

# 000 001 002 003 004 005 TOWARDS SAMPLING DATA STRUCTURES FOR TENSOR 006 PRODUCTS IN TURNSTILE STREAMS 007 008 009

010 **Anonymous authors**  
011 Paper under double-blind review  
012  
013  
014  
015  
016  
017  
018  
019

## 020 ABSTRACT 021

022 This paper studies the computational challenges of large-scale attention-based  
023 models in artificial intelligence by introducing innovative sampling methods in the  
024 streaming setting. Inspired by the classical definition of the  $\ell_2$  sampler and the  
025 recent progress of the attention scheme in Large Language Models (LLMs), we  
026 propose the definition of the attention sampler. Our approach significantly reduces  
027 the computational burden of traditional attention mechanisms. We demonstrate  
028 the effectiveness of the attention sampler from a theoretical perspective, including  
029 space and update time. Additionally, our framework exhibits scalability and broad  
030 applicability across various model architectures and domains.  
031

## 032 1 INTRODUCTION 033

034 In recent years, the field of artificial intelligence has witnessed a significant paradigm shift with the  
035 advent of attention-based models, particularly in the domains of natural language processing and  
036 computer vision (Vaswani et al., 2017; Devlin et al., 2019; Liu et al., 2019; Yang et al., 2019; Brown  
037 et al., 2020; Zhang et al., 2022; Chowdhery et al., 2023; Touvron et al., 2023a;b; Inc., 2023; Manyika,  
038 2023). At the heart of these models lies the attention mechanism (Vaswani et al., 2017), which is a  
039 powerful tool for enhancing the performance of deep learning networks. It enables models to focus  
040 on relevant parts of the input data, thereby facilitating context-aware processing.  
041

042 However, as these models scale in size and complexity Zeng et al. (2024); Reid et al. (2024); Zhang  
043 et al. (2024); Dubey et al. (2024); Abdin et al. (2024), the computational demands of the attention  
044 mechanism increase exponentially, posing significant challenges in terms of efficiency and scalability  
045 Fu (2024). In particular, traditional attention mechanisms used in Transformer models (Vaswani  
046 et al., 2017) require computing attention weights across all elements of the input sequence, leading  
047 to a *quadratic* increase in computational complexity with respect to the sequence length (Alman &  
048 Song, 2023; Kacham et al., 2023; Han et al., 2024; Zandieh et al., 2023; Alman & Song, 2024a;b;  
049 2025a). This computational burden becomes particularly pronounced in large-scale applications. It  
050 hinders the usage of attention-based models in resource-constrained settings and limits their real-time  
051 processing capabilities.  
052

053 To deal with this problem, the core question we ask in this paper is:  
054

055 *Instead of computing all entries, can we recover the most important ones in efficient space and time?*  
056

057 **Attention Samplers.** We adopt the classical idea of sampling a dataset, which selects important  
058 items to represent the entire dataset. Sampling is a central and effective technique for analyzing  
059 large-scale datasets, which has broad application in the field of big data (Vitter, 1985; Gemulla et al.,  
060 2008; Cohen et al., 2011; 2014), including network traffic monitoring (Mai et al., 2006; Huang et al.,  
061 2007; Thottan et al., 2010), database management (Haas & Swami, 1992; Haas, 2016; Cohen & Geri,  
062 2019), and data summarization (Frieze et al., 2004; Aggarwal et al., 2009; Mahabadi et al., 2019;  
063 Indyk et al., 2020; Mahabadi et al., 2020). A well-known example is the  $\ell_2$  sampler first asked by  
064 Cormode et al. (2005) and studied by Monemizadeh & Woodruff (2010): given a vector  $x \in \mathbb{R}^n$ , we  
065 sample an index  $i \in [n]$  with probability  $\frac{x_i^2}{\|x\|_2^2}$ .  
066

067 To address the challenges in implementing large-scale attention schemes, we seek to sample the  
068 most important coordinates in attention computation, reducing computational overhead and computer  
069

storage. Inspired by the classical definition of the  $g$ -sampler on a vector, we propose the following attention sampler, which is thoroughly investigated in this paper.

**Definition 1.1** (Attention sampler). Given matrix  $A \in \mathbb{R}^{n \times d}$ , vector  $x \in \mathbb{R}^d$ , and a distribution function  $g$ , the attention sampler samples index  $i \in [n]$  with probability  $p_i = \frac{g((Ax)_i)}{\sum_{j=1}^n g((Ax)_j)}$ .

The motivation of our definition lies in the internal structure of the attention mechanism. Given input matrices  $A_1$  and  $A_2$ , the linear attention matrix is defined as  $A_1 X A_2^\top$ , where  $X = QK^\top$  is the fused key and query matrix. For linear self-attention,  $A_1$  and  $A_2$  are identical. Utilizing a well-known tensor product construction (Alman & Song, 2024a), we simplify the expression of the attention matrix to a matrix-vector product. Let  $A = A_1 \otimes A_2$ , and let  $x = \text{vec}(X)$ , the vectorized linear attention matrix turns out to be  $\text{vec}(A_1 X A_2^\top) = Ax$ . Here,  $\text{vec}$  denotes the vector representation of a matrix by concatenating the rows. Therefore, our attention sampler detects the dominant entry in the linear attention matrix, providing an effective approximation of the attention scheme.

Given unlimited space and time, the sampling problem is trivial since one can compute each entry explicitly and sample an index with the corresponding probability. However, as mentioned earlier, we are not granted unlimited resource in real-world applications, for instance, in resource-constrained settings or in real-time processing. This inspires us to investigate the attention sampler in the streaming model, where the input matrix  $A$ , the weight vector  $x$ , or both  $A$  and  $x$  arrive sequentially in a data stream, and the goal is to report a valid sample *at all times* using efficient space and update time. We study the turnstile data streams where the data can either be inserted or deleted at each time.

As databases handle increasingly vast and dynamic real-time data, the streaming model has emerged as a vital framework for designing algorithms to process massive, constantly evolving datasets. Examples include real-time analysis of social media streams, sensor data for smart infrastructure, live video processing, detection of distributed denial of service (DDoS) attacks, and efficient indexing and querying in large-scale databases. In this work, we combine the streaming model with attention mechanisms and construct novel efficient attention samplers, which identifies the critical coordinates in attention computation. Our contributions can be summarized as follows:

- For the softmax distribution  $\langle \exp(Ax), \mathbf{1}_n \rangle^{-1} \exp(Ax)$ , we prove an  $\Omega(n)$  space streaming sampler algorithm lower bound. (See Theorem 4.4)
- As the softmax distribution has a strong lower bound, we then provide upper bounds for polynomial type samplers, i.e.,  $L_2$  sampling from  $Ax$ :
  1. For updating  $A$  and fixed  $x$ , our sampler takes  $d \text{poly}(\frac{1}{\epsilon}, n)$  bits of space and update time (see Theorem 5.3).
  2. For updating  $A$  and fixed  $x$ , our sampler takes  $d \text{poly}(\frac{1}{\epsilon}, n)$  bits of space and  $O(1)$  update time (see Theorem 5.5).
  3. For updating both  $A$  and  $x$ , our sampler takes  $d \text{poly}(\frac{1}{\epsilon}, n)$  bits of space and update time (see Theorem 5.7).
- For updating both  $A$  and  $x$ , we also provide a lower bound of  $\Omega(d)$  space (see Theorem 6.2).
- Toward tensor generalization, where we have updating  $A_1 \in \mathbb{R}^{n \times d}$  or  $A_2 \in \mathbb{R}^{n \times d}$  for  $A = A_1 \otimes A_2 \in \mathbb{R}^{n^2 \times d^2}$  and fixed  $x \in \mathbb{R}^{d^2}$ , we sample  $(i_1, i_2) = i \in [n^2]$  approximately according to the  $\ell_2$  sampling distribution on  $Ax \in \mathbb{R}^{n^2}$  using  $O(nd)$  space,  $O(n)$  update time (see Theorem 7.6). Note that the trivial result takes  $O(n^2)$  space.

**Hardness of softmax attention.** Our lower bound in the first result demonstrates the hardness for computing or approximating the softmax attention. This aligns with the lower bound in Alman & Song (2023), where they show that approximating softmax attentions up to small entry-wise error requires subquadratic time in  $n$  assuming the Strong Exponential Time Hypothesis. These challenges motivate us to explore polynomial attention mechanisms. Previous work has investigated the performance of polynomial attention from both theoretical and empirical perspectives. For instance, the PolySketchFormer (Kacham et al., 2023) demonstrates that polynomial attention achieves model quality comparable to softmax attentions with efficient low-dimensional approximations. Furthermore, polynomial attention schemes perform competitively in various vision and NLP tasks, including the linear attention in (Koohpayegani & Pirsavash, 2024) and the polynomial attention in (Saratchandran et al., 2024). Building on these insights, we obtain efficient polynomial attention

108 samplers in the streaming model, whose space and update time have no dependence on  $n$  factors,  
 109 effectively recovering the key components in the polynomial attention matrix.  
 110

111 **Streaming attention mechanism.** Our polynomial attention samplers work in the streaming model,  
 112 which matches the core idea of streaming Large Language Model (LLM) introduced and studied by  
 113 Xiao et al. (2024), and recently gained increasing focus in LLM research and big-data analysis (Strati  
 114 et al., 2024; Yao et al., 2024; Shikhar et al., 2025; Xiao et al., 2025). The motivation is from long (or  
 115 infinite) sequence generation, e.g., a chat bot having a day-long conversation. When we apply LLMs  
 116 in these scenarios, we often encounter a efficiency-performance trade-off (Xiao et al., 2024): during  
 117 the decoding stage, attention-based methods cache all Key and Value (KV) pairs, which requires  
 118 excessive memory usage; in contrast, under restricted memory, the performance collapses once the  
 119 sequence length exceeds the cache size. To deal with these drawbacks, their method is to train the  
 120 models with a finite attention window to work on text of infinite length. Unlike Xiao et al. (2024),  
 121 our model is dynamic and data-driven, supporting both model weights and input tokens to constantly  
 122 change. We note that our sampler provides a correct attention sample at all times using efficient space.  
 123 Thus, we identify the important coordinates in attention computing without probing each KV pair,  
 124 which has high potential in enhancing the performance of streaming LLMS.  
 125

126 **Sparse attention mechanism.** Another practical relevance of our attention sampler is sparse  
 127 attention mechanism. The attention matrix has been shown to be naturally sparse empirically and  
 128 theoretically (see e.g. (Deng et al., 2024b)). Based on this observation, researchers seek to reduce  
 129 computation by sampling the attention layers (Child et al., 2019; Kitaev et al., 2020; Wang et al.,  
 130 2020; Alman & Song, 2023; Brand et al., 2024; Deng et al., 2023c; Lai et al., 2025; Xiao et al.,  
 131 2025; Zhang et al., 2025). In general, they construct a *sparse mask* that selects the importance  
 132 entries in the attention multiplications while others are zeroed out. Then, they compute the partial  
 133 attention corresponding to those in the sparse mask. Specifically, (Xiao et al., 2025) explores the  
 134 sparse attention with streaming heads. Our attention sampler recovers large coordinates from the  
 135 attention matrix given specific streamed inputs and weights. Thus, the sampler serves as an efficient  
 136 subroutine in their sparse-attention sampling schemes, evaluating and enhancing the effectiveness of  
 137 their construction of the sparse mask.  
 138

139 **Streaming algorithms.** In addition, our sampler can be integrated into inner product computation  
 140 (see e.g. Woodruff & Zhou (2021)), which is a cornerstone for model training and attention computa-  
 141 tion. In fact, classical  $L_p$  samplers also serve as black-box subroutines in many other streaming  
 142 algorithms, including finding heavy hitters,  $F_p$  moment estimation, and cascaded norm approximation  
 143 Andoni et al. (2011); Jowhari et al. (2011); Jayaram & Woodruff (2021); Woodruff & Zhou (2021).  
 144 Therefore, our attention sampler can be applied to discover essential properties of the attention  
 145 scheme, e.g., the norm of the attention matrix.  
 146

## 2 RELATED WORK

147 In this section, we present related work in sampling and tensor sketch.  
 148

149 **On sampling.** Given a vector  $v \in \mathbb{R}^n$  and a distribution function  $g$ , recall that the classical  $g$ -  
 150 sampler samples index  $i \in [n]$  with probability  $p_i = \frac{g(v_i)}{\sum_{j=1}^n g(v_j)}$ . A well-known example is the  $L_p$   
 151 sampling defined by  $g(z) = |z|^p$  for  $p \geq 0$ . The existence of such a  $L_p$  sampler algorithms first posed  
 152 as a question by Cormode et al. (2005) in 2005. Monemizadeh & Woodruff (2010) partially answered  
 153 this question in the affirmative by giving an  $L_p$  sampler using polylogarithmic space for  $p \in [1, 2]$ ,  
 154 although the sampling probabilities were distorted by a multiplicative  $(1 + \epsilon)$  factor and an additive  
 155  $\frac{1}{\text{poly}(n)}$  factor. We note that the sampler is perfect if there is no  $\epsilon$ -multiplicative distortion; it is truly  
 156 perfect if there is no additive distortion, i.e., the sampling probability is exact. The space requirements  
 157 of the algorithm were subsequently improved (Andoni et al., 2011; Jowhari et al., 2011) and extended  
 158 to other choices of index domain  $U$  and weight function  $W$  (Cohen & Geri, 2019; Mahabadi et al.,  
 159 2020; 2022), while retaining a multiplicative distortion in the sampling probability. Surprisingly,  
 160 Jayaram & Woodruff (2021) showed that it is possible to achieve no perfect samplers while using  
 161 polylogarithmic space, while conversely Jayaram et al. (2022) showed that truly perfect samplers  
 162 would require linear space, essentially closing the line of work studying the space complexity of

*L<sub>p</sub>* samplers for  $p \in [1, 2]$ . It should be noted however, achieving such guarantees (no additive distortion) in sub-polynomial update time while retaining the space guarantees remains an intriguing open question (Jayaram et al., 2022). For the other regime of  $p > 2$ , recently, Woodruff et al. (2025) complemented the results by providing efficient perfect  $L_p$  samplers for  $p > 2$ . Swartworth et al. (2025) achieved perfect samplers with polylogarithmic update time for  $p > 2$ , improving on the previous update time. For a more comprehensive background on samplers, we refer to the survey by Cormode & Jowhari (2019).

**On tensors.** In the realm of tensor decomposition, the canonical polyadic (CP) decomposition, specifically the CANDECOMP/PARAFAC method, stands out for its unique ability to break down tensors into rank-1 tensors in a singular way, distinct from matrix decomposition (Harshman, 1970; Song et al., 2016). This method, having applications in computational neuroscience, data mining, and statistical learning (Wang et al., 2015), emphasizes the rigidity and uniqueness of tensor decomposition. Earlier studies (Tsourakakis, 2010; Phan et al., 2013; Choi & Vishwanathan, 2014; Huang et al., 2013; Kang et al., 2012; Wang et al., 2014; Bhojanapalli & Sanghavi, 2015) have delved into efficient tensor decomposition methods. Subsequent works introduced methods for fast orthogonal tensor decomposition using random linear sketching techniques (Wang et al., 2015) and explored symmetric orthogonally decomposable tensors' properties, integrating spectral theory (Robeva, 2016; Robeva & Seigal, 2017). Additionally, importance sampling for quicker decomposition was proposed (Song et al., 2016). (Deng et al., 2023a) studies the tensor cycle low rank approximation problem.

In algebraic statistics, tensor decompositions are linked to probabilistic models, particularly in determining latent variable models' identifiability through low-rank decompositions of specific moment tensors (Allman et al., 2009a; Rhodes & Sullivant, 2012). Kruskal's theorem (Kruskal, 1977) was pivotal in ascertaining the precision of model parameter identification. However, this approach, assuming an infinite sample size, does not provide the minimum sample size for learning model parameters within given error bounds. A more robust uniqueness guarantee is needed to ensure that the low-rank decomposition of an empirical moment tensor approximates that of an actual moment tensor, thus offering more insight into empirical moment tensors' decomposition.

**Roadmap.** In Section 3, we provide some standard notations and definitions in literature. In Section 4, we study the exponential sampler. In Section 5, we study the streaming upper for the  $\ell_2$  sampling problem, i.e., sampling coordinates from a vector  $Ax$ , where  $A$  and  $x$  may be updated across a data stream. In Section 6, we present lower bounds for the same  $\ell_2$  sampling problem. In Section 7, we discuss the tensor sampling problem.

### 3 PRELIMINARIES

For any positive integer  $n$ , we use  $[n]$  to denote the set  $\{1, 2, \dots, n\}$ . We use  $\mathbb{E}[\cdot]$  to denote the expectation. We use  $\Pr[\cdot]$  to denote the probability. We use  $\mathbf{1}_n$  to denote a length- $n$  vector where all the entries are ones. Given two length- $n$  vectors, we use  $\langle x, y \rangle$  to denote the inner product between  $x$  and  $y$ , i.e.,  $\langle x, y \rangle := \sum_{i=1}^n x_i y_i$ . For a vector  $x \in \mathbb{R}^n$ , we use  $\exp(x) \in \mathbb{R}^n$  to denote a vector that has length  $n$  and the  $i$ -th entry is  $\exp(x_i)$ . For a matrix  $A$ , we use  $\exp(A)$  to denote the matrix that  $(i, j)$ -th coordinate is  $\exp(A_{i,j})$ . For a vector  $x$ , we use  $\|x\|_2 := (\sum_{i=1}^n x_i^2)^{1/2}$ . We use  $\|x\|_1 := \sum_{i=1}^n |x_i|$ . We use  $\|x\|_0$  to denote the  $\ell_0$  norm of  $x$ , which is the number of nonzero entries in  $x$ . We use  $\|x\|_\infty$  to denote the  $\ell_\infty$  norm of  $x$ , which is  $\max_{i \in [n]} |x_i|$ .

Let  $n_1, n_2, d_1, d_2$  be positive integers. Let  $A \in \mathbb{R}^{n_1 \times d_1}$  and  $B \in \mathbb{R}^{n_2 \times d_2}$ . We define the Kronecker product between matrices  $A$  and  $B$ , denoted  $A \otimes B \in \mathbb{R}^{n_1 n_2 \times d_1 d_2}$ , as  $(A \otimes B)_{(i_1-1)n_2+i_2, (j_1-1)d_2+j_2}$  is equal to  $A_{i_1, j_1} B_{i_2, j_2}$ , where  $i_1 \in [n_1], j_1 \in [d_1], i_2 \in [n_2], j_2 \in [d_2]$ .

We use  $\text{poly}(n)$  to denote  $n^C$  where  $C > 1$  is some constant. For any function  $f$ , we use  $\tilde{O}(f)$  to denote  $f \cdot \text{poly}(\log f)$ . For two sets  $A$  and  $B$ , we use  $A \cap B$  to denote their intersection. We use  $|A \cap B|$  to denote the cardinality of  $A \cap B$ . We use  $A \cup B$  to denote the union of  $A$  and  $B$ .

**TensorSketch.** We next define TensorSketch (Pagh, 2013), which has been extensively used in many sketching and optimization problems (Diao et al., 2018; Song et al., 2019; Diao et al., 2019; Ahle et al., 2020; Song et al., 2021; 2024; 2022; Zhang, 2022; Song et al., 2023b). Song et al. (2022)

216 defined **TensorSparse** by composing Sparse embedding (Nelson & Nguyn, 2013; Cohen, 2016)  
 217 with a tensor operation (Pagh, 2013).

218 **Definition 3.1** (**TensorSparse**, see Definition 7.6 in Song et al. (2022)). Let  $h_1, h_2 : [n] \times [s] \rightarrow [m/s]$  be  
 219  $O(\log 1/\delta)$ -wise independent hash functions and let  $\sigma_1, \sigma_2 : [n] \times [s] \rightarrow \{\pm 1\}$  be  
 220  $O(\log 1/\delta)$ -wise independent random sign functions. Then, the degree two tensor sparse trans-  
 221 form,  $S : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}^m$  is given as:

$$222 \quad R_{r,(i,j)} = \exists k \in [s] : \sigma_1(i, k)\sigma_2(j, k)/\sqrt{s} \cdot \mathbf{1}[(h_1(i, k) + h_2(j, k)) \bmod m/s] + (k-1)m/s = r$$

224 For  $s = 1$ , the above definition becomes **TensorSketch** (Pagh, 2013).

## 226 4 EXPONENTIAL SAMPLER

228 In this section, we define and consider exponential samplers. We then show strong space lower  
 229 bounds for achieving such a data structure when the input dataset arrives in a data stream.

231 Let us firstly describe the offline version:

232 **Definition 4.1** (Exponential sampler). Given matrix  $A \in \mathbb{R}^{n \times d}$  and  $x \in \mathbb{R}^d$ , the goal is to sample  
 233 index  $i \sim [n]$  with probability  $p_i = \langle \exp(Ax), \mathbf{1}_n \rangle^{-1} \cdot \exp(Ax)_i$ , where  $\mathbf{1}_n$  denotes a length- $n$   
 234 vector,  $\exp(Ax) \in \mathbb{R}^n$  denotes a length- $n$  vector with  $\exp(Ax)_i = \exp((Ax)_i)$ , and  $\exp(z)$  is the  
 235 usual exponential function.

236 Now, consider  $y = Ax \in \mathbb{R}^n$ , where either  $A$  or  $x$ , or both are arriving in a data stream, we use the  
 237 following definition for each of the various cases:

238 **Definition 4.2.** Let  $C > 0$  be any fixed constant and let  $C_0 \in [n^{-C}, n^C]$ . Let  $y$  be a vector. Then  
 239 the exponential sampler outputs an index  $j^*$  such that for all  $i \in [n]$ ,  $\Pr[j^* = i] = C_0 \cdot \frac{\exp(y_i)}{\langle \exp(y), \mathbf{1}_n \rangle}$ .

241 We first recall the (two-party) set-disjointness communication problem  $\text{SetDisj}_n$ , in which two parties  
 242 Alice and Bob have subsets  $A$  and  $B$ , respectively, of  $[n]$ . Note that we can equivalently view  $A$   
 243 and  $B$  as binary vectors in  $n$ -dimensional space, serving as the indicator vector for whether each  
 244 index  $i \in [n]$  is in the player’s input subset. The task for the players is to determine whether there  
 245 exists a common element in their intersection, i.e., whether there exists  $i \in [n]$  such that  $i \in (A \cap B)$   
 246 or equivalently,  $A_i = B_i = 1$ . In fact, the problem promises that either the inputs are completely  
 247 disjoint,  $|A \cap B| = 0$  or the inputs contain only a single coordinate in their intersection,  $|A \cap B| = 1$ .  
 248 We recall the following standard communication complexity result of set-disjointness.

249 **Theorem 4.3** (Kalyanasundaram & Schnitger (1992); Razborov (1992); Bar-Yossef et al. (2004)).  
 250 Any protocol that solves the set-disjointness problem  $\text{SetDisj}_n$  with probability at least  $\frac{3}{4}$  requires  
 251  $\Omega(n)$  bits of total communication.

252 We show that even a sampler that relaxes the probability distribution defined in Definition 4.2 up to a  
 253 factor of  $n^C$  is infeasible in the streaming model.

254 **Theorem 4.4.** Let  $y \in \mathbb{R}^n$  that arrives as a data stream and let  $C > 0$  be a constant. Then any  
 255 algorithm that samples an index  $i \in [n]$  with probability proportional to  $p_i = \frac{\exp(y_i)}{\langle \exp(y), \mathbf{1}_n \rangle}$  must use  
 256  $\Omega(n)$  bits of space, even if the sampling probabilities are allowed to be distorted by as large as  $n^C$   
 257 and even if  $\|y\|_\infty = O(\log n)$ .

258 *Proof.* Let  $A, B \in \{0, 1\}^n$  be input vectors from the set disjointness problem, so that the goal is  
 259 to determine whether there exists  $i \in [n]$  such that  $A_i = B_i = 0$ . Observe that Alice and Bob can  
 260 multiply  $A$  and  $B$  by  $100C \log n$  for some constant  $C > 0$ . Now, note that in the disjoint case, we have  
 261 that  $\|A + B\|_\infty = 100C \log n$  and in the non-disjoint case, we have that  $\|A + B\|_\infty = 200C \log n$ .  
 262 In particular, in the non-disjoint case, there exists  $i \in [n]$  such that  $A_i + B_i = 200C \log n$  and for all  
 263  $j \neq i$ , we have that  $A_j + B_j \leq 100C \log n$ . Hence, in the non-disjoint case, any exponential sampler  
 264 will output  $i$  with probability proportional to  $\exp(200C \log n)$  and output  $j \neq i$  with probability  
 265 proportional to  $n \cdot \exp(100C \log n)$ . Even if the sampling probabilities are distorted by a factor of  
 266  $n^C$ , any exponential sampler would output  $i$  with probability at least  $\frac{3}{4}$ .

267 Thus, Alice and Bob can use such a data structure to sample an index  $i$  and then check whether  
 268  $A_i = B_i = 1$ . In particular, Alice can first create a data stream encoding the vector  $A$ , run the

270 sampling algorithm on the data stream, and then pass the state of the algorithm to Bob. Bob can  
 271 then create another portion of the data stream encoding an addition of the vector  $B$ , take the state of  
 272 the algorithm from Alice, run the sampling algorithm on the portion of the data stream, and query  
 273 the algorithm for an index  $i$ . Bob can then take the index and pass it to Alice, and the two parties  
 274 can finally communicate whether  $A_i = B_i = 1$ , thereby solving set-disjointness with probability at  
 275 least  $\frac{3}{4}$ . Note that the communication of the protocol is the space used by the sampling algorithm.  
 276 Therefore by Theorem 4.3, such a sampler must use  $\Omega(n)$  bits of space.  $\square$

## 278 5 $\ell_2$ SAMPLER UPPER BOUND WITH $A$ AND $x$

280 In this section, we describe a standard data structure for  $\ell_2$  sampling. We start with providing the  
 281 definition of  $\ell_2$  sampler as follows,

283 **Definition 5.1.** Let  $n$  denote a positive integer. Let  $\epsilon \geq 0$  denote a parameter. In  $\ell_2$  sampling, we  
 284 receive each coordinate of  $y \in \mathbb{R}^n$  in a turnstile data stream, and the goal is to output an index  $I \in [n]$   
 285 at all times such that for each  $j \in [n]$ ,  $\Pr[I = j] = (1 \pm \epsilon) \cdot \frac{|y_j|^2}{\|y\|_2^2} + 1/\text{poly}(n)$ .

287 We describe various instantiations of the  $\ell_2$  sampler for sampling entries from a vector  $Ax \in \mathbb{R}^n$ ,  
 288 based upon whether the matrix  $A \in \mathbb{R}^{n \times d}$  is updated during the data stream, whether the vector  
 289  $x \in \mathbb{R}^d$  is updated during the data stream, or both.

### 290 5.1 $A$ IS UPDATED DURING THE STREAMING AND $x$ IS FIXED

292 In this section, we describe the construction of an  $\ell_2$  sampler for sampling coordinates of the vector  
 293  $Ax \in \mathbb{R}^n$ , in the setting where the vector  $x \in \mathbb{R}^d$  is fixed, but the entries of  $A \in \mathbb{R}^{n \times d}$  are evolving  
 294 as the data stream progresses.

295 **Definition 5.2** (Updating  $A$  and fixed  $x$ ). In this setting, we assume  $x \in \mathbb{R}^d$  is fixed, we receive  
 296 updates to the entries of  $A \in \mathbb{R}^{n \times d}$  in a turnstile data stream. Then for  $y = Ax$ , we want a data  
 297 structure that produces the  $\ell_2$  sampling guarantee for  $y$ .

299 We remark that a turnstile data stream means that each update of the data stream can increase or  
 300 decrease a single entry of  $A$ .

301 In this work, we are interested in the regime of  $n \gg d$ . Then we have the following guarantee:

303 **Theorem 5.3.** Suppose  $y = Ax$ , for  $x \in \mathbb{R}^d$ , which is fixed, and  $A \in \mathbb{R}^{n \times d}$ , which is defined by  
 304 a turnstile stream. There exists an  $\ell_2$ -attention sampler that uses  $d \log n + \text{poly}(\frac{1}{\epsilon}, \log n)$  bits of  
 305 space and returns  $I \in [n]$  such that  $\Pr[I = j] = (1 \pm \epsilon) \cdot \frac{|y_j|^2}{\|y\|_2^2} + 1/\text{poly}(n)$ . The update time of the  
 306 data structure is  $d \text{ poly}(\frac{1}{\epsilon}, \log n)$ .

308 *Proof.* Recall that existing approximate  $\ell_2$  samplers, e.g., Algorithm 2 maintains a linear sketch  $\Phi y$ ,  
 309 where  $\Phi \in \mathbb{R}^{m \times n}$ , for  $m = \text{poly}(\frac{1}{\epsilon}, \log n)$ . We have  $y = Ax$ , where  $x \in \mathbb{R}^d$  is fixed but  $A \in \mathbb{R}^{n \times d}$   
 310 is defined through turnstile updates. Nevertheless, we can maintain the state of  $\Phi Ax$ . In particular,  
 311 whenever we receive an update in  $A_{i,j}$  by  $\Delta$ , then we can compute  $\Phi e_i e_j^\top \Delta x$  to update the sketch  
 312  $\Phi Ax$ . To analyze the space complexity, observe that storing  $\Phi Ax$  requires  $O(m)$  words of space and  
 313  $x$  requires  $d$  words of space, which is  $d \log n + \text{poly}(\frac{1}{\epsilon}, \log n)$  bits of space in total. Moreover, each  
 314 update to  $A_{i,j}$  can change all entries of  $\Phi Ax$ , so the update time is  $O(md) = d \text{ poly}(\frac{1}{\epsilon}, \log n)$ .  
 315  $\square$

### 318 5.2 $x$ IS UPDATED DURING THE STREAMING AND $A$ IS FIXED

320 We next consider the setting where the vector  $x \in \mathbb{R}^d$  is updated as the data stream progresses, but  
 321 the entries of  $A \in \mathbb{R}^{n \times d}$  are fixed.

322 **Definition 5.4** (Fixed  $A$  and updating  $x$ ). We assume  $A \in \mathbb{R}^{n \times d}$  is fixed, we receive updates to  
 323  $x \in \mathbb{R}^d$  in a turnstile data stream. Then for  $y = Ax$ , we want a data structure that produces the  $\ell_2$   
 324 sampling guarantee for  $y$ .

324 We have the following algorithmic guarantees for this setting:  
 325

326 **Theorem 5.5.** *Suppose  $y = Ax$ , for  $A \in \mathbb{R}^{n \times d}$ , which is fixed, and  $x \in \mathbb{R}^n$ , which is defined by  
 327 a turnstile stream. There is an  $\ell_2$ -attention sampler that uses  $d \cdot \text{poly}(\frac{1}{\epsilon}, \log n)$  bits of space and  
 328 returns  $I \in [n]$  such that  $\Pr[I = j] = (1 \pm \epsilon) \cdot \frac{|y_j|^2}{\|y\|^2} + 1/\text{poly}(n)$ . The update time of the data  
 329 structure is  $O(1)$ .*

330 *Proof.* Again recall that existing approximate  $\ell_2$  samplers, e.g., Algorithm 2 maintains a linear sketch  
 331  $\Phi y$ , where  $\Phi \in \mathbb{R}^{m \times n}$ , for  $m = \text{poly}(\frac{1}{\epsilon}, \log n)$ . Since  $y = Ax$ , but  $A \in \mathbb{R}^{n \times d}$  is too large to store,  
 332 while  $x \in \mathbb{R}^n$  is defined through turnstile updates, we can instead maintain the sketch  $\Phi A$  and the  
 333 vector  $x$  and compute  $\Phi Ax = \Phi y$  after the stream concludes. Note that storing  $\Phi A$  requires  $O(md)$   
 334 words of space and  $x$  requires  $d$  words of space, which is  $d \cdot \text{poly}(\frac{1}{\epsilon}, \log n)$  bits of space in total.  
 335 Moreover, each update to  $x$  changes a single entry, so the update time is  $O(1)$ .  $\square$

### 336 5.3 BOTH $A$ AND $x$ ARE UPDATED DURING THE STREAMING

337 Finally, we consider the setting where both the vector  $x \in \mathbb{R}^d$  and the entries of  $A \in \mathbb{R}^{n \times d}$  can be  
 338 changed by updates from the data stream.

339 **Definition 5.6** (Updating  $A$  and updating  $x$ ). In this setting, we receive updates to both  $A \in \mathbb{R}^{n \times d}$   
 340 and  $x \in \mathbb{R}^d$  in a turnstile data stream. Then for  $y = Ax$ , we want a data structure that provides the  $\ell_2$   
 341 sampling guarantee for  $y$ .

342 We have the following guarantees:

343 **Lemma 5.7** (Upper Bound). *Suppose  $y = Ax$ , for  $A \in \mathbb{R}^{n \times d}$  and  $x \in \mathbb{R}^d$ , which are each defined in  
 344 a stream through turnstile updates. There exists an  $\ell_2$ -attention sampler that uses  $d \cdot \text{poly}(\frac{1}{\epsilon}, \log n)$   
 345 bits of space and returns  $I \in [n]$  such that  $\Pr[I = j] = (1 \pm \epsilon) \cdot \frac{|y_j|^2}{\|y\|^2} + 1/\text{poly}(n)$ . The update time  
 346 is  $\text{poly}(\frac{1}{\epsilon}, \log n)$ .*

347 *Proof.* As before, recall that existing approximate  $\ell_2$  samplers, e.g., Algorithm 2 maintains a linear sketch  
 348  $\Phi y$ , where  $\Phi \in \mathbb{R}^{m \times n}$ , for  $m = \text{poly}(\frac{1}{\epsilon}, \log n)$ . Since  $y = Ax$ , but now both  $A \in \mathbb{R}^{n \times d}$   
 349 and  $x \in \mathbb{R}^d$  are defined through turnstile updates, we can instead maintain the sketch  $\Phi A$  and the  
 350 vector  $x$  and compute  $\Phi Ax = \Phi y$  after the stream concludes. Observe that maintaining  $\Phi A$  requires  
 351  $O(md)$  words of space and  $x$  requires  $d$  words of space, which is  $d \cdot \text{poly}(\frac{1}{\epsilon}, \log n)$  bits of space in  
 352 total. Each update to  $A$  can change all  $m$  entries of in a single column of  $\Phi A$ , while each update to  $x$   
 353 changes a single entry. Hence, the update time is  $\text{poly}(\frac{1}{\epsilon}, \log n)$ .  $\square$

## 354 6 $\ell_2$ SAMPLER LOWER BOUND (WITH $A$ AND $x$ )

355 In this section, we give lower bounds for  $\ell_2$  sampling from a vector  $y = A^{\otimes p}x$ , when either  $A$  or  $x$   
 356 are updated in a data stream. We show that in any of these cases, the general problem is substantially  
 357 more difficult than the previous case where  $p = 1$ .

358 We first recall the Index problem for one-way communication. In the  $\text{INDEX}_n$  problem, Alice receives  
 359 a vector  $v \in \{0, 1\}^n$  and Bob receives a coordinate  $i \in [n]$ . The goal is for Bob to compute  $v_i$  with  
 360 probability at least  $\frac{3}{4}$ , given some message  $\Pi$  from Alice. We recall the following communication  
 361 complexity lower bounds for Index.

362 **Theorem 6.1** (Kremer et al. (1999)). *Any protocol that solves  $\text{INDEX}_n$  with probability at least  $\frac{3}{4}$   
 363 requires  $\Omega(n)$  bits of communication.*

364 **Lemma 6.2** (Lower Bound). *Any streaming algorithm that solves problem defined as Definition 5.6  
 365 will require  $\Omega(d)$  space.*

366 *Proof.* Suppose Alice receives a vector  $v \in \{0, 1\}^d$ . Then Alice creates the diagonal matrix  $M \in$   
 367  $\{0, 1\}^{d \times d}$  so that the  $j$ -th diagonal entry of  $A$  is  $v_j$ , for all  $j \in [n]$ . Finally, Alice creates  $A \in$   
 368  $\mathbb{R}^{(d+1) \times d}$  by appending the row consisting of  $\frac{1}{10^{10}}$  in all of its  $d$  entries to  $M$ . Suppose Bob receives  
 369 the coordinate  $i \in [d]$  and wants to determine  $v_i$ . Then Bob can set  $x$  to be the elementary vector

378  $e_i \in \mathbb{R}^d$ , which has a 1 in its  $i$ -th coordinate and zeros elsewhere. Observe that by construction,  $Ax$   
 379 is the  $i$ -th column of  $A$ . If  $v_i = 1$ , then the  $i$ -th column of  $A$  consists of a 1 in the  $i$ -th entry,  $\frac{1}{10^{10}}$  in  
 380 the  $(d+1)$ -st entry, and zeros elsewhere. Hence, a sampler with the desired properties will output  
 381  $i$  with probability at least  $\frac{3}{4}$ . Similarly, if  $v_i = 0$ , then the  $i$ -th column of  $A$  consists of  $\frac{1}{10^{10}}$  in  
 382 the  $(d+1)$ -st entry and zeros elsewhere. Thus, the sampler with the desired properties will output  $d+1$   
 383 with probability 1. Bob can therefore distinguish between these two cases with probability at least  $\frac{3}{4}$ ,  
 384 thereby solving  $\text{INDEX}_d$  with probability at least  $\frac{3}{4}$ . Therefore, by Theorem 6.1, such a sampler must  
 385 use at least  $\Omega(d)$  space.  $\square$

386

387 In fact, we show that if  $y = A^{\otimes p}x$ , where  $A \in \mathbb{R}^{n \times n}$  so that  $A^{\otimes p} \in \mathbb{R}^{n^p \times n^p}$  denotes the  $p$ -wise  
 388 self-tensor and  $x \in \mathbb{R}^{n^p}$ , then actually  $L_2$  sampling from  $y$  uses  $\Omega(n)$  bits of space.

389

390 **Lemma 6.3.** *Let  $A \in \mathbb{R}^{n \times n}$  and  $A^{\otimes p} \in \mathbb{R}^{n^p \times n^p}$  denote the  $p$ -wise self-tensor. Let  $y = A^{\otimes p}x$ , so  
 391 that  $x \in \mathbb{R}^{n^p}$ . Then even if all the entries of  $x$  arrive in a data stream followed by all the entries of  $A$ ,  
 392  $L_2$  sampling from  $y$  requires  $\Omega(n)$  bits of space.*

393

394 *Proof.* Let  $S \in \{0, 1\}^n$  be an instance of  $\text{INDEX}_n$ . Suppose Alice creates the diagonal matrix  $A$   
 395 with exactly  $S$  being the entries across its diagonal, i.e.,  $A_{1,1} = S_1, \dots, A_{n,n} = S_n$ . Bob has an  
 396 index  $i \in [n]$ , and sets the vector  $x$  to be the elementary vector  $\mathbf{e}_j$ , where  $j = i \cdot n^{p-1}$ . Then by  
 397 construction  $Ax$  is the all zeros vector if  $S_i = 0$  and otherwise there is a nonzero entry, which allows  
 398 Alice and Bob to solve  $\text{INDEX}_n$ . Hence,  $L_2$  sampling from  $y$  requires  $\Omega(n)$  bits of space.  $\square$

399

## 7 THE TENSOR VERSION PROBLEM

400

401 In this section, we further consider sampling from a tensor product. We provide the tensor notations  
 402 and objects.

403

404 **Definition 7.1.** Let  $A_1 \in \mathbb{R}^{n \times d}$ , let  $A_2 \in \mathbb{R}^{n \times d}$ , we define  $\mathbf{A} = A_1 \otimes A_2 \in \mathbb{R}^{n^2 \times d^2}$ . Let  $x \in \mathbb{R}^{d^2}$ .  
 405 Let  $A_i \in \mathbb{R}^{n \times d^2}$  denote the  $i$ -th block of  $\mathbf{A}$ .

406

407 **Definition 7.2** (fixed  $x$ , Streaming Sampler for one of  $A_1$  and  $A_2$  is updating.). We assume  $x \in \mathbb{R}^{d^2}$   
 408 is fixed. We assume that (1) one of  $A_1$  and  $A_2$  is updating, (2) one of  $A_1$  and  $A_2$  is fixed. Let  $y = \mathbf{A}x$ ,  
 409 we want  $\ell_2$  sampling guarantee for sampling one coordinate in  $y_i \in \mathbb{R}^{n^2}$  for all  $i \in [n^2]$ .

410

411 To motive this model, recall that the tensor product  $(A_1 \otimes A_2)x$  equals to the linear cross-attention  
 412 matrix  $A_1 Q K^\top A_2^\top$ , where  $W_Q = A_1 Q$  is the projected query matrix and  $W_K = A_2 K$  is the  
 413 projected key matrix. Our model addresses a practical scenario involving real-time contextual  
 414 processing with a static reference dataset. In this setting,  $W_k$  is precomputed by the language model,  
 415 representing a static dataset such as embeddings of a knowledge base, user profiles, or multimedia  
 416 features. Then, the rows of matrix  $A_1$  arrive as a data stream, representing real-time data queries.  
 417 Thus, our attention sampler efficiently captures the important entries in the dynamic query dataset.

418

419 We use the following formulation of Nisan’s pseudorandom generator to derandomize our algorithm.

420

421 **Theorem 7.3** (Nisan’s PRG, Nisan (1992)). *Suppose  $\mathcal{A}$  is an algorithm that requires  $S = \Omega(\log n)$   
 422 bits of space and  $R$  random bits. Then there exists a pseudorandom generator for  $\mathcal{A}$  that succeeds  
 423 with probability  $1 - 1/\text{poly}(n)$  and uses  $O(S \log R)$  bits of space.*

424

425 **Algorithm 1** We build on algorithm based on  $S(x_1 \otimes x_2)$

---

426 1: Suppose we use  $O(nd)$  space to store  $A_1$  and  $A_2$  (Avoid  $n^2$  time/space)  
 427 2: Suppose we receive an update  $q \in [2]$ ,  $i \in [n]$ ,  $j \in [d]$ ,  $\Delta$   
 428 3: Suppose we have hash function  $g$  to access uniform number  
 429 4: **if**  $q = 1$  **then**  
 430 5:    $p \leftarrow g(i(n-1) + 1, \dots, in) \{p \in \mathbb{R}^n\}$   
 431 6:    $y \leftarrow y + \Phi \Delta(e_{[i(n-1)+1,in]} \circ (A_2)_{*,j})/p \{\Phi_1 \text{ is decided by } h_1, \sigma_1\}$   
 7: **else**  
 8:    $y_2 \leftarrow y_2 + \Phi_2 e_i \Delta \{\Phi_2 \text{ is decided by } h_2, \sigma_2\}$   
 9: **end if**

---

In the following Lemma, we state a streaming algorithm to solve tensor related sampling problem. We consider the situation that one of  $A_1$  and  $A_2$  is fixed, and the other one is updated in streaming fashion. We show the following estimation guarantees using the standard CountSketch analysis, c.f., Charikar et al. (2004); Jowhari et al. (2011). We defer the proof to appendix C.

**Lemma 7.4** (Tensor  $\ell_2$  Tail Estimation). *Let  $y = (A_1 \otimes A_2)x \in \mathbb{R}^{n^2}$ . Let only one of  $A_1$  and  $A_2$  be updated in streaming. Let  $w = \frac{y_i}{\sqrt{u_i}}$  for a constant  $u_i \in [0, 1]$  generated uniformly at random. There is an algorithm  $\mathcal{A}$  that that uses  $O(nd) + \text{poly}(\frac{1}{\epsilon}, \log n)$  space, uses  $O(n)$  update time, and estimates each element of  $w$  up to additive error  $\epsilon \cdot \|z\|_2$ , where  $z$  denotes the tail vector of  $w$  without the largest  $\frac{1}{\epsilon^2}$  entries in magnitude. Specifically, for all  $i \in [n^2]$ , we have  $|\widehat{w}_i - w_i| \leq \epsilon \cdot \|z\|_2$ .*

We state the following lemma as a structural property that will allow us to achieve our tensor product sampler. We remark that the proof is a simple adaptation of existing proofs for approximate  $\ell_p$  sampling (Jowhari et al., 2011). Thus we defer the proof to Appendix C.

**Lemma 7.5.** *Let  $y = (A_1 \otimes A_2)x \in \mathbb{R}^{n^2}$  and let  $w \in \mathbb{R}^{n^2}$  so that  $w_i = \frac{y_i}{\sqrt{u_i}}$  for a constant  $u_i \in [0, 1]$  generated uniformly at random. Let  $z$  denote the tail vector of  $w$  without the largest  $\frac{1}{\epsilon^2}$  entries in magnitude. Let  $\widehat{Z}$  be a 2-approximation to  $\|z\|_2$  and  $\widehat{Y}$  be a 2-approximation to  $\|y\|_2$ , then we have  $\Pr[\widehat{Z} > \sqrt{(C \log n)/\epsilon} \cdot \widehat{Y}] \leq O(\epsilon) + \frac{1}{\text{poly}(n)}$ .*

Finally, we describe the guarantees of our tensor-based sampler, deferring the proof to Appendix D.

**Theorem 7.6.** *Let  $y = (A_1 \otimes A_2)x \in \mathbb{R}^{n^2}$  and let  $w \in \mathbb{R}^{n^2}$  so that for each  $i \in [n^2]$ ,  $w_i = \frac{y_i}{\sqrt{u_i}}$  for a constant  $u_i \in [0, 1]$  generated uniformly at random. Let  $z$  denote the tail vector of  $w$  without the largest  $\frac{1}{\epsilon^2}$  entries in magnitude. Suppose there exists:*

1. *An algorithm  $\mathcal{A}_1$  that provides a 2-approximation to  $\|y\|_2$  with probability  $1 - \frac{1}{n^2}$ .*
2. *An algorithm  $\mathcal{A}_2$  that provides a 2-approximation to  $\|z\|_2$  with probability  $1 - \frac{1}{n^2}$ .*
3. *An algorithm  $\mathcal{A}_3$  that estimates each element of  $w$  up to additive error  $\epsilon \cdot \|z\|_2$ ,  $|\widehat{w}_i - w_i| \leq \epsilon \cdot \|z\|_2$ , for all  $i \in [n^2]$ .*

*Then there exists a data structure that uses  $\text{poly}(\frac{1}{\epsilon}, \log n)$  bits of space and outputs each index  $i$  with probability  $p_i = (1 \pm \epsilon) \cdot \frac{y_i^2}{\|y\|_2^2} \pm \frac{1}{\text{poly}(n)}$ .*

We remark that the algorithms  $\mathcal{A}_1$  and  $\mathcal{A}_2$  in the context of Theorem 7.6 can be achieved using the standard AMS  $\ell_2$  norm estimator (Alon et al., 1999). Moreover, algorithm  $\mathcal{A}_3$  in the context of Theorem 7.6 can be achieved using the standard CountSketch algorithm (Charikar et al., 2004).

## 8 CONCLUSIONS

To achieve efficient attention mechanisms, we introduce the attention sampler and study its behavior in the streaming model. We established efficient polynomial samplers under various streaming settings, when the input matrix, the weight vector, or both evolve dynamically, and we complement the results by proving space lower bounds. Our framework identify the critical components in attention computation, offering a foundation for efficient simulations of large-scale attention schemes, which is central to modern machine learning and LLMs.

For future directions, from a theoretical perspective, given the  $\Omega(n)$  lower bound on exponential samplers in general circumstances, it would be valuable to explore whether we can achieve  $o(n)$  space under certain assumptions, e.g., restricting the entries in the attention matrix to  $o(\log n)$ . From a practical perspective, it would be beneficial to evaluate our sampler’s performance by implementing it in existing sparse attention schemes and streaming attention schemes.

486 ETHIC STATEMENT  
487488 This paper does not involve human subjects, personally identifiable data, or sensitive applications.  
489 We do not foresee direct ethical risks. We follow the ICLR Code of Ethics and affirm that all aspects  
490 of this research comply with the principles of fairness, transparency, and integrity.  
491492 REPRODUCIBILITY STATEMENT  
493494 We ensure reproducibility of our theoretical results by including all formal assumptions, definitions,  
495 and complete proofs in the appendix. The main text states each theorem clearly and refers to the  
496 detailed proofs. No external data or software is required.  
497498 REFERENCES  
499

500 Marah I Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany  
501 Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat S. Behl, Alon Benhaim, Misha  
502 Bilenko, Johan Bjorck, Sébastien Bubeck, Martin Cai, Caio César Teodoro Mendes, Weizhu  
503 Chen, Vishrav Chaudhary, Parul Chopra, Allie Del Giorno, Gustavo de Rosa, Matthew Dixon,  
504 Ronen Eldan, Dan Iter, Amit Garg, Abhishek Goswami, Suriya Gunasekar, Emmann Haider,  
505 Junheng Hao, Russell J. Hewett, Jamie Huynh, Mojan Javaheripi, Xin Jin, Piero Kauffmann, Nikos  
506 Karampatziakis, Dongwoo Kim, Mahoud Khademi, Lev Kurilenko, James R. Lee, Yin Tat Lee,  
507 Yuanzhi Li, Chen Liang, Weishung Liu, Eric Lin, Zeqi Lin, Piyush Madan, Arindam Mitra, Hardik  
508 Modi, Anh Nguyen, Brandon Norick, Barun Patra, Daniel Perez-Becker, Thomas Portet, Reid  
509 Pryzant, Heyang Qin, Marko Radmilac, Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli  
510 Saarikivi, Amin Saied, Adil Salim, Michael Santacroce, Shital Shah, Ning Shang, Hiteshi Sharma,  
511 Xia Song, Masahiro Tanaka, Xin Wang, Rachel Ward, Guanhua Wang, Philipp Witte, Michael  
512 Wyatt, Can Xu, Jiahang Xu, Sonali Yadav, Fan Yang, Ziyi Yang, Donghan Yu, Chengruidong  
513 Zhang, Cyril Zhang, Jianwen Zhang, Li Lyra Zhang, Yi Zhang, Yue Zhang, Yunan Zhang, and  
514 Xiren Zhou. Phi-3 technical report: A highly capable language model locally on your phone.  
515 *CoRR*, abs/2404.14219, 2024. doi: 10.48550/ARXIV.2404.14219. URL <https://doi.org/10.48550/arXiv.2404.14219>.

516 Ankit Aggarwal, Amit Deshpande, and Ravi Kannan. Adaptive sampling for k-means clustering.  
517 In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques,*  
518 *12th International Workshop, APPROX and 13th International Workshop, RANDOM. Proceedings*,  
519 pp. 15–28, 2009.

520 Thomas D Ahle, Michael Kapralov, Jakob BT Knudsen, Rasmus Pagh, Ameya Velingker, David P  
521 Woodruff, and Amir Zandieh. Oblivious sketching of high-degree polynomial kernels. In *Pro-*  
522 *ceedings of the Fourteenth Annual ACM-SIAM Symposium on Discrete Algorithms*, pp. 141–160.  
523 SIAM, 2020.

525 Elizabeth S. Allman, Catherine Matias, and John A. Rhodes. Identifiability of parameters in latent  
526 structure models with many observed variables. *The Annals of Statistics*, 37(6A), dec 2009a. doi:  
527 10.1214/09-aos689. URL <https://doi.org/10.1214%2F09-aos689>.

528 Elizabeth S. Allman, Sonja Petrović, John A. Rhodes, and Seth Sullivant. Identifiability of 2-tree  
529 mixtures for group-based models, 2009b. URL <https://arxiv.org/abs/0909.1854>.

531 Josh Alman and Zhao Song. Fast attention requires bounded entries. In *Advances in Neural*  
532 *Information Processing Systems 36: Annual Conference on Neural Information Processing Systems,*  
533 *NeurIPS*, 2023.

534 Josh Alman and Zhao Song. How to capture higher-order correlations? generalizing matrix soft-  
535 max attention to kronecker computation. In *The Twelfth International Conference on Learning*  
536 *Representations, ICLR*, 2024a.

538 Josh Alman and Zhao Song. The fine-grained complexity of gradient computation for training large  
539 language models. In *The Thirty-eighth Annual Conference on Neural Information Processing*  
540 *Systems*, 2024b. URL <https://openreview.net/forum?id=up4tWnwR0l>.

540 Josh Alman and Zhao Song. Fast rope attention: Combining the polynomial method and fast fourier  
 541 transform. *arXiv preprint arXiv:2505.11892*, 2025a.  
 542

543 Josh Alman and Zhao Song. Only large weights (and not skip connections) can prevent the perils of  
 544 rank collapse. *arXiv preprint arXiv:2505.16284*, 2025b.  
 545

546 Noga Alon, Yossi Matias, and Mario Szegedy. The space complexity of approximating the frequency  
 547 moments. *J. Comput. Syst. Sci.*, 58(1):137–147, 1999.  
 548

549 Alexandr Andoni, Robert Krauthgamer, and Krzysztof Onak. Streaming algorithms via precision  
 550 sampling. In *IEEE 52nd Annual Symposium on Foundations of Computer Science, FOCS*, pp.  
 551 363–372, 2011.  
 552

553 Pranjal Awasthi and Anupam Gupta. Improving length-generalization in transformers via task hinting.  
 554 *arXiv preprint arXiv:2310.00726*, 2023.  
 555

556 Ziv Bar-Yossef, T. S. Jayram, Ravi Kumar, and D. Sivakumar. An information statistics approach to  
 557 data stream and communication complexity. *J. Comput. Syst. Sci.*, 68(4):702–732, 2004.  
 558

559 Srinadh Bhojanapalli and Sujay Sanghavi. A new sampling technique for tensors. In *arXiv preprint*.  
 560 [https://arxiv.org/pdf/1502.05023](https://arxiv.org/pdf/1502.05023.pdf), 2015.  
 561

562 Jan van den Brand, Zhao Song, and Tianyi Zhou. Algorithm and hardness for dynamic attention  
 563 maintenance in large language models. In *Forty-first International Conference on Machine  
 564 Learning, ICML*. OpenReview.net, 2024.  
 565

566 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
 567 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
 568 few-shot learners. *Advances in Neural Information Processing Systems 33: Annual Conference on  
 569 Neural Information Processing Systems, NeurIPS*, pp. 1877–1901, 2020.  
 570

571 Moses Charikar, Kevin C. Chen, and Martin Farach-Colton. Finding frequent items in data streams.  
 572 *Theor. Comput. Sci.*, 312(1):3–15, 2004.  
 573

574 Beidi Chen, Tri Dao, Eric Winsor, Zhao Song, Atri Rudra, and Christopher Ré. Scatterbrain: Unifying  
 575 sparse and low-rank attention. *Advances in Neural Information Processing Systems (NeurIPS)*, 34:  
 576 17413–17426, 2021.  
 577

578 Beidi Chen, Tri Dao, Kaizhao Liang, Jiaming Yang, Zhao Song, Atri Rudra, and Christopher Ré.  
 579 Pixelated butterfly: Simple and efficient sparse training for neural network models. In *International  
 580 Conference on Learning Representations, ICLR*, 2022.  
 581

582 Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse  
 583 transformers. *arXiv preprint arXiv:1904.10509*, 2019.  
 584

585 Joon Hee Choi and S. Vishwanathan. Dfacto: Distributed factorization of tensors. In Zoubin  
 586 Ghahramani, Max Welling, Corinna Cortes, Neil D. Lawrence, and Kilian Q. Weinberger (eds.),  
 587 *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information  
 588 Processing Systems, NeurIPS*, pp. 1296–1304, 2014.  
 589

590 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam  
 591 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh,  
 592 Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam  
 593 Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James  
 594 Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Lev-  
 595 skaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin  
 596 Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph,  
 597 Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M.  
 598 Dai, Thanumalayan Sankaranarayana Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon  
 599 Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark  
 600 Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean,  
 601 Slav Petrov, and Noah Fiedel. Palm: Scaling language modeling with pathways. *J. Mach. Learn.  
 602 Res.*, 24:240:1–240:113, 2023.  
 603

594     Edith Cohen and Ofir Geri. Sampling sketches for concave sublinear functions of frequencies. In  
 595       *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information*  
 596       *Processing Systems, NeurIPS*, pp. 1361–1371, 2019.  
 597

598     Edith Cohen, Nick G. Duffield, Haim Kaplan, Carsten Lund, and Mikkel Thorup. Efficient stream  
 599       sampling for variance-optimal estimation of subset sums. *SIAM J. Comput.*, 40(5):1402–1431,  
 600       2011.

601     Edith Cohen, Nick G. Duffield, Haim Kaplan, Carsten Lund, and Mikkel Thorup. Algorithms and  
 602       estimators for summarization of unaggregated data streams. *J. Comput. Syst. Sci.*, 80(7):1214–1244,  
 603       2014.

604

605     Michael B. Cohen. Nearly tight oblivious subspace embeddings by trace inequalities. In *Proceedings*  
 606       *of the Twenty-Seventh Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, Arlington,  
 607       VA, USA, January 10-12, 2016, pp. 278–287, 2016.

608

609     Michael B Cohen, Yin Tat Lee, and Zhao Song. Solving linear programs in the current matrix  
 610       multiplication time. In *Proceedings of the 51st Annual ACM Symposium on Theory of Computing*  
 611       (*STOC*). <https://arxiv.org/pdf/1810.07896.pdf>, 2019.

612     Graham Cormode and Hossein Jowhari.  $l_p$  samplers and their applications: A survey. *ACM Comput.*  
 613       *Surv.*, 52(1):16:1–16:31, 2019.

614

615     Graham Cormode, S. Muthukrishnan, and Irina Rozenbaum. Summarizing and mining inverse  
 616       distributions on data streams via dynamic inverse sampling. In *Proceedings of the 31st International*  
 617       *Conference on Very Large Data Bases*, pp. 25–36. ACM, 2005.

618

619     Giannis Daras, Nikita Kitaev, Augustus Odena, and Alexandros G Dimakis. Smyrf-efficient attention  
 620       using asymmetric clustering. *Advances in Neural Information Processing Systems (NeurIPS)*, 33:  
 621       6476–6489, 2020.

622

623     Mostafa Dehghani, Stephan Gouws, Oriol Vinyals, Jakob Uszkoreit, and Łukasz Kaiser. Universal  
 624       transformers. *arXiv preprint arXiv:1807.03819*, 2018.

625

626     Yichuan Deng, Wenyu Jin, Zhao Song, Xiaorui Sun, and Omri Weinstein. Dynamic kernel sparsifiers.  
 627       *arXiv preprint arXiv:2211.14825*, 2022a.

628

629     Yichuan Deng, Zhao Song, Omri Weinstein, and Ruizhe Zhang. Fast distance oracles for any  
 630       symmetric norm. In *Advances in Neural Information Processing Systems 35: Annual Conference*  
 631       *on Neural Information Processing Systems 2022, NeurIPS*, 2022b.

632

633     Yichuan Deng, Yeqi Gao, and Zhao Song. Solving tensor low cycle rank approximation. In *IEEE*  
 634       *International Conference on Big Data, BigData 2023, Sorrento, Italy, December 15-18, 2023*, pp.  
 635       6–16. IEEE, 2023a.

636

637     Yichuan Deng, Zhihang Li, and Zhao Song. Attention scheme inspired softmax regression. *arXiv*  
 638       *preprint arXiv:2304.10411*, 2023b.

639

640     Yichuan Deng, Sridhar Mahadevan, and Zhao Song. Randomized and deterministic attention sparsifi-  
 641       cation algorithms for over-parameterized feature dimension. *arxiv preprint: arxiv 2304.03426*,  
 642       2023c.

643

644     Yichuan Deng, Zhao Song, and Shenghao Xie. Convergence of two-layer regression with nonlinear  
 645       units. *arXiv preprint arXiv:2308.08358*, 2023d.

646

647     Yichuan Deng, Zhihang Li, Sridhar Mahadevan, and Zhao Song. Zero-th order algorithm for softmax  
 648       attention optimization. In *IEEE International Conference on Big Data, BigData 2024, Washington,*  
 649       *DC, USA, December 15-18, 2024*, pp. 24–33. IEEE, 2024a.

650

651     Yichuan Deng, Zhao Song, and Chiwun Yang. Attention is naturally sparse with gaussian distributed  
 652       input. *arXiv preprint arXiv:2404.02690*, 2024b.

648 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep  
 649 bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of*  
 650 *the North American Chapter of the Association for Computational Linguistics: Human Language*  
 651 *Technologies, NAACL-HLT*, pp. 4171–4186. Association for Computational Linguistics, 2019.

652 Huaian Diao, Zhao Song, Wen Sun, and David P. Woodruff. Sketching for kronecker product  
 653 regression and p-splines. In *International Conference on Artificial Intelligence and Statistics,*  
 654 *AISTATS 2018, 9-11 April 2018, Playa Blanca, Lanzarote, Canary Islands, Spain*, volume 84 of  
 655 *Proceedings of Machine Learning Research*, pp. 1299–1308. PMLR, 2018.

656 Huaian Diao, Rajesh Jayaram, Zhao Song, Wen Sun, and David Woodruff. Optimal sketching  
 657 for kronecker product regression and low rank approximation. *Advances in neural information*  
 658 *processing systems*, 32, 2019.

659 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha  
 660 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn,  
 661 Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston  
 662 Zhang, Aurélien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Rozière, Bethany Biron,  
 663 Bin Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris  
 664 McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton  
 665 Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David  
 666 Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes,  
 667 Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip  
 668 Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme  
 669 Nail, Grégoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu,  
 670 Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel M. Kloemann, Ishan Misra, Ivan  
 671 Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet  
 672 Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi,  
 673 Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph  
 674 Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani,  
 675 Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, and et al. The llama 3 herd of models. *CoRR*,  
 676 abs/2407.21783, 2024. URL <https://doi.org/10.48550/arXiv.2407.21783>.

677 Benjamin L. Edelman, Surbhi Goel, Sham M. Kakade, and Cyril Zhang. Inductive biases and variable  
 678 creation in self-attention mechanisms. In *International Conference on Machine Learning, ICML*,  
 679 volume 162 of *Proceedings of Machine Learning Research*, pp. 5793–5831. PMLR, 2022.

680 Alan M. Frieze, Ravi Kannan, and Santosh S. Vempala. Fast monte-carlo algorithms for finding  
 681 low-rank approximations. *J. ACM*, 51(6):1025–1041, 2004.

682 Yao Fu. Challenges in deploying long-context transformers: A theoretical peak performance analysis.  
 683 *CoRR*, abs/2405.08944, 2024. doi: 10.48550/ARXIV.2405.08944. URL <https://doi.org/10.48550/arXiv.2405.08944>.

684 Yeqi Gao, Lianke Qin, Zhao Song, and Yitan Wang. A sublinear adversarial training algorithm. *arXiv*  
 685 *preprint arXiv:2208.05395*, 2022.

686 Yeqi Gao, Sridhar Mahadevan, and Zhao Song. An over-parameterized exponential regression. *arXiv*  
 687 *preprint arXiv:2303.16504*, 2023.

688 Yeqi Gao, Zhao Song, and Junze Yin. An iterative algorithm for rescaled hyperbolic functions  
 689 regression. In *International Conference on Artificial Intelligence and Statistics, AISTATS*, volume  
 690 258 of *Proceedings of Machine Learning Research*, pp. 2548–2556. PMLR, 2025.

691 Rainer Gemulla, Wolfgang Lehner, and Peter J. Haas. Maintaining bounded-size sample synopses of  
 692 evolving datasets. *VLDB J.*, 17(2):173–202, 2008.

693 Yuzhou Gu and Zhao Song. A faster small treewidth sdp solver. *arXiv preprint arXiv:2211.06033*,  
 694 2022.

695 Yuzhou Gu, Zhao Song, Junze Yin, and Lichen Zhang. Low rank matrix completion via robust  
 696 alternating minimization in nearly linear time. In *The Twelfth International Conference on Learning*  
 697 *Representations, ICLR*. OpenReview.net, 2024.

702 Peter J. Haas. Data-stream sampling: Basic techniques and results. In *Data Stream Management -*  
 703 *Processing High-Speed Data Streams*, Data-Centric Systems and Applications, pp. 13–44. Springer,  
 704 2016.

705 Peter J. Haas and Arun N. Swami. Sequential sampling procedures for query size estimation. In  
 706 *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pp. 341–350,  
 707 1992.

708 Insu Han, Rajesh Jayaram, Amin Karbasi, Vahab Mirrokni, David P. Woodruff, and Amir Zandieh.  
 709 Hyperattention: Long-context attention in near-linear time. In *The Twelfth International Conference*  
 710 *on Learning Representations, ICLR*. OpenReview.net, 2024.

711 Richard A Harshman. Foundations of the parafac procedure: Models and conditions for an "explanatory"  
 712 multimodal factor analysis. 1970.

713 Furong Huang, Niranjan U. N, Mohammad Umar Hakeem, Prateek Verma, and Animashree Anand-  
 714 kumar. Fast detection of overlapping communities via online tensor methods on gpus. *CoRR*,  
 715 abs/1309.0787, 2013.

716 Ling Huang, XuanLong Nguyen, Minos N. Garofalakis, Joseph M. Hellerstein, Michael I. Jordan,  
 717 Anthony D. Joseph, and Nina Taft. Communication-efficient online detection of network-wide  
 718 anomalies. In *INFOCOM. 26th IEEE International Conference on Computer Communications, Joint Conference of the IEEE Computer and Communications Societies*, pp. 134–142, 2007.

719 Adobe Inc. Adobe firefly. <https://www.adobe.com/sensei/generative-ai/firefly.html>, 2023.

720 Piotr Indyk, Sepideh Mahabadi, Shayan Oveis Gharan, and Alireza Rezaei. Composable core-sets for  
 721 determinant maximization problems via spectral spanners. In *Proceedings of the 2020 ACM-SIAM*  
 722 *Symposium on Discrete Algorithms, SODA*, pp. 1675–1694, 2020.

723 Rajesh Jayaram and David P. Woodruff. Perfect  $l_p$  sampling in a data stream. *SIAM J. Comput.*, 50  
 724 (2):382–439, 2021.

725 Rajesh Jayaram, David P. Woodruff, and Samson Zhou. Truly perfect samplers for data streams and  
 726 sliding windows. In *PODS '22: International Conference on Management of Data*, pp. 29–40,  
 727 2022.

728 Shunhua Jiang, Zhao Song, Omri Weinstein, and Hengjie Zhang. Faster dynamic matrix inverse for  
 729 faster lps. *arXiv preprint arXiv:2004.07470*, 2021.

730 Hossein Jowhari, Mert Saglam, and Gábor Tardos. Tight bounds for  $l_p$  samplers, finding duplicates  
 731 in streams, and related problems. In *Proceedings of the 30th ACM SIGMOD-SIGACT-SIGART*  
 732 *Symposium on Principles of Database Systems, PODS*, pp. 49–58, 2011.

733 Praneeth Kacham, Vahab Mirrokni, and Peilin Zhong. Polysketchformer: Fast transformers via  
 734 sketches for polynomial kernels. *arXiv preprint arXiv:2310.01655*, 2023.

735 Bala Kalyanasundaram and Georg Schnitger. The probabilistic communication complexity of set  
 736 intersection. *SIAM J. Discret. Math.*, 5(4):545–557, 1992.

737 U. Kang, Evangelos E. Papalexakis, Abhay Harpale, and Christos Faloutsos. Gigatensor: scaling  
 738 tensor analysis up by 100 times - algorithms and discoveries. In *KDD*, pp. 316–324, 2012.

739 Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are rnns:  
 740 Fast autoregressive transformers with linear attention. In *International conference on machine*  
 741 *learning*, pp. 5156–5165. PMLR, 2020.

742 Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. *arXiv*  
 743 *preprint arXiv: 2001.04451*, 2020.

744 Soroush Abbasi Koohpayegani and Hamed Pirsiavash. Sima: Simple softmax-free attention for  
 745 vision transformers. In *IEEE/CVF Winter Conference on Applications of Computer Vision, WACV*,  
 746 pp. 2595–2605. IEEE, 2024.

756 Ilan Kremer, Noam Nisan, and Dana Ron. On randomized one-round communication complexity.  
 757 *Comput. Complex.*, 8(1):21–49, 1999.  
 758

759 Joseph B. Kruskal. Three-way arrays: rank and uniqueness of trilinear decompositions, with  
 760 application to arithmetic complexity and statistics. *Linear Algebra and its Applications*, 18(2):95–  
 761 138, 1977. ISSN 0024-3795. doi: [https://doi.org/10.1016/0024-3795\(77\)90069-6](https://doi.org/10.1016/0024-3795(77)90069-6). URL <https://www.sciencedirect.com/science/article/pii/0024379577900696>.

763 Xunhao Lai, Jianqiao Lu, Yao Luo, Yiyuan Ma, and Xun Zhou. Flexprefill: A context-aware  
 764 sparse attention mechanism for efficient long-sequence inference. In *The Thirteenth International  
 765 Conference on Learning Representations*. OpenReview.net, 2025.

766

767 Yin Tat Lee, Zhao Song, and Qiuyi Zhang. Solving empirical risk minimization in the current matrix  
 768 multiplication time. In *COLT*. <https://arxiv.org/pdf/1905.04447.pdf>, 2019.

769

770 Zhihang Li, Zhao Song, and Tianyi Zhou. Solving regularized exp, cosh and sinh regression problems.  
 771 *arXiv preprint arXiv:2303.15725*, 2023.

772

773 Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike  
 774 Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining  
 775 approach. *arXiv preprint arXiv:1907.11692*, 2019.

776

777 Sepideh Mahabadi, Piotr Indyk, Shayan Oveis Gharan, and Alireza Rezaei. Composable core-sets  
 778 for determinant maximization: A simple near-optimal algorithm. In *Proceedings of the 36th  
 779 International Conference on Machine Learning, ICML 2019*, volume 97, pp. 4254–4263, 2019.

780

781 Sepideh Mahabadi, Ilya P. Razenshteyn, David P. Woodruff, and Samson Zhou. Non-adaptive adaptive  
 782 sampling on turnstile streams. In *Proceedings of the 52nd Annual ACM SIGACT Symposium on  
 783 Theory of Computing, STOC*, pp. 1251–1264, 2020.

784

785 Sepideh Mahabadi, David P. Woodruff, and Samson Zhou. Adaptive sketches for robust regression  
 786 with importance sampling. In *Approximation, Randomization, and Combinatorial Optimization.  
 787 Algorithms and Techniques, APPROX/RANDOM*, pp. 31:1–31:21, 2022.

788

789 Arvind V. Mahankali, David P. Woodruff, and Ziyu Zhang. Near-linear time and fixed-parameter  
 790 tractable algorithms for tensor decompositions. In *15th Innovations in Theoretical Computer  
 791 Science Conference, ITCS*, volume 287 of *LIPICS*, pp. 79:1–79:23. Schloss Dagstuhl - Leibniz-  
 792 Zentrum für Informatik, 2024.

793

794 Jianning Mai, Chen-Nee Chuah, Ashwin Sridharan, Tao Ye, and Hui Zang. Is sampled data sufficient  
 795 for anomaly detection? In *Proceedings of the 6th ACM SIGCOMM Internet Measurement  
 796 Conference, IMC*, pp. 165–176, 2006.

797

798 James Manyika. An overview of bard: an early experiment with generative ai. Technical report, Tech.  
 799 rep., Technical report, Google AI, 2023.

800

801 Gary Marcus, Ernest Davis, and Scott Aaronson. A very preliminary analysis of dall-e 2. *arXiv  
 802 preprint arXiv:2204.13807*, 2022.

803

804 Morteza Monemizadeh and David P. Woodruff. 1-pass relative-error  $l_p$ -sampling with applications.  
 805 In *Proceedings of the Twenty-First Annual ACM-SIAM Symposium on Discrete Algorithms, SODA*,  
 806 pp. 1143–1160. SIAM, 2010.

807

808 Jelani Nelson and Huy L Nguyn. Osnap: Faster numerical linear algebra algorithms via sparser  
 809 subspace embeddings. In *2013 ieee 54th annual symposium on foundations of computer science*,  
 810 pp. 117–126. IEEE, 2013.

811

812 Noam Nisan. Pseudorandom generators for space-bounded computation. *Comb.*, 12(4):449–461,  
 813 1992.

814

815 Rasmus Pagh. Compressed matrix multiplication. *ACM Transactions on Computation Theory  
 816 (TOCT)*, 5(3):1–17, 2013.

810 Anh-Huy Phan, Petr Tichavsky, and Andrzej Cichocki. Low complexity damped gauss–newton  
 811 algorithms for candecomp/parafac. *SIAM Journal on Matrix Analysis and Applications*, 34(1):  
 812 126–147, 2013.

813

814 Lianke Qin, Rajesh Jayaram, Elaine Shi, Zhao Song, Danyang Zhuo, and Shumo Chu. Adore:  
 815 Differentially oblivious relational database operators. *VLDB*, 2022a.

816

817 Lianke Qin, Aravind Reddy, Zhao Song, Zhaozhuo Xu, and Danyang Zhuo. Adaptive and dynamic  
 818 multi-resolution hashing for pairwise summations. In *IEEE International Conference on Big Data,  
 819 Big Data*, pp. 115–120. IEEE, 2022b.

820 Lianke Qin, Zhao Song, and Yitan Wang. Fast submodular function maximization. *arXiv preprint  
 821 arXiv:2305.08367*, 2023a.

822

823 Lianke Qin, Zhao Song, Lichen Zhang, and Danyang Zhuo. An online and unified algorithm  
 824 for projection matrix vector multiplication with application to empirical risk minimization. In  
 825 *International Conference on Artificial Intelligence and Statistics*, pp. 101–156. PMLR, 2023b.

826 Lianke Qin, Zhao Song, and Ruizhe Zhang. A general algorithm for solving rank-one matrix sensing.  
 827 In *International Conference on Artificial Intelligence and Statistics*, volume 238 of *Proceedings of  
 828 Machine Learning Research*, pp. 757–765. PMLR, 2024.

829

830 Alexander A. Razborov. On the distributional complexity of disjointness. *Theor. Comput. Sci.*, 106  
 831 (2):385–390, 1992.

832

833 Aravind Reddy, Zhao Song, and Lichen Zhang. Dynamic tensor product regression. In *Advances in  
 834 Neural Information Processing Systems 35: Annual Conference on Neural Information Processing  
 835 Systems 2022, NeurIPS*, 2022.

836 Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy P. Lillicrap, Jean-  
 837 Baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittweiser, Ioannis  
 838 Antonoglou, Rohan Anil, Sebastian Borgeaud, Andrew M. Dai, Katie Millican, Ethan Dyer,  
 839 Mia Glaese, Thibault Sottiaux, Benjamin Lee, Fabio Viola, Malcolm Reynolds, Yuanzhong Xu,  
 840 James Molloy, Jilin Chen, Michael Isard, Paul Barham, Tom Hennigan, Ross McIlroy, Melvin  
 841 Johnson, Johan Schalkwyk, Eli Collins, Eliza Rutherford, Erica Moreira, Kareem Ayoub, Megha  
 842 Goel, Clemens Meyer, Gregory Thornton, Zhen Yang, Henryk Michalewski, Zaheer Abbas,  
 843 Nathan Schucher, Ankesh Anand, Richard Ives, James Keeling, Karel Lenc, Salem Haykal,  
 844 Siamak Shakeri, Pranav Shyam, Aakanksha Chowdhery, Roman Ring, Stephen Spencer, Eren  
 845 Sezener, and et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens  
 846 of context. *CoRR*, abs/2403.05530, 2024. doi: 10.48550/ARXIV.2403.05530. URL <https://doi.org/10.48550/arXiv.2403.05530>.

847

848 John A Rhodes and Seth Sullivant. Identifiability of large phylogenetic mixture models. *Bulletin of  
 849 mathematical biology*, 74:212–231, 2012.

850

851 Elina Robeva. Orthogonal decomposition of symmetric tensors. *SIAM Journal on Matrix Analysis  
 852 and Applications*, 37(1):86–102, 2016.

853

854 Elina Robeva and Anna Seigal. Singular vectors of orthogonally decomposable tensors. *Linear and  
 855 Multilinear Algebra*, 65(12):2457–2471, 2017.

856

857 Hemanth Saratchandran, Jianqiao Zheng, Yiping Ji, Wenbo Zhang, and Simon Lucey. Rethinking  
 softmax: Self-attention with polynomial activations. *CoRR*, abs/2410.18613, 2024.

858

859 Sambal Shikhar, Mohammed Irfan Kurpath, Sahal Shaji Mullappilly, Jean Lahoud, Fahad Shahbaz  
 860 Khan, Rao Muhammad Anwer, Salman H. Khan, and Hisham Cholakkal. Llmvox: Autoregressive  
 861 streaming text-to-speech model for any LLM. In *Findings of the Association for Computational  
 862 Linguistics, ACL*, pp. 20481–20493. Association for Computational Linguistics, 2025.

863

Charlie Snell, Ruiqi Zhong, Dan Klein, and Jacob Steinhardt. Approximating how single head  
 attention learns. *arXiv preprint arXiv:2103.07601*, 2021.

864 Zhao Song, David P. Woodruff, and Huan Zhang. Sublinear time orthogonal tensor decomposition. In  
 865 *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information*  
 866 *Processing Systems (NIPS) 2016, December 5-10, 2016, Barcelona, Spain*, pp. 793–801, 2016.  
 867

868 Zhao Song, David P. Woodruff, and Peilin Zhong. Relative error tensor low rank approximation. In  
 869 *Proceedings of the Thirtieth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA*, pp.  
 870 2772–2789. SIAM, 2019.

871 Zhao Song, David Woodruff, Zheng Yu, and Lichen Zhang. Fast sketching of polynomial kernels of  
 872 polynomial degree. In *International Conference on Machine Learning*, pp. 9812–9823. PMLR,  
 873 2021.

874 Zhao Song, Zhaozhuo Xu, and Lichen Zhang. Speeding up sparsification using inner product search  
 875 data structures. *arXiv preprint arXiv:2204.03209*, 2022.

876 Zhao Song, Xin Yang, Yuanyuan Yang, and Lichen Zhang. Sketching meets differential privacy: Fast  
 877 algorithm for dynamic kronecker projection maintenance. In *International Conference on Machine*  
 878 *Learning, ICML*, volume 202 of *Proceedings of Machine Learning Research*, pp. 32418–32462.  
 879 PMLR, 2023a.

880 Zhao Song, Mingquan Ye, and Lichen Zhang. Streaming semidefinite programs:  $o(\sqrt{n})$  passes, small  
 881 space and fast runtime. *Manuscript*, 2023b.

882 Zhao Song, Lichen Zhang, and Ruizhe Zhang. Training multi-layer over-parametrized neural network  
 883 in subquadratic time. In *15th Innovations in Theoretical Computer Science Conference, ITCS*,  
 884 volume 287 of *LIPICS*, pp. 93:1–93:15. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2024.

885 Foteini Strati, Sara McAllister, Amar Phanishayee, Jakub Tarnawski, and Ana Klimovic. Déjàvu:  
 886 Kv-cache streaming for fast, fault-tolerant generative LLM serving. In *Forty-first International*  
 887 *Conference on Machine Learning, ICML*. OpenReview.net, 2024.

888 William Swartworth, David P. Woodruff, and Samson Zhou. Perfect lp sampling with polylogarithmic  
 889 update time. In *IEEE 66th Annual Symposium on Foundations of Computer Science (FOCS)*, 2025.

890 Marina Thottan, Guanglei Liu, and Chuanyi Ji. Anomaly detection approaches for communication  
 891 networks. In *Algorithms for Next Generation Networks*, Computer Communications and Networks,  
 892 pp. 239–261. Springer, 2010.

893 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
 894 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and  
 895 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.

896 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay  
 897 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation  
 898 and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.

899 Charalampos E. Tsourakakis. MACH: fast randomized tensor decompositions. In *SDM*, pp. 689–700,  
 900 2010.

901 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz  
 902 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing*  
 903 *systems*, 30, 2017.

904 Jeffrey Scott Vitter. Random sampling with a reservoir. *ACM Trans. Math. Softw.*, 11(1):37–57, 1985.

905 James Vuckovic, Aristide Baratin, and Remi Tachet des Combes. A mathematical theory of attention.  
 906 *arXiv preprint arXiv:2007.02876*, 2020.

907 Chi Wang, Xueqing Liu, Yanglei Song, and Jiawei Han. Scalable moment-based inference for latent  
 908 dirichlet allocation. In *ECML-PKDD*, pp. 290–305, 2014.

909 Sinong Wang, Belinda Z Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-attention with  
 910 linear complexity. *arXiv preprint arXiv:2006.04768*, 2020.

918 Yining Wang, Hsiao-Yu Tung, Alexander J Smola, and Anima Anandkumar. Fast and guaranteed  
 919 tensor decomposition via sketching. In *Advances in Neural Information Processing Systems (NIPS)*,  
 920 pp. 991–999. <https://arxiv.org/pdf/1506.04448.pdf>, 2015.

921  
 922 Colin Wei, Yining Chen, and Tengyu Ma. Statistically meaningful approximation: a case study on  
 923 approximating turing machines with transformers. *arXiv preprint arXiv:2107.13163*, 2021.

924 David P. Woodruff and Samson Zhou. Tight bounds for adversarially robust streams and sliding  
 925 windows via difference estimators. In *62nd Annual Symposium on Foundations of Computer  
 926 Science, FOCS*, pp. 1183–1196, 2021.

927 David P. Woodruff, Shenghao Xie, and Samson Zhou. Perfect sampling in turnstile streams beyond  
 928 small moments. *Proc. ACM Manag. Data*, 3(2):106:1–106:27, 2025.

929  
 930 Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming  
 931 language models with attention sinks. In *The Twelfth International Conference on Learning  
 932 Representations, ICLR*. OpenReview.net, 2024.

933 Guangxuan Xiao, Jiaming Tang, Jingwei Zuo, Junxian Guo, Shang Yang, Haotian Tang, Yao Fu, and  
 934 Song Han. Duoattention: Efficient long-context LLM inference with retrieval and streaming heads.  
 935 In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore,  
 936 April 24-28, 2025*. OpenReview.net, 2025.

937  
 938 Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le.  
 939 Xlnet: Generalized autoregressive pretraining for language understanding. *Advances in neural  
 940 information processing systems*, 32, 2019.

941 Yao Yao, Zuchao Li, and Hai Zhao. Sirllm: Streaming infinite retentive LLM. In Lun-Wei Ku, Andre  
 942 Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association  
 943 for Computational Linguistics (Volume 1: Long Papers)*, ACL, pp. 2611–2624. Association for  
 944 Computational Linguistics, 2024.

945  
 946 Guanghao Ye. Fast algorithm for solving structured convex programs. *The University of Washington,  
 Undergraduate Thesis*, 2020.

947  
 948 Amir Zandieh, Insu Han, Majid Daliri, and Amin Karbasi. Kdeformer: Accelerating transformers via  
 949 kernel density estimation. In *ICML*. arXiv preprint arXiv:2302.02451, 2023.

950 Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Diego Rojas, Guanyu Feng, Hanlin  
 951 Zhao, Hanyu Lai, Hao Yu, Hongning Wang, Jiadai Sun, Jiajie Zhang, Jiale Cheng, Jiayi Gui,  
 952 Jie Tang, Jing Zhang, Juanzi Li, Lei Zhao, Lindong Wu, Lucen Zhong, Mingdao Liu, Minlie  
 953 Huang, Peng Zhang, Qinkai Zheng, Rui Lu, Shuaiqi Duan, Shudan Zhang, Shulin Cao, Shuxun  
 954 Yang, Weng Lam Tam, Wenyi Zhao, Xiao Liu, Xiao Xia, Xiaohan Zhang, Xiaotao Gu, Xin Lv,  
 955 Xinghan Liu, Xinyi Liu, Xinyue Yang, Xixuan Song, Xunkai Zhang, Yifan An, Yifan Xu, Yilin  
 956 Niu, Yuantao Yang, Yueyan Li, Yushi Bai, Yuxiao Dong, Zehan Qi, Zhaoyu Wang, Zhen Yang,  
 957 Zhengxiao Du, Zhenyu Hou, and Zihan Wang. Chatglm: A family of large language models from  
 958 GLM-130B to GLM-4 all tools. *CoRR*, abs/2406.12793, 2024. doi: 10.48550/ARXIV.2406.12793.  
 959 URL <https://doi.org/10.48550/arXiv.2406.12793>.

960 Jingzhao Zhang, Sai Praneeth Karimireddy, Andreas Veit, Seungyeon Kim, Sashank Reddi, Sanjiv  
 961 Kumar, and Suvrit Sra. Why are adaptive methods good for attention models? *Advances in Neural  
 962 Information Processing Systems*, 33:15383–15393, 2020.

963 Jintao Zhang, Chendong Xiang, Haofeng Huang, Jia Wei, Haocheng Xi, Jun Zhu, and Jianfei Chen.  
 964 Spargeattn: Accurate sparse attention accelerating any model inference. *CoRR*, abs/2502.18137,  
 965 2025. doi: 10.48550/ARXIV.2502.18137. URL <https://doi.org/10.48550/arXiv.2502.18137>.

966  
 967 Lichen Zhang. Speeding up optimizations via data structures: Faster search, sample and maintenance.  
 968 Master’s thesis, Carnegie Mellon University, 2022.

969  
 970 Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuhui Chen, Christopher  
 971 Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language  
 972 models. *arXiv preprint arXiv:2205.01068*, 2022.

972 Xinrong Zhang, Yingfa Chen, Shengding Hu, Zihang Xu, Junhao Chen, Moo Khai Hao, Xu Han,  
973 Zhen Leng Thai, Shuo Wang, Zhiyuan Liu, and Maosong Sun.  $\infty$ bench: Extending long context  
974 evaluation beyond 100k tokens. *CoRR*, abs/2402.13718, 2024. URL <https://doi.org/10.48550/arXiv.2402.13718>.

975

976 Haoyu Zhao, Abhishek Panigrahi, Rong Ge, and Sanjeev Arora. Do transformers parse while  
977 predicting the masked word? In *Proceedings of the 2023 Conference on Empirical Methods in*  
978 *Natural Language Processing, EMNLP 2023*, pp. 16513–16542. Association for Computational  
979 Linguistics, 2023.

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1000

1001

1002

1003

1004

1005

1006

1007

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

# 1026 1027 1028 1029 1030 1031 1032 1033 1034 1035 1036 1037 1038 1039 1040 1041 1042 1043 1044 1045 1046 1047 1048 1049 1050 1051 1052 1053 1054 1055 1056 1057 1058 1059 1060 1061 1062 1063 1064 1065 1066 1067 1068 1069 1070 1071 1072 1073 1074 1075 1076 1077 1078 1079 Appendix

**Roadmap.** In Section A, we provide additional related work. In Section B, we briefly discuss the background on  $\ell_2$  sampler. In Section C, we show that how to use the tail bound to obtain sampling result. In Section D, we present the tensor sampling result. In Section E, we discuss the LLM usage of the paper.

## A ADDITIONAL RELATED WORK

**On sketching.** The application of sketching and sampling techniques in numerical linear algebra has been remarkably effective, revolutionizing a broad spectrum of core tasks. These methods are crucial in linear programming (LP), as evidenced by Cohen et al. (2019); Jiang et al. (2021); Ye (2020); Gu & Song (2022), and have significantly impacted tensor approximation (Song et al., 2019; Mahankali et al., 2024; Deng et al., 2023a). Sketching and sampling techniques also have been widely applied in matrix completion (Gu et al., 2024), matrix sensing (Qin et al., 2024; Deng et al., 2023b), submodular function maximization (Qin et al., 2023a), dynamic sparsification (Deng et al., 2022a), dynamic tensor product regression (Reddy et al., 2022), and semi-definite programming (Song et al., 2023a). Additionally, sketching has been pivotal in iterative sparsification problems (Song et al., 2022), adversarial training (Gao et al., 2022), kernel density estimation (Qin et al., 2022b), solving the distance oracle problem (Deng et al., 2022b), and empirical risk minimization (Lee et al., 2019; Qin et al., 2023b). Its applications furthermore extends to relational databases (Qin et al., 2022a) and Large Language Model (LLM) research (Deng et al., 2023c;b; Gao et al., 2025; Li et al., 2023).

**On theoretical attention.** A comprehensive body of research, including studies (Child et al., 2019; Kitaev et al., 2020; Wang et al., 2020; Daras et al., 2020; Katharopoulos et al., 2020; Chen et al., 2021; 2022; Zandieh et al., 2023; Alman & Song, 2023; Brand et al., 2024; Deng et al., 2023c; Kacham et al., 2023; Alman & Song, 2024a; Han et al., 2024; Awasthi & Gupta, 2023; Marcus et al., 2022; Alman & Song, 2024b; 2025a;b), has progressively shed light on the complexities and optimization of attention matrix computation. This exploration has been further enriched by insights into the effectiveness of attention mechanisms in Transformers (Dehghani et al., 2018; Vuckovic et al., 2020; Zhang et al., 2020; Edelman et al., 2022; Snell et al., 2021; Wei et al., 2021; Deng et al., 2023d; 2024a). Among these, Zhao et al. (2023) revealed the adeptness of mid-scale masked language models in identifying syntactic elements, paving the way for innovations like partial parse tree reconstructions. Inspired the exponential mechanism in attention structure, Gao et al. (2023) provide an analysis which shows exponential regression within the over-parameterized neural tangent kernel framework can converge. In the over-constrained setting, several work show the convergence for attention inspired regression problem (Li et al., 2023; Deng et al., 2023b).

## B $\ell_2$ SAMPLER

We give the full details of the standard  $L_2$  sampler from Jowhari et al. (2011); Mahabadi et al. (2020) in Algorithm 2. The proof of correctness is verbatim from Jowhari et al. (2011); Mahabadi et al. (2020). The challenge is how to implement the data structures of  $y$ , which is implicitly defined as  $(A_1 \otimes A_2)x$ . By comparison, in the standard setting of  $\ell_2$  samplers Monemizadeh & Woodruff (2010); Andoni et al. (2011); Jowhari et al. (2011); Jayaram & Woodruff (2021); Mahabadi et al. (2020),  $y$  is given as a data stream.

## C FROM TAIL TO SAMPLING

**Lemma C.1** (Restatement of Lemma 7.4). *Let  $y = (A_1 \otimes A_2)x \in \mathbb{R}^{n^2}$ . Let only one of  $A_1$  and  $A_2$  be updated in streaming. Let  $w = \frac{y_i}{\sqrt{u_i}}$  for a constant  $u_i \in [0, 1]$  generated uniformly at random. There is an algorithm  $\mathcal{A}$  that that uses  $O(nd) + \text{poly}(\frac{1}{\epsilon}, \log n)$  space, uses  $O(n)$  update time, and estimates each element of  $w$  up to additive error  $\epsilon \cdot \|z\|_2$ , where  $z$  denotes the tail vector of  $w$  without the largest  $\frac{1}{\epsilon^2}$  entries in magnitude. Specifically, for all  $i \in [n^2]$ , we have  $|\hat{w}_i - w_i| \leq \epsilon \cdot \|z\|_2$ .*

---

1080 **Algorithm 2** Standard  $\ell_2$  Sampler, e.g., extension of Jowhari et al. (2011) to  $p = 2$

---

1081 1: For each  $i \in [n]$ , let  $u_i \in [0, 1]$  be chosen uniformly at random

1082 2:  $w_i \leftarrow \frac{y_i}{\sqrt{u_i}}$

1083 3: Let  $z$  denote the tail vector of  $w$  without the largest  $\frac{1}{\epsilon^2}$  entries in magnitude

1084 4: Let  $\hat{Y}$  be a 2-approximation of  $\|y\|_2$

1085 5: Let  $\hat{Z}$  be a 2-approximation of  $\|z\|_2$

1086 6:  $i \leftarrow \operatorname{argmax}_{i \in [n]} |\hat{w}_i|$

1087 7: Let  $C > 0$  be a large constant determined by the additive failure probability  $\frac{1}{\operatorname{poly}(n)}$

1088 8: **if**  $\hat{Z} > \sqrt{\frac{C \log n}{\epsilon}} \cdot \hat{Y}$  or  $|w_i| < \sqrt{\frac{C \log n}{\epsilon}} \cdot \hat{Y}$  **then**

1089 9:     Return FAIL

1090 10: **else**

1091 11:     Return  $i$  with estimate  $\sqrt{u_i} \cdot \hat{w}_i$

1092 12: **end if**

---

1093

1094

1095

1096 *Proof.* Consider hash function  $h_1, h_2 : [n] \rightarrow [b]$ . Consider random sign functions  $\sigma_1, \sigma_2 : [n] \rightarrow \{-1, +1\}$ . We consider a fixed index  $i_1, i_2 \in [n]$ . Let  $j = h_1(i_1) + h_2(i_2) \pmod{b}$ . Let  $h^{-1}(j)$  denote the all the pairs  $(i_1, i_2) \in [n] \times [n]$  such that  $h_1(i_1) + h_2(i_2) \pmod{b} = j$ . Note that  $\hat{y}_i$  induced by  $h$  is  $\hat{w}_i = w_i + \sum_{l \in h^{-1}(j) \setminus \{i\}} s_i s_l w_{l_1} w_{l_2}$ . For ease of presentation, we write  $\sigma_i = \sigma_{1, i_1} \sigma_{2, i_2}$  and  $\sigma_l = \sigma_{1, l_1} \sigma_{2, l_2}$ .

1101 
$$\mathbb{E}[\hat{w}_i] = \mathbb{E}\left[w_i + \sum_{l \in h^{-1}(j) \setminus \{i\}} \sigma(i)\sigma(l)w_l\right] = \mathbb{E}[w_i] + \sum_{l \in h^{-1}(j) \setminus \{i\}} \mathbb{E}[\sigma(i) \cdot \sigma(l)] \cdot w_l$$

1102 
$$= w_i + \sum_{l \in h^{-1}(j) \setminus \{i\}} \mathbb{E}[\sigma(i)] \cdot \mathbb{E}[\sigma(l)] \cdot w_l = w_i,$$

1103 where the first step follows from definition, the second step follows from linearity of expectation, the third step follows from  $\sigma(i)$  and  $\sigma(l)$  are independent, the forth step follows from  $\mathbb{E}[\sigma(l)] = 0$ .

1104 We now upper bound the variance of  $\hat{w}_i - y_i$  by analyzing  $\mathbb{E}[(\hat{y}_i)^2]$ . Let  $\mathcal{H}$  be the set of the top  $\frac{1}{\epsilon^2}$  items and let  $\mathcal{E}$  be the event that none of the items in  $\mathcal{H}$  are mapped to  $h(i)$ , i.e.,  $h(a) \neq h(i)$  for all  $a \in \mathcal{H}$ .

1105 Observe that for  $b = \frac{100}{\epsilon^2}$ , we have that  $\Pr[\mathcal{E}] \geq 0.9$ . Then we have:

1106 
$$\mathbb{E}[(\hat{w}_i - w_i)^2 \mid \mathcal{E}] = \mathbb{E}\left[\left(\sum_{l \in [n]^2 \setminus \mathcal{H}, l \in h^{-1}(j)} \sigma(i)\sigma(l)w_l\right)^2\right] = \mathbb{E}\left[\sum_{l \in [n]^2 \setminus \mathcal{H}, l \in h^{-1}(j)} w_l^2\right]$$

1107 
$$= \frac{1}{b} \cdot \sum_{l \in [n]^2 \setminus \mathcal{H}, l \in h^{-1}(j)} w_l^2 \leq \frac{1}{b} \cdot (w_1^2 + \dots + w_{n^2}^2 - \sum_{l \in \mathcal{H}} w_l^2)$$

1108 
$$= 100\epsilon^2 \cdot \|z\|_2^2,$$

1109 for  $b = \frac{100}{\epsilon^2}$ , since  $z$  is the vector corresponding to  $y$  that removes the entries in  $\mathcal{H}$ . By Cheby-  
shev's inequality, we have that  $\Pr[|\hat{w}_i - w_i| \geq \epsilon \cdot \|z\|_2 \mid \mathcal{E}] \leq \frac{1}{10}$ . Since  $\Pr[\mathcal{E}] \geq 0.9$ , then  
1110  $\Pr[|\hat{w}_i - w_i| \geq \epsilon \cdot \|z\|_2] \leq 0.2$ , for a fixed hash function  $h$ . By taking the median of  $O(\log n)$  estima-  
1111 tions corresponding to  $O(\log n)$  different hash functions  $h$ , we have that  $\Pr[|\hat{w}_i - w_i| \geq \epsilon \cdot \|z\|_2] \leq$   
1112  $\frac{1}{n^{10}}$ . Thus by a union bound over  $i \in [n] \times [n]$ , we have that with probability at least  $1 - \frac{1}{n^5}$ , we have  
1113 for all  $i \in [n]$ ,  $|\hat{w}_i - w_i| \geq \epsilon \cdot \|z\|_2$ .  $\square$

1114 **Lemma C.2** (Restatement of Lemma 7.5). *Let  $y = (A_1 \otimes A_2)x \in \mathbb{R}^{n^2}$  and let  $w \in \mathbb{R}^{n^2}$  so that*  
1115  *$w_i = \frac{y_i}{\sqrt{u_i}}$  for a constant  $u_i \in [0, 1]$  generated uniformly at random. Let  $z$  denote the tail vector*  
1116 *of  $w$  without the largest  $\frac{1}{\epsilon^2}$  entries in magnitude. Let  $\hat{Z}$  be a 2-approximation to  $\|z\|_2$  and  $\hat{Y}$  be a*  
1117 *2-approximation to  $\|y\|_2$ . Then*

1118 
$$\Pr\left[\hat{Z} > \sqrt{\frac{C \log n}{\epsilon}} \cdot \hat{Y}\right] \leq O(\epsilon) + \frac{1}{\operatorname{poly}(n)}.$$

1134 *Proof.* Let  $\mathcal{E}_1$  denote the event that  $\widehat{Z}$  is a 2-approximation to  $\|z\|_2$  and  $\widehat{Y}$  is a 2-approximation to  
 1135  $\|y\|_2$ , so that  
 1136

$$1137 \quad \Pr[\mathcal{E}_1] \geq 1 - \frac{1}{\text{poly}(n)}. \\ 1138$$

1140 Conditioned on  $\mathcal{E}_1$ , it suffices to bound the probability that  
 1141

$$1142 \quad 4\|z\|_2 > \sqrt{\frac{C \log n}{\epsilon}} \cdot \|y\|_2. \\ 1143$$

1144 Let  $j \in [n^2]$  be a fixed index and let  $u_j$  be fixed.  
 1145

1146 Let  $T = \sqrt{\epsilon} \cdot \|y\|_2$  and for each  $i \in [n^2]$ , we define the indicator random variable  $W_i = 1$  if  $|w_i| > T$   
 1147 and  $W_i = 0$  otherwise, if  $|w_i| \leq T$ . Note that  $W_i$  is an indicator random variable for whether the  
 1148 coordinate  $w_i$  in the vector  $w$  is “heavy” in magnitude.  
 1149

1150 We then define

$$1151 \quad Z_i = \frac{w_i^2}{T^2} \cdot (1 - W_i) \\ 1152$$

1153 to be the scaled contribution of the small entries of  $z$ , and observe that  $Z_i \in [0, 1]$ .  
 1154

1155 Let

$$1156 \quad W = \sum_{i \in [n^2], i \neq j} w_i \\ 1158$$

1159 denote the total number of heavy indices besides possibly index  $j$  and  $Z = \sum_{i \in [n^2], i \neq j} Z_i$  denote  
 1160 the total scaled contribution of the light indices besides possibly index  $j$ . Let  $v$  denote the vector  
 1161 containing the heavy indices, so that  $v_i = w_i$  for  $W_i = 1$  and  $v_i = 0$  otherwise for  $W_i = 0$ . Note that  
 1162  $v$  has sparsity at most  $Y + 1$  and moreover  $U^2 Z = \|w - v\|_2^2$ . We also have that  $\|z\|_2 \leq \|w - v\|_2$   
 1163 unless  $W \geq \frac{2}{\epsilon^2}$ .  
 1164

1165 Let  $\mathcal{E}_2$  denote the event that  $W \geq \frac{2}{\epsilon^2}$  and let  $\mathcal{E}_3$  denote the event that  $Z \geq \frac{C \log n}{16T^2\epsilon} \cdot \|y\|_2^2$ . Observe  
 1166 that if neither  $\mathcal{E}_2$  nor  $\mathcal{E}_3$  occur, then we have  $4\|z\|_2 \leq \sqrt{\frac{C \log n}{\epsilon}} \cdot \|y\|_2$ , as desired. Thus it remains to  
 1167 bound the probability of the failure events  $\mathcal{E}_2$  and  $\mathcal{E}_3$ .  
 1168

1169 We have  $\mathbb{E}[W_i] = \frac{\|w\|_2^2}{T^2}$ , so that  $\mathbb{E}[W] \leq \frac{1}{\epsilon}$ . By Markov’s inequality, we have that  $\Pr[\mathcal{E}_2] \leq \frac{\epsilon}{2}$ .  
 1170

1171 We now upper bound  $\Pr[\mathcal{E}_3]$ . Recall that  $Z_i = \frac{w_i^2}{T^2} \cdot (1 - W_i) = \frac{w_i^2}{T u_i^2} \cdot (1 - W_i)$ , since  $w_i = \frac{y_i}{\sqrt{u_i}}$ .  
 1172

1173 Observe that  $Z_i > 0$  only if  $|w_i| < T$ , i.e., if  $u_i \geq \frac{y_i^2}{\epsilon \cdot \|y\|_2^2}$ , since  $T = \sqrt{\epsilon} \cdot \|y\|_2$ . For  $\epsilon \in (0, 1)$ , we  
 1174 thus have

$$1175 \quad \mathbb{E}[Z_i] \leq \int_{y_i^2/\|y\|_2^2}^1 z_i du_i \\ 1176 \\ 1177 \\ 1178 \\ 1179 \\ 1180 \quad = \int_{y_i^2/\|y\|_2^2}^1 \frac{y_i^2}{u_i} \frac{1}{T^2} du_i.$$

1181 Now, let  $\mathcal{E}_4$  be the event that  $u_i \geq \frac{1}{n^{C/2}}$  for all  $i \in [n^2]$ , so that  $\Pr[\mathcal{E}_4] \geq 1 - \frac{1}{n^{C/2-2}}$ .  
 1182

1183 Then

$$1184 \quad \mathbb{E}[Z_i | \mathcal{E}_4] \leq \frac{1}{1 - \frac{1}{n^{C/2-2}}} \int_{1/n^{C/2}}^1 \frac{y_i^2}{u_i} \frac{1}{T^2} du_i \\ 1185 \\ 1186 \\ 1187 \leq \frac{C \log n}{T^2} y_i^2.$$

1188 Thus, we have

$$\begin{aligned}
 \mathbb{E}[Z \mid \mathcal{E}_4] &= \sum_{i \in [n^2]} \mathbb{E}[Z_i \mid \mathcal{E}_4] \\
 &= \sum_{i \in [n^2]} \frac{C \log n}{T^2} y_i^2 \\
 &\leq \sum_{i \in [n^2]} \frac{C \log n}{\epsilon} \frac{y_i^2}{\|y\|_2^2} \\
 &= \frac{C \log n}{\epsilon}.
 \end{aligned}$$

1200 Thus by Markov's inequality, the probability that  $Z$  is larger than  $\frac{C \log n}{16T^2\epsilon} \cdot \|y\|_2^2 = \frac{C \log n}{16\epsilon^2}$  is at most  
 1201  $\frac{\epsilon}{16}$ . The claim then follows from taking a union bound over the events  $\mathcal{E}_1, \neg\mathcal{E}_2, \neg\mathcal{E}_3, \neg\mathcal{E}_4$ .  $\square$

## D TENSOR SAMPLING

1205 **Theorem D.1** (Restatement of Theorem 7.6). *Let  $y = (A_1 \otimes A_2)x \in \mathbb{R}^{n^2}$  and let  $w \in \mathbb{R}^{n^2}$  so that  
 1206 for each  $i \in [n^2]$ ,  $w_i = \frac{y_i}{\sqrt{u_i}}$  for a constant  $u_i \in [0, 1]$  generated uniformly at random. Let  $z$  denote  
 1207 the tail vector of  $w$  without the largest  $\frac{1}{\epsilon^2}$  entries in magnitude. Suppose there exists:*

- 1209 1. An algorithm  $\mathcal{A}_1$  that provides a 2-approximation to  $\|y\|_2$  with probability  $1 - \frac{1}{n^2}$ .
- 1210 2. An algorithm  $\mathcal{A}_2$  that provides a 2-approximation to  $\|z\|_2$  with probability  $1 - \frac{1}{n^2}$ .
- 1212 3. An algorithm  $\mathcal{A}_3$  that estimates each element of  $w$  up to additive error  $\epsilon \cdot \|z\|_2$ ,

$$|\widehat{w}_i - w_i| \leq \epsilon \cdot \|z\|_2,$$

1215 for all  $i \in [n^2]$ .

1216 Then there exists a data structure that uses  $\text{poly}(\frac{1}{\epsilon}, \log n)$  bits of space and outputs each index  $i$   
 1217 with probability  $p_i$ , such that

$$(1 - \epsilon) \cdot \frac{y_i^2}{\|y\|_2^2} - \frac{1}{\text{poly}(n)} \leq p_i \leq (1 + \epsilon) \cdot \frac{y_i^2}{\|y\|_2^2} + \frac{1}{\text{poly}(n)}.$$

1222 *Proof.* Let  $i$  be fixed and let  $\mathcal{E}$  denote the event that  $u_i < \frac{\epsilon}{C \log n} \frac{y_i^2}{\widehat{Y}^2}$ , so that  $|w_i| > \sqrt{\frac{C \log n}{\epsilon}} \cdot \widehat{Y}$ .

1224 Let  $\mathcal{E}_1$  denote the event that  $\widehat{Y}$  is a 2-approximation to  $\|y\|_2$ ,  $\widehat{Z}$  is a 2-approximation to  $\|z\|_2$ , and  
 1225  $|\widehat{w}_i - w_i| \leq \epsilon \cdot \|z\|_2$  for all  $i \in [n]$ . Let  $\mathcal{E}_2$  denote the event that  $\widehat{Z} > \sqrt{\frac{C \log n}{\epsilon}} \cdot \widehat{Y}$  and let  $\mathcal{E}_3$  denote  
 1227 the event that multiple indices  $j$  satisfy  $|w_j| > \sqrt{\frac{C \log n}{\epsilon}} \cdot \widehat{Y}$ . Finally, let  $\mathcal{E}_4$  denote the event that  
 1228  $|\widehat{w}_i| < \sqrt{\frac{C \log n}{\epsilon}} \cdot \widehat{Y}$ .

1230 Intuitively,  $\mathcal{E}_1$  is a good event, i.e., correctness of the data structures, which we would like to hold.  
 1231 On the other hand,  $\mathcal{E}_2, \mathcal{E}_3, \mathcal{E}_4$  are bad events that distort the sampling probabilities, which we would  
 1232 like to avoid.

1233 We first note that  $\mathcal{E}_1$  holds with high probability due to the correctness of the CountSketch and  
 1234  $L_2$ -norm estimation data structures. We next note that by Lemma 7.5, the probability that  $\mathcal{E}_2$  occurs  
 1235 is  $O(\epsilon)$ .

1237 Next, note that the probability that for a fixed  $j \in [n]$ ,  $u_j$  satisfies  $\frac{y_j^2}{u_j} \geq \frac{C \log n}{\epsilon} \cdot \widehat{Y}$  is at most  
 1238  $\frac{\epsilon}{C' \log n} \frac{y_j^2}{\|y\|_2^2}$  for some constant  $C'$ . Thus summing over all  $j \in [n]$ , the probability that there exist  
 1239 an additional  $j \in [n]$  for which  $|w_j| > \sqrt{\frac{C \log n}{\epsilon}} \cdot \widehat{Y}$  is  $O(\epsilon)$ . Thus the probability that  $\mathcal{E}_3$  occurs is  
 1240  $O(\epsilon)$ .

Finally, conditioned on  $-\mathcal{E}_2$ , we have that  $\widehat{Z} \leq \sqrt{\frac{C \log n}{\epsilon}} \cdot \widehat{Y}$ . Then conditioning on  $\mathcal{E}_1$ , we have  $\|z\|_2 \leq \widehat{Z}$  and thus  $|\widehat{w}_i - w_i| \leq \epsilon \widehat{Z} \leq \sqrt{C \epsilon \log n} \widehat{Y}$ , so that  $\mathcal{E}_4$  can only occur for  $\sqrt{\frac{C \log n}{\epsilon}} \cdot \widehat{Y} \leq |w_i| \leq \sqrt{\frac{C \log n}{\epsilon}} \cdot \widehat{Y}$ , which is at most probability  $O\left(\frac{\epsilon^2}{C \log n} \frac{y_i^2}{\widehat{Y}^2}\right)$ , over the randomness of  $u_i$ .

In summary, we observe that conditioned on some value being output, the probability that item  $i$  is selected is proportional to the event that the events  $\mathcal{E}$  and  $\mathcal{E}_1$  occur, and none of the events  $\mathcal{E}_2, \mathcal{E}_3, \mathcal{E}_4$  occur. The probability that  $\mathcal{E}$  occurs is  $\frac{\epsilon}{C \log n} \frac{y_i^2}{\widehat{Y}^2}$ , which  $u_i$  is chosen uniformly at random. Due to the event  $\mathcal{E}_1$ , the sampling probability is distorted additively by  $\frac{1}{\text{poly}(n)}$ , while due to the events  $\mathcal{E}_2, \mathcal{E}_3, \mathcal{E}_4$ , the sampling probability is distorted multiplicatively by  $(1 + \epsilon)$ . Thus conditioned on the event that some index is returned, the probability  $p_i$  that index  $i$  is returned satisfies

$$(1 - \epsilon) \cdot \frac{y_i^2}{\|y\|_2^2} - \frac{1}{\text{poly}(n)} \leq p_i \leq (1 + \epsilon) \cdot \frac{y_i^2}{\|y\|_2^2} + \frac{1}{\text{poly}(n)},$$

as desired.  $\square$

## E LLM USAGE DISCLOSURE

LLMs were used only to polish language, such as grammar and wording. These models did not contribute to idea creation or writing, and the authors take full responsibility for this paper's content.