation for the original HSR data structure, initially proposed in [26]. We are confident that our theoretical insights will inspire the development of future algorithmic designs.

B. Full Background and Definition

In this section, we display more fundamental concepts. In Section B.1, we introduce a modified version of Softmax attention that operates on a specific subset of indices. It defines the top-r nearest neighbors Softmax attention, which focuses on the most relevant entries in the attention matrix. In Section B.2, we describe the massive activation property for attention mechanisms. In Section B.3, we introduce several important probability properties and bounds. In Section B.4, we detail the time complexity and performance of half-space reporting (HSR) data structures.

B.1. Softmax Attention with Index Set

Recall that we have already provided the definition of ReLU attention in Definition 1.2. Here, we present the key concepts of Softmax attention. For Softmax attention, since we only calculate the "massive activated" entries to get our approximated results, we introduce the formal definition:

Definition B.1 (Input with index set). Let $K \in \mathbb{R}^{n \times d}$ and $V \in \mathbb{R}^{n \times d}$ be defined in Definition 1.1. Let $R \subseteq [n]$ be an index set of size $|R| = r \in [n]$. Let $\overline{R} := [n] \setminus R$ be the complementary set, where $|\overline{R}| = n - r$.

We define

$$\widehat{K} := K_R \in \mathbb{R}^{r \times d} \quad \widehat{V} := V_R \in \mathbb{R}^{r \times d} \quad \overline{K} := K_{\overline{R}} \in \mathbb{R}^{(n-r) \times d} \quad \overline{V} := V_{\overline{R}} \in \mathbb{R}^{(n-r) \times d}$$

as the submatrix of K and V, i.e., whose row index is in R or \overline{R} , respectively.

We consider calculating the Softmax attention on the "massive activation" index set, where we define the "massive activation" index set as the top-*r* indices, where the formal definition is as follows: **Definition B.2** (Top-*r* indices Softmax attention). Let $q \in \mathbb{R}^d$, $K, V \in \mathbb{R}^{n \times d}$ be defined in Definition 1.1. Let $NN(r, q, K) \subseteq [n]$ denote the indices of top-*r* entries of qK, where |NN(r, q, K)| = r. Let $\widehat{K}, \widehat{V} \in \mathbb{R}^{r \times d}$ and $\overline{K}, \overline{V} \in \mathbb{R}^{(n-r) \times d}$ be defined in Definition B.1. We define the top-*r* nearest neighbors (NN) Softmax attention computation $\widehat{Attn}_s(q, K, V) \in \mathbb{R}^d$ as follows:

$$\widehat{\mathsf{Attn}}_s(q,K,V) := \mathsf{Softmax}(q\widehat{K}^\top)\widehat{V} = \widehat{\alpha}^{-1}\widehat{u}\widehat{V} \in \mathbb{R}^d$$

where

$$\widehat{u} := \exp(q\widehat{K}^{\top}) \in \mathbb{R}^r \quad and \quad \widehat{\alpha} := \langle \widehat{u}, \mathbf{1}_r \rangle \in \mathbb{R}.$$

Furthermore, we define $\overline{u} := \exp(q\overline{K}^{\top}) \in \mathbb{R}^{n-r}$, $\overline{\alpha} := \langle \overline{u}, \mathbf{1}_{n-r} \rangle \in \mathbb{R}$, and $u := \exp(qK^{\top}) \in \mathbb{R}^{n+1}$, $\alpha := \langle u, \mathbf{1}_{n+1} \rangle \in \mathbb{R}$.

In Definition B.2, we view the "massive activated" entries as the top-*r* entries. Therefore, we only calculate the Softmax attention based on $\hat{K}, \hat{V} \in \mathbb{R}^{r \times d}$, instead of $K, V \in \mathbb{R}^{n \times d}$.

B.2. Massive Activation

Now, we introduce our observations on the properties of the attention scores (the inner products of query vectors and key vectors). This further facilitates the error analysis of the top-r indices Softmax attention. To begin with, we provide the definition of the massive activation property as follows:

Definition B.3 (Massive activation property). Let $\gamma \in [0,1]$, $\beta_1 \ge \beta_2 \ge 0$. Let $\mathsf{NN}(r,q,K) \subseteq [n]$ denote the indices of top-*r* entries of qK^{\top} . We define $(\gamma, \beta_1, \beta_2)$ massive activation for a query $q \in \mathbb{R}^d$ and key cache $K \in \mathbb{R}^{n \times d}$, if the following conditions hold:

- The top- n^{γ} entries are massive, i.e., $\frac{1}{n^{\gamma} \cdot ||q||_2} \sum_{i \in NN(n^{\gamma}, q, K)} \langle q, K_i \rangle \geq \beta_1 \log(n)$.
- The remaining terms are upper bounded, i.e, $\forall i \in [n] \setminus \mathsf{NN}(n^{\gamma}, q, K), \frac{1}{\|q\|_2} \langle q, K_i \rangle \leq \beta_2 \log(n).$

An intuitive understanding of Definition B.3 is that the summation of "massive activated" entries dominates the summation of all entries, and the entries we ignored only contribute little to the final summation. Therefore, it is reasonable for us to omit those non "massive activated" entries.

Remark B.4. There are many distributions satisfying the property in Definition B.3, such as (1) K drawing from any subexponential distribution, e.g., multivariate Laplace distributions, (2) K drawing from any mixture of Gaussian distribution with $n^{1-\gamma}$ Gaussian clusters.

B.3. Probability Tools

We state several fundamental properties and bounds for some common distributions.

Fact B.5 (Weighted summation of Gaussian). If the following conditions hold:

- Let $x \in \mathbb{R}^d$ be a fixed vector and $y \in \mathbb{R}^d$ be a random vector.
- For $i \in [d]$, let x_i denote the *i*-th entry of x.
- Suppose for $i \in [d]$, $y_i \sim \mathcal{N}(0, \sigma^2)$.

Then the inner product of x and y, $\langle x, y \rangle$ conforms Gaussian distribution $\mathcal{N}(0, \|x\|_2^2 \sigma^2)$. Namely, we have $\langle x, y \rangle \sim \mathcal{N}(0, \|x\|_2^2 \sigma^2)$.

Fact B.6 (Independence between $\langle x, y_i \rangle$ and $\langle x, y_j \rangle$). *If the following conditions hold:*

- Let $x \in \mathbb{R}^d$ be a fixed vector.
- Let $y_1, y_2, \dots, y_n \in \mathbb{R}^d$ be *n* random vectors.
- For any $i, j \in [n], i \neq j$, y_i and y_j are independent.

Then, for any $i, j \in [n], i \neq j, \langle x, y_i \rangle$ *and* $\langle x, y_i \rangle$ *are independent.*

We provide tail bounds for chi-square and Gaussian distributed random variables:

Lemma B.7 (Chi-square tail bound, Lemma 1 in [101]). Let $X \sim \mathcal{X}_k^2$ be a chi-squared distributed random variable with k degrees of freedom. Each one has zero means and σ^2 variance.

Then, it holds that

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

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

Fact B.8 (Gaussian tail bound). Suppose we have a random variable $x \sim \mathcal{N}(\mu, \sigma)$.

Then, for $t \in \mathbb{R}$ *, we have*

$$\Pr[x \ge \mu + t] \le \exp(-\frac{t^2}{2\sigma^2})$$

Proof. We can show

$$\Pr[x \ge \mu + t] = \Pr[x - \mu \ge t]$$

$$= \Pr[e^{x - \mu} \ge e^{t}]$$

$$= \inf_{\lambda \ge 0} \Pr[e^{\lambda(x - \mu)} \ge e^{\lambda t}]$$

$$\leq \inf_{\lambda \ge 0} \frac{\mathbb{E}[e^{\lambda(x - \mu)}]}{e^{\lambda t}}$$
(1)

where the first step, the second step follows from basic algebra, the third step follows from that the inequality holds for any $\lambda > 0$, and the fourth step follows from Markov's inequality.

Then we consider the numerator, and we use $y = x - \mu$ to simplify the calculation, we have

$$\mathbb{E}[e^{\lambda y}] = \int_{\mathbb{R}} e^{\lambda y} \frac{e^{-y^2/2\sigma^2}}{\sqrt{2\pi\sigma}} dy$$

$$= \int_{\mathbb{R}} \frac{e^{-(y-\lambda/\sigma^2)^2 \cdot \frac{1}{2\sigma^2}} e^{\lambda^2 \sigma^2/2}}{\sqrt{2\pi\sigma}} dy$$

$$= e^{\frac{\lambda^2 \sigma^2}{2}} \int_{\mathbb{R}} \frac{e^{-(y-\lambda/\sigma^2) \cdot \frac{1}{2\sigma^2}}}{\sqrt{2\pi\sigma}} dy$$

$$= e^{\frac{\lambda^2 \sigma^2}{2}}$$
(2)

where the first step follows from the definition of the moment generating function, the second and the third steps follow from basic algebra, and the fourth step follows from the property of the probability density function.

Then we have

$$\Pr[x \ge \mu + t] \le \inf_{\lambda \ge 0} \exp(\frac{\lambda^2 - \sigma^2}{2} - \lambda t)$$
$$\le \exp(-\frac{t^2}{2\sigma^2})$$

where the first step follows from Eq. (1) and Eq.(2), the second step follows from the calculation of infimum. \Box

The Bernstein's inequality for bounding sums of independent random variables is:

Lemma B.9 (Bernstein inequality [102]). Assume Z_1, \dots, Z_n are *n* i.i.d. random variables. $\forall i \in [n]$, $\mathbb{E}[Z_i] = 0$ and $|Z_i| \leq M$ almost surely. Let $Z = \sum_{i=1}^n Z_i$. Then,

$$\Pr\left[Z > t\right] \le \exp\left(-\frac{t^2/2}{\sum_{j=1}^n \mathbb{E}[Z_j^2] + Mt/3}\right), \forall t > 0.$$

B.4. Half-Space Reporting (HSR) Data Structures

| Algorithm 3 Half Space Report Data Structure | | | | | |
|--|--|--|--|--|--|
| 1: data structure HalfSpaceReport | | | | | |
| 2: $INIT(S, n, d)$ | \triangleright Initialize the data structure with a set <i>S</i> of <i>n</i> points in \mathbb{R}^d | | | | |
| 3: $QUERY(a, b)$ | $\triangleright a, b \in \mathbb{R}^d$. Output the set $\{x \in S : \operatorname{sgn}(\langle a, x \rangle - b) \ge 0\}$ | | | | |
| 4: end data structure | | | | | |

We restate the result from [26] for solving the half-space range reporting problem. The half-space range reporting problem is a fundamental problem in computational geometry and can be formally defined as follows:

Definition B.10 (Half-space range reporting [26, 103]). *Given a set* S *of* n *points in* \mathbb{R}^d *with initialization, we have an operation* Query(H): *given a half-space* $H \subset \mathbb{R}^d$, *output all of the points in* S *that contain in* H, *i.e.*, $S \cap H$.

The time complexity of the HSR data structure is:

Theorem B.11 (Agarwal, Eppstein and Matousek [26]). Let *d* be a fixed constant. Let *t* be a parameter between *n* and $n^{\lfloor d/2 \rfloor}$. There is a dynamic data structure for half-space reporting that uses $O_{d,\epsilon}(t^{1+\epsilon})$ space and pre-processing time, $O_{d,\epsilon}(\frac{n}{t^{1/\lfloor d/2 \rfloor}} \log n + k)$ time per query where *k* is the output size and $\epsilon > 0$ is any fixed constant, and $O_{d,\epsilon}(t^{1+\epsilon}/n)$ amortized update time.

As a direct corollary, we have

Corollary B.12 (HSR data-structure time complexity [26], formal version of Corollary 3.1). Let \mathcal{T}_{init} denote the pre-processing time to build the data structure, \mathcal{T}_{query} denote the time per query, and \mathcal{T}_{update} time per update. Given a set of *n* points in \mathbb{R}^d , the half-space range reporting problem can be solved with the following performances:

- Part 1. $\mathcal{T}_{init}(n,d) = O_d(n\log n), \mathcal{T}_{query}(n,d,k) = O(dn^{1-1/\lfloor d/2 \rfloor} + dk).$
- Part 2. $\mathcal{T}_{\text{init}}(n,d) = O(n^{\lfloor d/2 \rfloor}), \mathcal{T}_{\text{query}}(n,d,k) = O(d\log(n) + dk).$

C. ReLU Attention Prompt Prefilling

In this section, we focus on optimizing the standard ReLU attention calculation. By leveraging a HSR data structure and assuming sparsity, the time complexity can be reduced to $O(n^{1+4/5}d)$.

Lemma C.1 (General full attention computation framework, formal version of Lemma 6.3). *If the following conditions hold:*

- Let $Q \in \mathbb{R}^{m \times d}$ and $K, V \in \mathbb{R}^{n \times d}$ be defined as Definition 1.2.
- Assume each entry of K is from Gaussian $\mathcal{N}(0, \sigma_k^2)$, and each entry of Q is from Gaussian $\mathcal{N}(0, \sigma_a^2)$.
- Let $\delta \in (0, 1)$ denote the failure probability.
- Let $\sigma_a = 4 \cdot (1 + d^{-1} \log(m/\delta))^{1/2} \cdot \sigma_q \sigma_k$.
- Let $b = \sigma_a \cdot \sqrt{0.4 \log n}$.

• Let HSR data structure be defined as Part 1 in Corollary B.12.

There exists an algorithm (Algorithm 2), with at least $1 - \delta$ probability, computes full attention of Q, K, V in $O(mn^{1-1/\lfloor d/2 \rfloor} + mn^{4/5})$ time.

Proof. For $i \in [m]$, let $\tilde{k}_i := |\tilde{S}_{i,\text{fire}}|$ denote the number of non-zero entries in *i*-th row of $A \in \mathbb{R}^{m \times n}$.

The running time for INFERENCE procedure can be written as

$$\mathcal{T}_{\mathsf{init}}(n,d) + \sum_{i=1}^{m} \mathcal{T}_{\mathsf{query}}(n,d,\widetilde{k}_i) + O(d\sum_{i=1}^{m}\widetilde{k}_i) + O(d\sum_{i=1}^{m}\widetilde{k}_i)$$

The first term $\mathcal{T}_{init}(n, d)$ corresponds to the initialization of the HSR data structure. Since we use the Part 1 result from Corollary B.12, the running time for initialization is $\mathcal{T}_{init}(n, d) = O_d(n \log n)$.

The second term $\sum_{i=1}^{m} \mathcal{T}_{query}(n, d, \tilde{k}_i)$ comes from the HSR query operation (Line 11). Since we use Part 1 result from Corollary B.12, we have

$$\sum_{i=1}^{m} \mathcal{T}_{query}(n, d, \widetilde{k}_i) = O(mn^{1-1/\lfloor d/2 \rfloor}d + d\sum_{i=1}^{m} \widetilde{k}_i)$$
$$= O(mn^{1-1/\lfloor d/2 \rfloor}d + mn^{4/5}d)$$

where the first step follows from $\mathcal{T}_{query}(n, d, \tilde{k}_i) = O(dn^{1-\lfloor d/2 \rfloor} + d\tilde{k}_i)$ (Part 1 of Corollary B.12), the second step follows from with high probability \tilde{k}_i at most $n^{4/5}$ (Lemma E.3).

The third term $O(\sum_{i=1}^{m} \tilde{k}_i)$ corresponds to calculating $A_{j,i}$ (Line 13). By Lemma E.3, we have the third term is $O(mn^{4/5})$.

The fourth term $O(\sum_{i=1}^{m} \tilde{k}_i)$ corresponds to calculating $D^{-1}AV$. Since for *i*-th row of A, there are \tilde{k}_i non-zero entries. Therefore, it takes $O(\sum_{i=1}^{m} \tilde{k}_i)$ time for calculating $D^{-1}A$. Therefore, it takes $O(d \sum_{i=1}^{m} \tilde{k}_i)$ time to calculate $D^{-1}AV$. By Lemma E.3, with high probability, \tilde{k}_i is at most $n^{4/5}$. Therefore, we have the third term as $O(mn^{4/5}d)$.

To sum up, the overall running time is $O(mn^{1-1/\lfloor d/2 \rfloor}d + mn^{4/5}d)$.

We can now derive a more specific result for the full ReLU attention computation:

Theorem C.2 (Running time of full ReLU attention computation, formal version of Lemma 5.1). *If the following conditions hold:*

- Let ReLU attention be defined as Definition 1.2.
- Assume each entry of K is from Gaussian $\mathcal{N}(0, \sigma_k^2)$, and each entry of Q is from Gaussian $\mathcal{N}(0, \sigma_q^2)$.

- Let $\delta \in (0, 1)$ denote the failure probability.
- Let $\sigma_a = 4 \cdot (1 + d^{-1} \log(m/\delta))^{1/2} \cdot \sigma_q \sigma_k$.
- Let $b = \sigma_a \cdot \sqrt{0.4 \log n}$.
- Suppose we have $Q, K, V \in \mathbb{R}^{n \times d}$.

There exists an algorithm (Algorithm 2), with probability at least $1 - \delta$, takes $O(n^{2-1/\lfloor d/2 \rfloor}d + n^{1+4/5}d)$ time to compute the full ReLU attention of Q, K, V.

Proof. By Lemma C.1, we have that the FullAttentionComputation data structure (Algorithm 2) can run Inference to calculate the ReLU attention, in $O(m^{1-\lfloor d/2 \rfloor}nd + mn^{4/5}d)$ time.

By our assumption, we have $Q \in \mathbb{R}^{n \times d}$. For each calculation, we only need to call FullAttention-Computation.Inference(K, Q, V, n, n, d) for once.

Then, we have the ReLU attention calculation run in $O(n^{1+4/5}d)$ time.

D. ReLU Attention Generation Decoding

In this section, we present a theoretical analysis of the time complexity of ReLU attention generation using an HSR data structure.

Lemma D.1 (General attention generation framework, formal version of Lemma 6.2). *If the following conditions hold:*

- Let $Q \in \mathbb{R}^{m \times d}$ and $K, V \in \mathbb{R}^{n \times d}$ be defined as Definition 1.2.
- Assume each entry of K is from Gaussian $\mathcal{N}(0, \sigma_k^2)$, and each entry of Q is from Gaussian $\mathcal{N}(0, \sigma_a^2)$.
- Let $\delta \in (0, 1)$ denote the failure probability.
- Let $\sigma_a = 4 \cdot (1 + d^{-1} \log(m/\delta))^{1/2} \cdot \sigma_q \sigma_k$.
- Let $b = \sigma_a \cdot \sqrt{0.4 \log n}$.
- Let HSR data structure be defined as Part 2 in Corollary B.12.

Then, there exists an algorithm (Algorithm 1), with at least $1 - \delta$ *probability, has the following performance:*

- **Part 1.** The INIT procedure runs in $O(n^{\lfloor d/2 \rfloor})$ time.
- **Part 2.** For each query, the INFERENCE procedure runs in $O(mn^{4/5}d)$ time.

Proof. Proof of Part 1.

The INIT procedure only runs the initialization of the HSR data structure. Since we use Part 2 result from Corollary B.12, the running time of INIT procedure is $\mathcal{T}_{init}(n,d) = O(n^{\lfloor d/2 \rfloor})$.

Proof of Part 2.

For $i \in [m]$, let $\widetilde{k}_i := |\widetilde{S}_{i,\text{fire}}|$ denote the number of non-zero entries in *i*-th row of $A \in \mathbb{R}^{m \times n}$.

The running time for INFERENCE procedure can be written as

$$\sum_{i=1}^{m} \mathcal{T}_{\mathsf{query}}(n, d, \widetilde{k}_i) + O(d\sum_{i=1}^{m} \widetilde{k}_i) + O(d\sum_{i=1}^{m} \widetilde{k}_i)$$

The first term $\sum_{i=1}^{m} \mathcal{T}_{query}(n, d, \tilde{k}_i)$ corresponds to the HSR query operation (Line 16). Since we use the Part 2 result from Corollary B.12, we have

$$\sum_{i=1}^{m} \mathcal{T}_{query}(n, d, \tilde{k}_i) = O(md \log n + d \sum_{i=1}^{m} \tilde{k}_i)$$
$$= O(md \log n + mn^{4/5}d)$$
$$= O(mn^{4/5}d)$$

where the first step follows from $\mathcal{T}_{query}(n, d, k) = O(d \log n + dk)$ in Part 2 of Corollary B.12, the second step follows from with high probability, \tilde{k}_i is at most $n^{4/5}$ (Lemma E.3), the third step follows from $\log n < n^{4/5}$.

The second term $O(d \sum_{i=1}^{m} \tilde{k}_i)$ corresponds to calculating $A_{i,j}$ (Line 18). There are *m* iterations, and in each iteration, it calculates \tilde{k}_i entries of *A*. Then, the second term is $O(d \sum_{i=1}^{m} \tilde{k}_i)$. By Lemma E.3, with high probability, \tilde{k}_i is at most $n^{4/5}$. Therefore, we have the second term as $O(mn^{4/5}d)$.

Similar to the proof of Lemma C.1 this term is $O(mn^{4/5}d)$.

To sum up, we have the overall running time for INFERENCE procedure is $O(mn^{4/5}d)$.

We now derive a comprehensive sparsity analysis for the ReLU attention mechanism:

Theorem D.2 (Running time of full ReLU attention generation, formal version of Theorem 4.1). *If the following conditions hold:*

- Let ReLU attention be defined as Definition 1.2.
- Assume each entry of K is from Gaussian $\mathcal{N}(0, \sigma_k^2)$, and each entry of Q is from Gaussian $\mathcal{N}(0, \sigma_a^2)$.
- Let $\delta \in (0, 1)$ denote the failure probability.
- Let $\sigma_a = 4 \cdot (1 + d^{-1} \log(m/\delta))^{1/2} \cdot \sigma_q \sigma_k$.
- Let $b = \sigma_a \cdot \sqrt{0.4 \log n}$.
- Suppose we have KV Cache $K, V \in \mathbb{R}^{n \times d}$. We want to generate a m length answer, where $n \gg m$.

There exists an algorithm (Algorithm 1), with at least $1 - \delta$ *probability, takes O*($mn^{4/5}d$) *time to generate the answer.*

Proof. We make use of the AttentionGeneration data structure (Algorithm 1) in Lemma D.1.

The generation process is an auto-regressive procedure, we define the following notations for better understanding. For $i \in [m]$, let $q_i, k_i \in \mathbb{R}^d$ denote the query vector of the *i*-th iteration, respectively. Note that q_i need to attend on both $K \in \mathbb{R}^{n \times d}$ and $\{k_1, k_2, \dots, k_{i-1}\}$.

For calculating the attention between q_i and $K \in \mathbb{R}^{n \times d}$, we just need to call AttentionGeneration .Inference $(q_i, 1)$ for once. Therefore, the running time for this part is $O(n^{4/5}d)$ time.

For calculating the attention between q_i and $\{k_1, k_2, \dots, k_{i-1}, k_i\}$, it takes $O(i \cdot d)$ time.

Therefore, for a single query q_i , the running time for getting the attention matrix $A \in \mathbb{R}^{1 \times (n+i)}$ is $(n^{4/5} + i) \cdot d$. Since there are only $n^{4/5} + i$ non-zero entries in A, it takes $n^{4/5} + i$ time to calculate $D^{-1}A$. Then, it takes $(n^{4/5} + i) \cdot d$ time to calculate $D^{-1}AV$. Since $i \leq m$, the total running time for calculating attention for a single query q_i is $O((n^{4/5} + m) \cdot d)$.

There are *m* queries in total. The running time for *m* queries is $O(mn^{4/5}d + m^2d)$.

Since we have $n \gg m$, the overall running time for the generation is $O(mn^{4/5}d)$.

E. Sparsity Analysis

To begin our analysis, we first examine the application of Bernstein's inequality to the matrix *K*: **Lemma E.1** (Bernstein on *K*). *If the following conditions hold:*

- Let the ReLU attention be defined as Definition 1.2.
- Let $Q \in \mathbb{R}^{m \times d}$ and $K, V \in \mathbb{R}^{n \times d}$ be defined as Definition 1.2.
- Let $b \in \mathbb{R}$ denote the threshold of ReLU activation, as defined in Definition 1.2.
- For $i \in [m]$, let k_i denote the number of non-zero entries in *i*-th row of $A \in \mathbb{R}^{m \times n}$.
- Assume each entry of K is from Gaussian $\mathcal{N}(0, \sigma_k^2)$
- Let $x \in \mathbb{R}^d$ denote a single row of $Q \in \mathbb{R}^{m \times d}$.
- Let $\sigma_a = \|x\|_2 \sigma_k / \sqrt{d}$.

Then, we can show that, with probability at least $1 - \exp(-\Omega(n \cdot \exp(-\frac{b^2}{2\sigma_a^2})))$, the number of non-zero entries \tilde{k}_i is at most $2n \cdot \exp(-\frac{b^2}{2\sigma_a^2})$. Namely, we have

$$\Pr[\widetilde{k}_i \leq 2n \cdot \exp(-\frac{b^2}{2\sigma_a^2})] \geq 1 - \exp(-\Omega(n \cdot \exp(-\frac{b^2}{2\sigma_a^2})))$$

Proof. For simplicity, for $i \in [n], j \in [d]$, we use $K_{i,j} \in \mathbb{R}$ to denote the (i, j)-th entry of $K \in \mathbb{R}^{n \times d}$. Let $r_i \in \{0, 1\}$ be the indicator function of $\langle x, K_{i,*} \rangle$. Then, we have $\tilde{k}_i = \sum_{j=1}^n r_j$. Since r_i is an indicator function, then we have

 $|r_i| \leq 1.$

By assumption, we have $K_{i,j} \sim \mathcal{N}(0, \sigma_k^2)$.

Let
$$\sigma_a = \|x\|_2 \cdot \sigma_k / \sqrt{d}$$
.

By the property of Gaussian distribution (Fact B.5), we have $\langle x, K_{i,*} \rangle \sim \mathcal{N}(0, d \cdot \sigma_a^2)$ and $\langle x, K_{i,*} \rangle / \sqrt{d} \sim \mathcal{N}(0, \sigma_a^2)$.

For any $i, j \in [n]$, by Fact B.6, we have $\langle x, K_{i,*} \rangle$ and $\langle x, K_{j,*} \rangle$ are independent, which implies r_i and r_j are independent.

By the tail bound of Gaussian distribution (Fact B.8), we have

$$\begin{aligned} \Pr[r_i = 1] &= \Pr[\langle x, K_{i,*} \rangle / \sqrt{d} \ge b] \\ &\leq \exp(-\frac{b^2}{2\sigma_a^2}), \end{aligned}$$

which implies

$$\mathbb{E}[r_i] \le \exp(-\frac{b^2}{2\sigma_a^2}),\tag{3}$$

and

$$\mathbb{E}[r_i^2] \le \exp(-\frac{b^2}{2\sigma_a^2})$$

which implies

$$\sum_{i=1}^{n} \mathbb{E}[r_i^2] \le n \cdot \exp(-\frac{b^2}{2\sigma_a^2})$$

Since we have $\widetilde{k}_i = \sum_{j=1}^n r_j$, by Eq. (3), we have

$$E[\widetilde{k}_i] \le n \cdot \exp(-\frac{b^2}{2\sigma_a^2}).$$

Let $k_0 := n \cdot \exp(-\frac{b^2}{2\sigma_a^2})$. By the Bernstein inequality (Lemma B.9), we have

$$\Pr[\tilde{k}_i \ge k_0 + t] \le \exp(-\frac{t^2/2}{k_0 + t/3})$$
(4)

We choose $t = k_0$, then we have

$$\Pr[k_i \ge 2k_0] \le \exp(-3k_0/8)$$

Then, we reach our conclusion: with probability at least $1 - \exp(-\Omega(n \cdot \exp(-\frac{b^2}{2\sigma_a^2})))$, the number of non-zero entries in each row of the attention matrix A is bounded by $\tilde{k}_i \leq 2n \cdot \exp(-\frac{b^2}{2\sigma_a^2})$.

We turn our attention to bounding $||x||_2$: Lemma E.2 ($||x||_2$ bound). *If the following conditions hold:*

- Let $Q \in \mathbb{R}^{m \times d}$ be defined as Definition 1.2.
- Let $x \in \mathbb{R}^d$ denote a single row of $Q \in \mathbb{R}^{m \times d}$.
- Assume each entry of Q is from $\mathcal{N}(0, \sigma_q^2)$.

Then, we can show that, for $t \ge 0$ with probability $1 - \exp(-t)$, $||x||_2$ is at most $\sqrt{3} \cdot (d+t)^{1/2} \cdot \sigma_q$. Namely, we have

$$\Pr[\|x\|_{2} \le \sqrt{3} \cdot (d+t)^{1/2} \cdot \sigma_{q}] \ge 1 - \exp(-t).$$

Proof. For simplicity, we use $x_i \in \mathbb{R}$ to denote the *i*-th entry of x.

By the assumption, we have $x_i \sim \mathcal{N}(0, \sigma_q^2)$.

Since $\|x\|_2^2 = \sum_{i=1}^d x_i^2$, by Chi-square tail bound (Lemma B.7), we have

$$\Pr[\|x\|_{2}^{2} - d\sigma_{q}^{2} \ge (2\sqrt{dt} + 2t)\sigma_{q}^{2}] \le \exp(-t)$$

which implies

$$\Pr[\|x\|_{2}^{2} \ge (2\sqrt{dt} + 2t + d)\sigma_{q}^{2}] \le \exp(-t).$$
(5)

Since we have $2\sqrt{dt} \le d + t$, Eq. (5) implies

$$\Pr[\|x\|_{2}^{2} \ge 3(d+t)\sigma_{q}^{2}] \le \exp(-t),$$

which is equivalent to

$$\Pr[\|x\|_{2} \ge \sqrt{3} \cdot (d+t)^{1/2} \cdot \sigma_{q}] \le \exp(-t)$$

| I | | _ | _ | |
|---|---|---|---|--|
| | 1 | | | |
| | | | | |

We can now present our formal sparsity analysis, which builds upon the previous lemmas: **Lemma E.3** (Sparsity analysis, formal version of Lemma 6.1). *If the following conditions hold:*

- Let the ReLU attention be defined as Definition 1.2.
- Let $Q \in \mathbb{R}^{m \times d}$ and $K, V \in \mathbb{R}^{n \times d}$ be defined as Definition 1.2.
- Let $b \in \mathbb{R}$ denote the threshold of ReLU activation, as defined in Definition 1.2.
- For $i \in [m]$, let \widetilde{k}_i denote the number of non-zero entries in *i*-th row of $A \in \mathbb{R}^{m \times n}$.
- Assume each entry of K is from Gaussian $\mathcal{N}(0, \sigma_k^2)$, and each entry of K is from Gaussian $\mathcal{N}(0, \sigma_q^2)$.
- Let $\delta \in (0, 1)$ denote the failure probability.
- Let $\sigma_a = 4 \cdot (1 + d^{-1} \log(m/\delta))^{1/2} \cdot \sigma_q \sigma_k$.
- Let $b = \sigma_a \cdot \sqrt{0.4 \log n}$.

Then, we can show that, with probability at least $1 - \delta$, for all $i \in [m]$, the number of non-zero entries of the *i*-th row \tilde{k}_i is at most $2n^{4/5}$.

Proof. This proof follows from applying union bound on Lemma E.1 and Lemma E.2.

By Lemma E.2, we have

$$\Pr[\|x\|_{2} \le \sqrt{3} \cdot (d+t)^{1/2} \cdot \sigma_{q}] \ge 1 - \exp(-t).$$
(6)

We choose $t = d + \log(m/\delta)$. Then, Eq. (6) implies

$$\Pr[\|x\|_{2} \le 4 \cdot (d + \log(m/\delta))^{1/2} \cdot \sigma_{q}] \ge 1 - \exp(-(d + \log(m/\delta))).$$
(7)

Let $\sigma_a = \|x\|_2 \cdot \sigma_k / \sqrt{d}$. By Eq.(7), we have $\sigma_a = 4 \cdot (1 + d^{-1} \log(m/\delta))^{1/2} \cdot \sigma_q \sigma_k$. By Lemma E.1, we have

$$\Pr[\widetilde{k}_i \le 2n \cdot \exp(-\frac{b^2}{2\sigma_a^2})] \ge 1 - \exp(-\Omega(n \cdot \exp(-\frac{b^2}{2\sigma_a^2}))).$$
(8)

Let $b = \sigma_a \cdot \sqrt{0.4 \log n}$. Then, Eq. (8) implies

$$\Pr[\tilde{k}_i \le 2n^{4/5}] \ge 1 - \exp(-O(n^{4/5}))$$
(9)

Since we have $n \gg d$, this implies

$$\exp(-O(n^{4/5})) \le \exp(-d) \tag{10}$$

Taking union bound over Eq. (7) and Eq. (9), we have

$$\Pr[\tilde{k}_{i} \leq 2n^{4/5}] \geq 1 - (\exp(-O(n^{4/5}) + \exp(-(d + \log(m/\delta)))))$$

= 1 - (exp(-O(n^{4/5}) + (\delta/m) \cdot exp(-d)))
\geq 1 - \delta/m. (11)

where the first step follows from the union bound, the second step follows from basic algebra, the third step follows from Eq. (10).

Since $x \in \mathbb{R}$ represents a single row of $Q \in \mathbb{R}^{m \times d}$, we already proved that for each fixed row of A, the \tilde{k}_i is at most $2n^{4/5}$ with probability $1 - \delta/m$.

Taking the union bound over *m* rows in *A*, then we can show that with probability $1 - \delta$, for all rows of *A*, that row's \tilde{k}_i is at most $2n^{4/5}$.

F. Running Time of Softmax Attention

In this section, we provide our results on reducing the running time of Softmax attention. We begin with introducing our result on Softmax attention generation.

Theorem F.1 (Running time of Softmax attention generation, formal version of Theorem 4.2). Let $Q \in \mathbb{R}^{m \times d}$, $K, V \in \mathbb{R}^{n \times d}$ and the Softmax attention Attn_s be defined in Definition 1.1. Let $\operatorname{NN}(r, q, K) \subseteq [n]$ and the Softmax attention with index set Attn_s be defined as Definition B.2. We choose the threshold $b \in \mathbb{R}$ in Algorithm 1 such that $R = \operatorname{NN}(n^{4/5}, q, K)$. Then, we can show that the Softmax attention with index set Attn_s achieves outstanding running time under the Softmax attention generation scenario: Suppose we have KV Cache $K, V \in \mathbb{R}^{n \times d}$. We want to generate a m length answer, where $m = \Theta(1)$. Algorithm 1 (replacing ReLU attention with Softmax attention) takes $O(mn^{4/5})$ time to generate the answer.

Proof. The Softmax attention generation scenario can be proved by substituting the ReLU attention Attn_r (Definition 1.2) with Softmax attention with index set $\widehat{\text{Attn}}_s$ (Definition B.2) in Algorithm 1 and Theorem 4.1.

Then, we move on to our result on Softmax full attention computation.

Theorem F.2 (Running time of Softmax full attention computation, formal version of Theorem 5.2). Let $Q \in \mathbb{R}^{m \times d}$, $K, V \in \mathbb{R}^{n \times d}$ and the Softmax attention Attn_s be defined in Definition 1.1. Let $NN(r, q, K) \subseteq [n]$ and the Softmax attention with index set $Attn_s$ be defined as Definition B.2. We choose the threshold $b \in \mathbb{R}$ in Algorithm 2 such that $R = NN(n^{4/5}, q, K)$. Then, we can show that the Softmax attention with index set $Attn_s$ achieves outstanding running time under full Softmax attention computation scenario: Suppose we have $m = \Theta(n)$. Algorithm 2 (replacing ReLU attention with Softmax attention) takes $O(n^{2-1/\lfloor d/2 \rfloor} d + n^{1+4/5} d)$ time to calculate the attention output. *Proof.* The Softmax full attention computation scenario can be proved by substituting the ReLU attention $Attn_r$ (Definition 1.2) with Softmax attention with index set $Attn_s$ (Definition B.2) in Algorithm 2 and Theorem 5.1.

G. Error Analysis of Softmax Attention

In this section, we provide an error analysis of the Softmax attention mechanism, deriving error bounds for the general case and a specific case with the massive activation property.

The following lemmas establish error bounds for Softmax attention when using index sets, formalizing the approximation error in attention computation.

Lemma G.1 (General error analysis of Softmax attention with index set, formal version of Lemma 6.5). *If the following conditions hold:*

- Let $Q \in \mathbb{R}^{m \times d}$, $K, V \in \mathbb{R}^{n \times d}$ and the Softmax attention $Attn_s$ be defined in Definition 1.1.
- Let $q \in \mathbb{R}^d$ denote a single row of $Q \in \mathbb{R}^{m \times d}$.
- Let α , $\overline{\alpha}$ and \widehat{Attn}_s be defined as Definition B.2.

Then we have

$$\|\mathsf{Attn}_s(q, K, V) - \widehat{\mathsf{Attn}}_s(q, K, V)\|_{\infty} \le \frac{2\overline{\alpha}}{\alpha} \cdot \|V\|_{\infty}$$

Proof. Recall that $\overline{R} = [n] \setminus R$ and $\widehat{K} = K_R \in \mathbb{R}^{r \times d}$ and $\widehat{V} = V_R \in \mathbb{R}^{r \times d}$ and $\overline{K} = K_{\overline{R}} \in \mathbb{R}^{(n-r) \times d}$ and $\overline{V} = V_{\overline{R}} \in \mathbb{R}^{(n-r) \times d}$ as defined in Definition B.1. Also, we have $\widehat{u} = \exp(q\widehat{K}^{\top}) \in \mathbb{R}^r$ and $\widehat{\alpha} = \langle \widehat{u}, \mathbf{1}_r \rangle \in \mathbb{R}$ and $\overline{u} = \exp(q\overline{K}^{\top}) \in \mathbb{R}^{n-r}$ and $\overline{\alpha} = \langle \overline{u}, \mathbf{1}_{n-r} \rangle \in \mathbb{R}$ as defined in Definition B.2. Then, we have

$$\begin{split} \|\mathsf{Attn}_{s}(q,K,V) - \widehat{\mathsf{Attn}}_{s}(q,K,V)\|_{\infty} \\ &= \|(\widehat{\alpha} + \overline{\alpha})^{-1}(\widehat{u}\widehat{V} + \overline{u}\overline{V}) - \widehat{\alpha}^{-1}\widehat{u}\widehat{V}\|_{\infty} \\ &\leq \|((\widehat{\alpha} + \overline{\alpha})^{-1} - \widehat{\alpha}^{-1})\widehat{u}\widehat{V}\|_{\infty} + \|(\widehat{\alpha} + \overline{\alpha})^{-1}\overline{u}\overline{V}\|_{\infty} \\ &\leq |(\widehat{\alpha} + \overline{\alpha})^{-1} - \widehat{\alpha}^{-1}| \cdot \|\widehat{u}\|_{1} \cdot \|\widehat{V}\|_{\infty} + (\widehat{\alpha} + \overline{\alpha})^{-1} \cdot \|\overline{u}\|_{1} \cdot \|\overline{V}\|_{\infty} \\ &= (\widehat{\alpha}^{-1} - (\widehat{\alpha} + \overline{\alpha})^{-1}) \cdot \widehat{\alpha} \cdot \|\widehat{V}\|_{\infty} + (\widehat{\alpha} + \overline{\alpha})^{-1} \cdot \overline{\alpha} \cdot \|\overline{V}\|_{\infty} \\ &\leq (\widehat{\alpha}^{-1} - (\widehat{\alpha} + \overline{\alpha})^{-1}) \cdot \widehat{\alpha} \cdot \|V\|_{\infty} + (\widehat{\alpha} + \overline{\alpha})^{-1} \cdot \overline{\alpha} \cdot \|V\|_{\infty} \\ &= 2(\widehat{\alpha} + \overline{\alpha})^{-1} \cdot \overline{\alpha} \cdot \|V\|_{\infty} \\ &= 2\alpha^{-1} \cdot \overline{\alpha} \cdot \|V\|_{\infty}, \end{split}$$

where the first step is by Definition B.2, the second step is by triangle inequality, the third step is by $||uV||_{\infty} \leq ||u||_1 \cdot ||V||_{\infty}$ for any vector u and conformable matrix V, and the fourth step is by definition of $\hat{\alpha}$ and $\overline{\alpha}$, i.e., $\hat{\alpha} = \langle \hat{u}, \mathbf{1}_r \rangle = ||\hat{u}||_1$ (note that each entry of \hat{u} is positive), the fifth step is by $\max\{\|\hat{V}\|_{\infty}, \|\overline{V}\|_{\infty}\} = \|V\|_{\infty}$, the sixth step in by simple calculation and the last step is by $\hat{\alpha} + \overline{\alpha} = \alpha$.

Building on this, we now present a more specific error analysis incorporating the massive activation property:

Theorem G.2 (Error analysis of Softmax attention with index set, formal version of Theorem 4.3). *If the following conditions hold:*

- Let $Q \in \mathbb{R}^{m \times d}$, $K, V \in \mathbb{R}^{n \times d}$ and the Softmax attention $Attn_s$ be defined in Definition 1.1.
- Let $q \in \mathbb{R}^d$ denote a single row of $Q \in \mathbb{R}^{m \times d}$.

- Let $\gamma \in [0, 1], \beta_1 \ge \beta_2 \ge 0$.
- Let the Softmax attention with index set \widehat{Attn}_s be defined as Definition B.2.
- Let $NN(r, q, K) \subseteq [n]$ denote the indices of top-r entries of qK.
- Let $R = \mathsf{NN}(n^{\gamma}, q, K) \subseteq [n]$, where $|R| = n^{\gamma}$.
- Assume the query q and key cache K have $(\gamma, \beta_1, \beta_2)$ massive activation property.

Then, we can show that

$$\|\widehat{\mathsf{Attn}}_s(q,K,V) - \mathsf{Attn}_s(q,K,V)\|_{\infty} \le \frac{2\|V\|_{\infty}}{n^{\gamma + (\beta_1 - \beta_2) \cdot \|q\|_2 - 1}}$$

Proof. Let $\alpha, \overline{\alpha}, \widehat{\alpha}$ be defined in Definition B.2. By Lemma G.1, we have

$$\|\mathsf{Attn}_s(q, K, V) - \widehat{\mathsf{Attn}}_s(q, K, V)\|_{\infty} \le \frac{2\overline{\alpha}}{\alpha} \cdot \|V\|_{\infty}.$$

By Definition B.3, we have

$$\widehat{\alpha} = \sum_{i \in \mathsf{NN}(n^{\gamma}, q, K)} \exp(\langle q, K_i \rangle)$$

$$\geq \sum_{i \in \mathsf{NN}(n^{\gamma}, q, K)} \exp(\|q\|_2 \beta_1 \log(n))$$

$$= n^{\gamma + \beta_1 \cdot \|q\|_2},$$

where the first step is by Definition of $\hat{\alpha}$, the second step is by Definition B.3 and Jensen inequality, and the last step is by simple calculation.

We also have

$$\overline{\alpha} = \sum_{i \in [n] \setminus \mathsf{NN}(n^{\gamma}, q, K)} \exp(\langle q, K_i \rangle)$$

$$\leq \sum_{i \in [n] \setminus \mathsf{NN}(n^{\gamma}, q, K)} \exp(\|q\|_2 \beta_2 \log(n))$$

$$\leq n^{1 + \beta_2 \cdot \|q\|_2},$$

where the first step is by Definition of $\overline{\alpha}$, the second step is by Definition B.3, and the last step is by simple calculation.

Finally, we finish the proof by the fact $\hat{\alpha} + \overline{\alpha} = \alpha$.