

SUBLINEAR TIME QUANTUM ALGORITHM FOR ATTENTION APPROXIMATION

Zhao Song

Simons Institute for the Theory of Computing, UC Berkeley
magic.linuxkde@gmail.com

Jianfei Xue

New York University
jx898@nyu.edu

Jiahao Zhang

ml.jiahaozhang02@gmail.com

Lichen Zhang

MIT CSAIL
lichenz@csail.mit.edu

ABSTRACT

Given the query, key and value matrices $Q, K, V \in \mathbb{R}^{n \times d}$, the attention module is defined as $\text{Att}(Q, K, V) = D^{-1}AV$ where $A = \exp(QK^\top/\sqrt{d})$ with $\exp(\cdot)$ applied entrywise, $D = \text{diag}(A\mathbf{1}_n)$. The attention module is the backbone of modern transformers and large language models, but explicitly forming the softmax matrix $D^{-1}A$ incurs $\Omega(n^2)$ time, motivating numerous approximation schemes that reduce runtime to $\tilde{O}(nd)$ via sparsity or low-rank factorization.

We propose a quantum data structure that approximates any row of $\text{Att}(Q, K, V)$ using only row queries to Q, K, V . Our algorithm preprocesses these matrices in $\tilde{O}(\epsilon^{-1}n^{0.5}(s_\lambda^{2.5} + s_\lambda^{1.5}d + \alpha^{0.5}d))$ time, where ϵ is the target accuracy, s_λ is the λ -statistical dimension of the exponential kernel defined by Q and K , and α measures the row distortion of V that is at most $d/\text{srnk}(V)$, the stable rank of V . Each row query can be answered in $\tilde{O}(s_\lambda^2 + s_\lambda d)$ time.

To our knowledge, this is the first quantum data structure that approximates rows of the attention matrix in sublinear time with respect to n . Our approach relies on a quantum Nyström approximation of the exponential kernel, quantum multivariate mean estimation for computing D , and quantum leverage score sampling for the multiplication with V .

1 INTRODUCTION

Transformers (Vaswani et al., 2017) have emerged as one of the most successful machine learning architectures in recent years, revolutionizing fields such as natural language processing (Devlin et al., 2019; Yang et al., 2019; Raffel et al., 2020; Brown et al., 2020; Jiao et al., 2020), computer vision (Carion et al., 2020; Dosovitskiy et al., 2021; Guo et al., 2022), speech recognition (Chorowski et al., 2015; Wang et al., 2021), robotics (Liu et al., 2022), and time series forecasting (Zhou et al., 2021). These models typically operate on sequences of length n , autoregressively predicting the next most likely token to produce an output of length n . In applications like large language models (LLMs), it has been widely observed that increasing the sequence length n significantly enhances generative performance. However, this benefit comes at a substantial computational cost: the core attention module has a quadratic time complexity in n , which severely limits both training and inference scalability.

Formally, let $Q, K, V \in \mathbb{R}^{n \times d}$ denote the query, key, and value embeddings. The attention module is defined as

$$\text{Att}(Q, K, V) = D^{-1}AV \in \mathbb{R}^{n \times d},$$

where $A = \exp(QK^\top/\sqrt{d}) \in \mathbb{R}^{n \times n}$ is computed entrywise, and $D = \text{diag}(A\mathbf{1}_n) \in \mathbb{R}^{n \times n}$. The matrix A is referred to as the *attention matrix*, and $D^{-1}A$ as the *softmax matrix*. Due to the $n \times n$ size of A , much recent research has focused on reducing the quadratic complexity by approximating attention through pattern-based sparse attention (Daras et al., 2020; Kitaev et al., 2020; Roy et al.,

2021; Sun et al., 2022; Child et al., 2019; Beltagy et al., 2020; Ainslie et al., 2020; Zaheer et al., 2020), linearizing the kernel through feature mapping (Katharopoulos et al., 2020; Choromanski et al., 2021; Wang et al., 2020; Peng et al., 2021), or various algorithmic and data structure optimizations (Zandieh et al., 2023; Alman & Song, 2023; Han et al., 2024; Kacham et al., 2024; Zandieh et al., 2024; van den Brand et al., 2024; Song et al., 2024; Kannan et al., 2025; Chu et al., 2024; Chen et al., 2025b; Indyk et al., 2025).

The theoretical goal in these efforts is to achieve a runtime that scales nearly linearly with n , allowing some approximation error. This is a natural target, since the input size to the attention module is $n \times d$. On a classical computer, any algorithm that approximates attention in time $\tilde{O}(nd)$ is considered optimal. But could this process be accelerated further using a quantum computer?

If our objective is to output the entire $n \times d$ matrix $\text{Att}(Q, K, V)$, then $\Omega(nd)$ time is unavoidable due to output size. However, in many transformer applications — particularly during inference (Pope et al., 2023; Brandon et al., 2024; Adnan et al., 2024; Zhang et al., 2024a; Feng et al., 2025; Liu et al., 2024b; Kumari et al., 2024; Behnam et al., 2025; Chen et al., 2025a;c; Indyk et al., 2025) — only *row queries* are needed. In this setting, we aim to preprocess Q, K, V into a data structure such that, for any index $i \in [n]$, the structure can return a vector $\tilde{r}_i \in \mathbb{R}^d$ that approximates the i -th row of $\text{Att}(Q, K, V)$. This model circumvents the $\Omega(nd)$ lower bound by focusing on partial output. Nonetheless, since each row of $\text{Att}(Q, K, V)$ is a convex combination of rows of V , achieving truly sublinear time in n still appears classically intractable.

In this work, we answer this question affirmatively. Specifically, we construct a quantum data structure that preprocesses Q, K, V using only *row queries*, and does so in time¹ $\tilde{O}(\epsilon^{-1}n^{0.5} \cdot \text{poly}(d, s_\lambda, \alpha))$, where s_λ is the *statistical dimension* of the exponential kernel matrix associated with Q and K , and α is a measure of the row distortion of V (see Definition C.2). Given any index $i \in [n]$, the data structure returns an approximation to the i -th row of $\text{Att}(Q, K, V)$ in time $\tilde{O}(s_\lambda^2 + s_\lambda d)$.

To our knowledge, this is the first quantum algorithm to implement the row query model in sublinear time. Prior works either require superlinear preprocessing time or impose structural assumptions (Gao et al., 2023). Our approach avoids both: it makes *no assumptions* on Q, K, V , making it broadly applicable in practice. Moreover, our construction is conceptually simple — it combines quantum techniques such as Grover search (Grover, 1996), Nyström kernel approximation, and quantum multivariate mean estimation (Cornelissen et al., 2022) to approximate each component of the attention module: D, A , and V .

Quantum Computation Model. We follow the standard quantum computation framework as in Apers & De Wolf (2022); Apers & Gribling (2023). The model allows quantum subroutines using $O(\log n)$ qubits, quantum queries to the input, and access to a quantum-read/classical-write RAM (QRAM) of $\text{poly}(n)$ bits. Each quantum read or classical write takes unit cost. We measure *time complexity* by the number of QRAM operations, and *query complexity* by the number of queries to the input. In our setting, we query rows of Q, K , and V , each requiring $O(d)$ time classically. For simplicity, we assume Q and K have been scaled by $1/d^{1/4}$, which can also be done via row queries in $O(d)$ time.

2 PRELIMINARY

Notation. Given symmetric matrices $A, B \in \mathbb{R}^{n \times n}$, we use $A - B \succeq 0$ to denote $A - B$ is a positive semidefinite (PSD) matrix, i.e., for any $x \in \mathbb{R}^n$, $x^\top(A - B)x \geq 0$. Given a matrix $M \in \mathbb{R}^{n \times n}$, we use $\exp(M)$ to denote the entrywise exponentiation operation. We use $\text{tr}[M]$ to denote the trace of M . For a real matrix A , we use A^\dagger to denote its Moore-Penrose pseudoinverse, and for a square, nonsingular real matrix M , we use M^{-1} to denote its inverse. For two vectors $x, y \in \mathbb{R}^n$, we use $x^\top y$ or $\langle x, y \rangle$ to denote the inner product of x and y . We use $\mathbf{0}_n$ and $\mathbf{1}_n$ to denote all-0’s and all-1’s vector. For a vector $x \in \mathbb{R}^n$, we use $\|x\|_2 = \sqrt{x^\top x}$ to denote its ℓ_2 norm, $\|x\|_\infty = \max_{i \in [n]} |x_i|$ to denote its ℓ_∞ norm. If M is a PSD matrix, then we use $\|x\|_M = \sqrt{x^\top M x}$ to denote the M -energy norm of x . For a matrix A , we use $\|A\|$ to denote its spectral norm and $\|A\|_\infty$ to denote its max row ℓ_1 norm, and $\|A\|_F$ to denote its Frobenius norm. Throughout the paper, we will also exclusively

¹We use $\tilde{O}(\cdot)$ to suppress polylogarithmic factors in n, d, s_λ , and $1/\epsilon$.

work with weighted sampling matrices, usually denoted by $S \in \mathbb{R}^{n \times s}$ for where s is the total number of samples taken, let $i(j)$ be the index of the i -th sample, then the i -th column of S is $\frac{1}{\sqrt{p_j}} e_j$, where p_j is the probability of choosing the index j . We use $\mathbb{E}[X]$ to denote the expectation of a random variable X . We use $\mathbb{I}[E]$ to denote the indicator of whether event E happens.

Numerical Linear Algebra. We rely on several primitives from numerical linear algebra for fast approximations and provable guarantees.

Definition 2.1 (Leverage score). *Let $A \in \mathbb{R}^{n \times d}$. The i -th leverage score of A is defined as*

$$\tau_i := a_i^\top (A^\top A)^{-1} a_i,$$

where a_i is the i -th row of A . Equivalently, let $A = U\Sigma V^\top$ be its SVD, then $\tau_i = \|u_i\|_2^2$, where u_i is the i -th row of U .

We will also work exclusively with *kernel matrices*. Given a dataset $X = \{x_1, \dots, x_n\} \subseteq \mathbb{R}^d$, we define the exponential kernel matrix $E \in \mathbb{R}^{n \times n}$ by $E_{i,j} = \exp(\langle x_i, x_j \rangle)$. Although E is generally full-rank, our algorithm depends only on a parameter called the λ -statistical dimension of E , which may be much smaller than n .

Definition 2.2 (Statistical dimension (Zhang, 2005; Hastie et al., 2009)). *Let $E \in \mathbb{R}^{n \times n}$ be a PSD matrix, and let $\lambda > 0$. The λ -statistical dimension of E is defined as $s_\lambda(E) := \text{tr}[E(E + \lambda I)^{-1}]$. When E is clear from context, we write s_λ for simplicity.*

Note that s_λ is a monotonically decreasing function of λ , and is closely related to the notion of ridge leverage scores.

Definition 2.3 (Ridge leverage score (Alaoui & Mahoney, 2015)). *Let $E \in \mathbb{R}^{n \times n}$ be a kernel matrix and let $\lambda > 0$. The λ -ridge leverage score of the data point x_i is defined as*

$$\tau_i^\lambda := (E(E + \lambda I)^{-1})_{i,i}.$$

If $E = BB^\top$ for some $B \in \mathbb{R}^{n \times n}$, then this can be equivalently written as

$$\tau_i^\lambda = b_i^\top (B^\top B + \lambda I)^{-1} b_i,$$

where b_i is the i -th row of B .

It is easy to see that $\sum_{i=1}^n \tau_i^\lambda = s_\lambda$. Moreover, Musco & Musco (2017) shows that Nyström approximations (Williams & Seeger, 2000) based on ridge leverage score sampling yield accurate spectral approximations to E .

Lemma 2.4 (Theorem 3 of Musco & Musco (2017)). *Let $s = O(s_\lambda \log(s_\lambda/\delta))$, $\lambda > 0$, and $\delta \in (0, 1)$. Let $E \in \mathbb{R}^{n \times n}$ be any kernel matrix. Let $S \in \mathbb{R}^{n \times s}$ be the λ -ridge leverage score sampling matrix. Then the Nyström approximation $\tilde{E} := ES(S^\top ES)^\dagger S^\top E$ satisfies $E \preceq \tilde{E} \preceq E + \lambda I$ with probability at least $1 - \delta$.*

Quantum Primitives. In this paper, we primarily leverage two quantum algorithmic primitives. The first is an efficient quantum sampling oracle based on Grover search.

Lemma 2.5 (Claim 3 in Apers & De Wolf (2022)). *Let n be a positive integer, and let $\{p_1, \dots, p_n\} \subseteq [0, 1]$ be a list of probabilities. There exists a quantum algorithm, $\text{QSAMPLE}(p)$, that generates a list of indices where each i is sampled independently with probability p_i , in time $\tilde{O}\left(\sqrt{n \sum_{i=1}^n p_i}\right) \cdot \mathcal{T}$, where \mathcal{T} denotes the time required to generate any individual p_i .*

The second primitive is a quantum procedure for approximating matrix-vector products using quantum multivariate mean estimation.

Lemma 2.6 (Theorem 5.1 of Apers & Gribling (2023)). *Let $\epsilon \in (0, 1)$, and let $A \in \mathbb{R}^{n \times d}$ and $v \in \mathbb{R}^n$. Suppose we are given quantum query access to the rows of A and the entries of v . Then there exists a quantum algorithm $\text{QMATVEC}(A, v, \epsilon)$ that outputs a vector $\tilde{\mu} \in \mathbb{R}^d$ such that, with probability at least $1 - 1/\text{poly}(n)$, $\|\tilde{\mu} - A^\top v\|_{(A^\top A)^{-1}} \leq \epsilon$, using $\tilde{O}\left(\epsilon^{-1} n^{0.5} d^{0.5} \|v\|_\infty\right)$ queries to A and v .*

3 TECHNICAL OVERVIEW

In this section, we provide an overview on the algorithmic techniques we utilize to approximate A , D and V , in sublinear time.

3.1 APPROXIMATE THE ATTENTION MATRIX VIA QUANTUM NYSTRÖM

To approximate the attention matrix A , we will make use of Nyström approximation (Williams & Seeger, 2000). However, recall that $A = \exp(QK^\top)$; for $Q \neq K$, the matrix itself is not even symmetric. This poses significant challenges for obtaining a good approximation. On the other hand, if we treat the queries and keys as the *dataset*, and form the exponential kernel matrix over them, then the resulting matrix is indeed a kernel matrix.

Specifically, let the dataset $X = \{q_1, \dots, q_n, k_1, \dots, k_n\}$, and consider $E \in \mathbb{R}^{2n \times 2n}$ where $E = \begin{bmatrix} \exp(QQ^\top) & \exp(QK^\top) \\ \exp(KQ^\top) & \exp(KK^\top) \end{bmatrix}$, then the attention matrix can be retrieved via $PE \begin{bmatrix} \mathbf{0}_n \\ \mathbf{1}_n \end{bmatrix}$ where $P \in \mathbb{R}^{n \times 2n}$ is the matrix consisting of the first n rows of the $2n \times 2n$ identity matrix, which selects the first n rows of E . Thus, once we obtain an approximation for E , we automatically obtain an approximation for A .

It remains to compute a Nyström approximation of E , as at first glance it is not clear how to even generate the ridge leverage score sampling matrix S in sublinear time. Musco & Musco (2017) shows that on a classical computer, it is possible to compute a *generalized* ridge leverage score sampling matrix using $\tilde{O}(ns_\lambda)$ evaluations of the kernel function and an additional $\tilde{O}(ns_\lambda^2)$ time, via a recursive sampling scheme:

- Uniformly sample half of the data points, then recursively compute the weighted sampling matrix $\tilde{S}^{n \times s}$ for the subset;
- Compute the *generalized ridge leverage score*, defined as $\tilde{\tau}_i^\lambda := b_i^\top (B^\top \tilde{S} \tilde{S}^\top B + \lambda I)^\dagger b_i$, and set $p_i = \min\{1, \tilde{\tau}_i^\lambda \cdot \log(s_\lambda/\delta)\}$;
- Output S as the weighted sampling matrix according to p_i .

The key ingredients in their algorithm are (1) the generalized ridge leverage score can be computed via kernel function evaluations instead of computing the factorization (see Definition A.4), and (2) sampling according to generalized ridge leverage score only increases the sample size by a constant factor, hence it does not affect the asymptotic runtime of the algorithm (see Lemma A.3).

For the simpler setting of leverage score sampling, Apers & Gribling (2023) shows that this recursive framework can benefit from quantum speedup, especially the Grover search sampler of Lemma 2.5, by noting that when sampling according to the leverage score, it is not necessary to compute or approximate all the scores; rather, it is enough to implement an oracle that can supply any approximate leverage score when needed.

For our application, however, this oracle is much more difficult to implement, as in the setting of Apers & Gribling (2023), one could directly query the row of B , which is not the case for the kernel setting. Nevertheless, we show how to implement such an oracle for generalized ridge leverage scores of kernels. The algorithm is detailed in Algorithm 1. Throughout this section, we let s denote the final sample size of the Nyström approximation.

The main idea is to utilize the identity $\tilde{\tau}_i^\lambda = \frac{1}{\lambda} (E - ES(S^\top ES + \lambda I)^{-1} S^\top E)_{i,i}$, where $E_{i,i}$ involves a single kernel evaluation $\mathsf{K}(x_i, x_i)$, and $S^\top ES$ requires only $O(s^2)$ kernel evaluations. Finally, the term $(ES(S^\top ES + \lambda I)^\dagger S^\top E)_{i,i}$ can be computed by evaluating the kernel between x_i and the sampled points in S , weighted appropriately, which requires $O(s)$ kernel evaluations. This shows that we can implement the oracle by precomputing $(S^\top ES + \lambda I)^\dagger$ in $O(s^2) \cdot \mathcal{T}_K + s^\omega$ time, where \mathcal{T}_K denotes the time for kernel evaluation and $\omega \approx 2.37$ is the matrix multiplication exponent (Duan et al., 2023; Williams et al., 2024; Alman et al., 2025). Each oracle query can then be answered in $O(s) \cdot \mathcal{T}_K + s^2$ time. By Lemma 2.5, the quantum sampler requires only $\tilde{O}(n^{0.5} s^{0.5})$ oracle calls, so the overall runtime is $\tilde{O}(n^{0.5} s^{1.5} \cdot (\mathcal{T}_K + s) + s^\omega)$. In our setting, the

Algorithm 1 Quantum Nyström approximation via recursive generalized ridge leverage score sampling.

```

1: procedure QNYSTRÖMKERNEL( $\{x_1, \dots, x_n\} \in (\mathbb{R}^d)^n, K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^m, \delta \in (0, 1), \lambda \in$ 
   ( $0, \infty$ ))
    $\triangleright \delta$  is the failure probability,  $\lambda$  is the ridge leverage score parameter.
2:    $s \leftarrow O(s_\lambda \log(s_\lambda/\delta))$ 
3:    $T \leftarrow O(\log(n/s))$ 
4:   Let  $S_0 \subset_{1/2} S_1 \subset_{1/2} \dots \subset_{1/2} S_T = [n]$   $\triangleright$  We use  $A \subset_{1/2} B$  to denote  $A$  is a uniform
   subset of half of the indices of  $B$ 
5:    $M_0 \leftarrow \{K(x_i, x_j)\}_{(i,j) \in S_0 \times S_0}$   $\triangleright |S_0| = s$ 
6:   Let  $D_0 \in \mathbb{R}^{n \times |S_0|}$  be the sampling matrix of  $S_0$ 
7:   for  $t = 1$  to  $T$  do
8:      $\tilde{M} \leftarrow (M_{t-1} + \lambda I_s)^{-1}$ 
9:      $\triangleright$  Let  $D_{t-1}^\top K_i := \{D_{t-1}(j) \cdot K(x_i, x_j)\}_{j \in D_{t-1}} \in \mathbb{R}^s$  for  $i \in S_t$  where  $D_{t-1}(j)$  is the
   weight corresponding to  $x_j$  specified by  $D_{t-1}$ 
10:    Implement oracle for  $q_i \leftarrow \frac{5}{\lambda} \cdot (K(x_i, x_i) - (D_{t-1}^\top K_i)^\top \tilde{M} D_{t-1}^\top K_i)$  for  $i \in S_t$ 
11:     $\triangleright p_i = \min\{1, 16q_i \log(2s/\delta)\}$ 
12:     $\tilde{D}_t \leftarrow \text{QSAMPLE}(p)$   $\triangleright \tilde{D}_t \in \mathbb{R}^{|S_t| \times s}$ 
13:     $D_t \leftarrow D_{S_t} \cdot \tilde{D}_t$   $\triangleright D_t \in \mathbb{R}^{n \times s}$ 
14:     $M_t \leftarrow \{D_t(i) D_t(j) \cdot K(x_i, x_j)\}_{(i,j) \in D_t \times D_t}$   $\triangleright M_t \in \mathbb{R}^{s \times s}$ 
15:  end for
16:  return  $D_T$ 
17: end procedure

```

kernel function $K(x_i, x_j) = \exp(\langle x_i, x_j \rangle)$ can be computed in $O(d)$ time, which gives a runtime of $\tilde{O}(n^{0.5} s^{1.5} (d + s) + s^\omega)$, sublinear in n .

It remains to analyze the approximation guarantee. Sampling according to generalized ridge leverage scores ensures that $E \preceq \tilde{E} \preceq E + \lambda I$, but this does not immediately imply a bound on the approximation error for $\exp(QK^\top)$. To address this, let $E = \begin{bmatrix} B & A \\ A^\top & C \end{bmatrix}$ and $\tilde{E} = \begin{bmatrix} \tilde{B} & \tilde{A} \\ \tilde{A}^\top & \tilde{C} \end{bmatrix}$.

Standard spectral approximation theory guarantees that $B \preceq \tilde{B} \preceq B + \lambda I$ and $C \preceq \tilde{C} \preceq C + \lambda I$. For the off-diagonal block we are interested in A , we cannot get such a strong spectral approximation guarantee; in fact, one can show that the best we could hope for is a symmetrization bound: $A + A^\top \preceq \tilde{A} + \tilde{A}^\top \preceq A + A^\top + 2\lambda I$. On the other hand, a weaker and a more handy bound can be exhibited: $\|A - \tilde{A}\| \leq \lambda$ and $\|A - \tilde{A}\|_F \leq \lambda\sqrt{n}$, and we will show these bounds are sufficient to derive the final approximation guarantees of our algorithm.

It is also worth noting that Algorithm 1 merely computes the weighted sampling matrix S , which can be stored compactly by recording the sampled indices and corresponding weights, but does not explicitly form the Nyström approximation $\tilde{E} = ES(S^\top ES)^\dagger S^\top E$. While $(S^\top ES)^\dagger$ can be computed and stored in $O(s^2 d + s^\omega)$ time, forming \tilde{E} would take $\Omega(ns)$ time, which is prohibitive due to output size. In what follows, we show that this restricted representation of S is nonetheless sufficient to approximate D, V , and $\text{Att}(Q, K, V)$.

We now compare our Nyström approximation scheme to a related method known as Nyströmformer (Xiong et al., 2021), which also integrates Nyström into the attention mechanism. Specifically,

they consider the attention matrix A and partition it as $A = \begin{bmatrix} X_1 & X_2 \\ X_3 & X_4 \end{bmatrix}$, aiming to approximate X_4

using the other three blocks. Given Nyström landmark points Q' and K' sampled from Q and K , they set $X_1 = \exp(Q'K'^\top)$, $X_2 = \exp(QK'^\top)$, and $X_3 = \exp(Q'K^\top)$. Since the number of landmarks is small, these blocks are all low-dimensional. Xiong et al. (2021) proves that X_4 can be efficiently approximated using X_1, X_2 , and X_3 in $O(nmd)$ time, where m is the number of landmarks. While Nyströmformer performs well in practice, it guarantees convergence to the true attention matrix only when all rows of Q and K are included as landmarks. In contrast, our Nyström scheme operates on

the exponential kernel matrix formed from Q and K , and achieves spectral approximation guarantees as long as the sample size is sufficiently large without needing to include all data points.

3.2 APPROXIMATE THE NORMALIZATION FACTOR VIA QUANTUM MEAN ESTIMATION

Recall that $D = \text{diag}(A\mathbf{1}_n)$, and each normalization factor only requires computing $a_i^\top \mathbf{1}_n$, where a_i is the i -th row of A . If we have access to \tilde{E} , then the i -th normalization factor could be estimated as $\tilde{E}_{i,*}^\top \begin{bmatrix} \mathbf{0}_n \\ \mathbf{1}_n \end{bmatrix}$. However, as discussed earlier, we cannot explicitly form \tilde{E} due to its size. To resolve this, we define $U := ES(S^\top ES)^\dagger/2 \in \mathbb{R}^{2n \times s}$. By the definition of the Nyström approximation, we have $\tilde{E} = UU^\top$. Moreover, U also exhibits a block structure $U = \begin{bmatrix} U_1 \\ U_2 \end{bmatrix}$ where $U_1, U_2 \in \mathbb{R}^{n \times s}$, and the desired approximate $\tilde{A} = U_1 U_2^\top$ can be obtained via these blocks. Given any vector $v \in \mathbb{R}^n$, if we can compute or approximate $U_2^\top v$, then the normalization factor for the i -th row can be estimated as $(U_1)_{i,*}^\top (U_2^\top v)$ where $(U_1)_{i,*} \in \mathbb{R}^s$ is the i -th row of U_1 . Fortunately, we can implement row queries to U_2 . We first precompute $(S^\top ES)^\dagger/2$ in $O(s^2 d + s^\omega)$ time, then each row $(U_2)_{i,*}$ of U_2 is computed via kernel evaluations between x_{i+n} and the points in S , followed by matrix-vector multiplication with $(S^\top ES)^\dagger/2$. This takes $O(s^2 + sd)$ time.

It remains to approximate $U_2^\top v$, which we cast as a multivariate mean estimation problem. Define the random variable $X = 2nv_i(U_2)_{*,i}$, where $i \in [n]$ is selected uniformly at random. It is easy to verify that $\mathbb{E}[X] = U_2^\top v$, and the variance is bounded. Therefore, one can apply the quantum multivariate mean estimation procedure of [Cornelissen et al. \(2022\)](#) to approximate $U_2^\top v$. To further reduce variance, [Apers & Gribling \(2023\)](#) proposes approximating the matrix-vector product in the $(U_2^\top U_2)^{-1}$ -energy norm. Following this idea, we apply Lemma 2.6 to output a vector $\tilde{\mu} \in \mathbb{R}^s$ such that $\|\tilde{\mu} - U_2^\top v\|_{(U_2^\top U_2)^{-1}} \leq \epsilon$, using $\tilde{O}(\epsilon^{-1} n^{0.5} s^{0.5} \|v\|_\infty)$ row queries to U_2 and v . In our application, we always have $\|v\|_\infty = 1$, and as noted above, each row query to U takes $O(s^2 + sd)$ time. We present the full algorithm below in Algorithm 2.

For the approximation guarantee, we prove that for any vector $x \in \mathbb{R}^s$, if $\|x\|_{(U_2^\top U_2)^{-1}} \leq \epsilon$, then $\|U_1 x\|_2 \leq \epsilon \cdot \|U_1 U_2^\top\|$. This is particularly useful for us, as we can set $x = U_2^\top v - \tilde{\mu}$, in which case $U_1 x = U_1 U_2^\top v - U_1 \tilde{\mu} = \tilde{A} v - U_1 \tilde{\mu}$, and the upper bound becomes $\epsilon \cdot \|\tilde{A}\| \leq \epsilon \cdot (\|A\| + \lambda)$. On the other hand, we can upper bound $\|(\tilde{A} - A)v\|_\infty$ using the matrix infinity norm, defined as $\|\tilde{A} - A\|_\infty = \max_{i \in [n]} \|\tilde{A}_{i,*} - A_{i,*}\|_1$. A simple argument shows that $\|\tilde{A} - A\|_\infty \leq \sqrt{n} \cdot \|\tilde{A} - A\| \leq \lambda \sqrt{n}$. A triangle inequality then yields the final approximation guarantee. If we define $\tilde{D} := \text{diag}(\tilde{A}\mathbf{1}_n)$, the above analysis provides a bound on $\|D - \tilde{D}\|$. However, in forming the attention module, it is more desirable to control $\|\tilde{D}^{-1}\|$. To achieve this, we prove a perturbation bound on matrix inversion that relates $\|\tilde{D}^{-1}\|$ to $\|D^{-1}\|$.

3.3 APPROXIMATE THE VALUE MATRIX VIA LEVERAGE SCORE SAMPLING

In preceding discussions, we have shown how to construct the sampling matrix for Nyström approximation and how to compute the normalization factor for any row $i \in [n]$. It remains to approximate V in sublinear time. Prior classical algorithms, such as [Zandieh et al. \(2023\)](#), propose using importance sampling based on the *joint row norm* of V and $D^{-1}A$. Specifically, the sampling probability for the i -th row is set as $p_i \geq 1/4 \cdot (\|e_i^\top D^{-1}A\|_2^2 + \gamma \cdot \|v_i\|_2^2) / (\|D^{-1}A\|_F^2 + \gamma \cdot \|V\|_F^2)$, where $\gamma = \|D^{-1}A\|_2^2 / \|V\|_2^2$. This method achieves a final sample size that is nearly linear in $d + \text{srnk}(D^{-1}A)$, where $\text{srnk}(D^{-1}A) = \|D^{-1}A\|_F^2 / \|D^{-1}A\|_2^2$ is the stable rank of the softmax matrix. While this approach is conceptually simple and easy to implement, it requires estimating the Frobenius norms of both V and $D^{-1}A$ to constant-factor accuracy. This is straightforward if we are allowed to read all entries of V , but becomes particularly challenging in sublinear time. Our solution is to instead use leverage score sampling on the matrix V , which can be implemented in sublinear time ([Apers & Gribling, 2023](#)).

Unlike the joint sampling distribution of [Zandieh et al. \(2023\)](#), which yields a *spectral norm approximate matrix multiplication* guarantee of the form $\|D^{-1}ASS^\top V\| \leq \epsilon \cdot \|D^{-1}A\| \cdot \|V\|$, leverage

Algorithm 2 Algorithm for estimating normalization factor.

```

1: data structure QROWNORM
2: begin members
3:    $s \in \mathbb{N}$ 
4:    $S \in (\mathbb{R}^2)^s$ 
5:    $N \in \mathbb{R}^{s \times s}$ 
6:    $\tilde{\mu} \in \mathbb{R}^s$ 
7: end members
8:
9: procedure PREPROCESS( $Q \in \mathbb{R}^{n \times d}, K \in \mathbb{R}^{n \times d}, \lambda \in (0, \infty), \epsilon \in (0, 1)$ )
10:    $s \leftarrow O(s_\lambda \log(s_\lambda n))$ 
11:    $S \leftarrow \text{QNYSTRÖMKERNEL}(Q \cup K, (x_i, x_j) \mapsto \exp(\langle x_i, x_j \rangle), 1/\text{poly}(n), \lambda)$   $\triangleright$ 
     Algorithm 1,  $S$  is a list of sampled indices and weights
12:    $N \leftarrow (S^\top E S)^\dagger / 2$ 
13:   Implement row oracle  $(U_2)_{j,*}$  as follows:
14:    $(\tilde{U}_2)_{j(k),*} \leftarrow S_k \cdot \exp(\langle x_{j+n}, x_k \rangle), \forall k \in S$   $\triangleright (\tilde{U}_2)_{j(k),*} \in \mathbb{R}^s$ 
15:    $(U_2)_{j,*} \leftarrow N \tilde{U}_2$   $\triangleright S$  stores pairs of indices and weights,  $S_k$  is the weight
     corresponding to index  $k$ ,  $(U_2)_{j,*} \in \mathbb{R}^s$ 
16:   Implement entry oracle for a vector  $v = \mathbf{1}_n \in \mathbb{R}^n$ 
17:    $\tilde{\mu} \leftarrow \text{QMATVEC}(U_2, v, \epsilon)$   $\triangleright \tilde{\mu} \in \mathbb{R}^s$ , Lemma 2.6
18: end procedure
19:
20: procedure QUERY( $i \in [n]$ )
21:    $b_i \leftarrow \langle (U_1)_{i,*}, \tilde{\mu} \rangle$   $\triangleright (U_1)_{i,*}$  is computed via row oracle
22:   return  $b_i$ 
23: end procedure
24: end data structure

```

score sampling has two key limitations: (1) it requires that V have orthonormal columns (Clarkson & Woodruff, 2017), and (2) it provides approximate matrix multiplication guarantees in Frobenius norm, i.e., $\|D^{-1} A S S^\top V\|_F \leq \epsilon \cdot \|D^{-1} A\|_F \cdot \|V\|_F$.

To address the first limitation, we introduce a new parameter called the *row distortion* of V , defined as $\alpha := d/\|V\|_F^2 \cdot \max_{i \in [n]} \|v_i\|_2^2 / \tau_i$. Intuitively, α measures the mismatch between the row density and row importance. Specifically, the ratio $\|v_i\|_2^2 / \|V\|_F^2$ quantifies how much row v_i contributes in ℓ_2^2 norm, while τ_i/d measures how linearly independent v_i is compared to other rows via τ_i .

Our main result is that by sampling $\tilde{O}(\epsilon^{-2}\alpha)$ rows of V according to its leverage score distribution, we obtain an approximate matrix multiplication guarantee in Frobenius norm. Note that $\alpha = 1$ if V has orthonormal columns, which recovers the result of Clarkson & Woodruff (2017). This sampling procedure can be implemented in $\tilde{O}(\epsilon^{-1} n^{0.5} \alpha^{0.5} d)$ time by making row queries to V .

3.4 MAIN RESULT

Now that we have described how to approximate each of the matrices D , A , and V , we are in a position to state our main result. We provide an overview of our algorithm below in Algorithm 3.

Theorem 3.1 (Informal version of Theorem D.2). *Let $Q, K, V \in \mathbb{R}^{n \times d}$ be the query, key and value matrices, let $\epsilon, \lambda > 0$. Let $E \in \mathbb{R}^{2n \times 2n}$ be the exponential kernel matrix on the dataset $Q \cup K$ and s_λ be the statistical dimension of E (Definition 2.2) and α be the row distortion of V (Definition C.2). Assume that $\|D^{-1}\| < \frac{1}{\epsilon \|A\| + \lambda \sqrt{n}}$ and let $\beta = \frac{1}{1 - (\epsilon \|A\| + \lambda \sqrt{n}) \|D^{-1}\|}$. There exists a quantum data structure that preprocesses Q, K, V through only row queries to these matrices and maintains matrices $\tilde{D}, \tilde{A}, \tilde{V}$ implicitly such that, with probability at least $1 - 1/\text{poly}(n)$,*

$$\|\tilde{D}^{-1} \tilde{A} \tilde{V} - \text{Att}(Q, K, V)\|_F \leq \epsilon \cdot (\beta \cdot \|D^{-1}\|) \cdot (\|A\|_F + \lambda \sqrt{n}) \cdot \|V\|_F.$$

Moreover, the data structure has the specification

Algorithm 3 Quantum data structure for attention row query.

```

1: data structure QATTENTION ▷ Theorem 3.1
2: begin members
3:    $s_E, s_V \in \mathbb{N}$ 
4:    $\tilde{V} \in \mathbb{R}^{s_V \times d}$ 
5:    $\tilde{N} \in \mathbb{R}^{s_E \times s_V}$ 
6:    $\tilde{L} \in \mathbb{R}^{s_E \times d}$ 
7:   QROWNORM QRN ▷ Algorithm 2
8: end members
9:
10: procedure PREPROCESS( $Q \in \mathbb{R}^{n \times d}, K \in \mathbb{R}^{n \times d}, V \in \mathbb{R}^{n \times d}, \lambda > 0, \epsilon > 0, \alpha \geq 1$ )
11:    $s_\lambda \leftarrow s_\lambda(E)$ 
12:    $s_V \leftarrow \tilde{O}(\epsilon^{-2}\alpha), s_E \leftarrow \tilde{O}(s_\lambda)$ 
13:   QRN.PREPROCESS( $Q, K, \lambda, \epsilon$ ) ▷ Algorithm 2
14:    $\tilde{S}_V \leftarrow \text{QLEVERAGE SCORE}(V, s_V)$  ▷  $\tilde{S}_V \in \mathbb{R}^{n \times s_V}$ , Lemma A.5
15:    $\tilde{V} \leftarrow \tilde{S}_V^\top V$  ▷  $\tilde{V} \in \mathbb{R}^{s_V \times d}$ 
16:    $S_E \leftarrow \text{QNYSTRÖM KERNEL}(Q \cup K, (x_i, x_j) \mapsto \exp(\langle x_i, x_j \rangle), 1/\text{poly}(n), \lambda)$ 
17:   ▷ Let  $x_1, \dots, x_{2n}$  denote the dataset  $Q \cup K$ 
18:    $\tilde{M} \leftarrow \{S_E(i)S_E(j) \cdot \exp(\langle x_i, x_j \rangle)\}_{(i,j) \in S_E \times S_E}$  ▷  $\tilde{M} \in \mathbb{R}^{s_E \times s_E}$ 
19:    $\tilde{R} \leftarrow \{S_E(i)\tilde{S}_V(j) \cdot \exp(\langle x_i, x_j \rangle)\}_{(i,j) \in S_E \times S_V}$  ▷  $\tilde{R} \in \mathbb{R}^{s_E \times s_V}, \tilde{R} = S_E^\top E \tilde{S}_V$ 
20:    $\tilde{N} \leftarrow \tilde{M}^\dagger \tilde{R}$  ▷  $\tilde{N} \in \mathbb{R}^{s_E \times s_V}$ 
21:    $\tilde{L} \leftarrow \tilde{N} \tilde{V}$  ▷  $\tilde{L} \in \mathbb{R}^{s_E \times d}$ 
22: end procedure
23:
24: procedure QUERY( $i \in [n]$ )
25:    $b_i \leftarrow \text{QRN.QUERY}(i)$  ▷ Algorithm 2
26:    $u_i \leftarrow \{\tilde{S}_E(j) \cdot \exp(\langle x_i, x_j \rangle)\}_{j \in S_E}$  ▷  $u_i \in \mathbb{R}^{s_E}$ 
27:   return  $\tilde{L}^\top u_i / b_i$ 
28: end procedure
29: end data structure

```

- It preprocesses Q, K, V in $\tilde{O}(\epsilon^{-1}n^{0.5}s_\lambda^{0.5})$ row queries to Q, K and $\tilde{O}(\epsilon^{-1}n^{0.5}\alpha^{0.5})$ row queries to V , and $\tilde{O}(\epsilon^{-1}n^{0.5}(s_\lambda^{2.5} + s_\lambda^{1.5}d + \alpha^{0.5}d))$ time;
- For any $i \in [n]$, it returns a vector $\tilde{r}_i = e_i^\top \tilde{D}^{-1} \tilde{A} \tilde{V}$ in $\tilde{O}(s_\lambda^2 + s_\lambda d)$ time.

We pause to make some remarks on Theorem 3.1. The preprocessing time scales with $n^{0.5}$, achieving a quadratic speedup with respect to n over any classical algorithm. Several parameters merit further discussion, in particular the statistical dimension s_λ and the approximation factor for $\|D^{-1}\|$, denoted by β . We summarize their relationships as functions of λ in Table 1. The row distortion factor α also affects the runtime, and in Appendix C, we prove that $\alpha \leq \frac{d}{\text{srank}(V)}$ where $\text{srank}(V) = \frac{\|V\|_F^2}{\|V\|^2}$ is the stable rank of V . This ensures $\alpha \leq d$ and becomes smaller if the value matrix V has close to d stable rank. We empirically verify that (1) the assumption on $\|D^{-1}\|$ is easy to satisfy with wide range of choices for ϵ , (2) the Frobenius norm of A is only a small constant factor of its spectral norm, (3) the row distortion $\alpha = O(1)$ and (4) the infinity norm of A is only a small constant factor of its spectral norm, implying in practice, the additive $\lambda\sqrt{n}$ term is likely to be $O(\lambda)$. We refer to Appendix E for a more detailed section.

λ	s_λ	$\frac{1}{\epsilon\ A\ + \lambda\sqrt{n}}$	β
↑	↓	↓	↑
↓	↑	↑	↓

Table 1: Parameters s_λ , $\frac{1}{\epsilon\|E\| + \lambda\sqrt{n}}$, and β as functions of λ .

Bit Complexity of Our Algorithm. Our discussions and results above are grounded in the assumption that arithmetic operations are performed in infinite precision, while this is usually adopted in the analysis of classical algorithms, QRAM model only allows $O(\log n)$ qubits and $\text{poly}(n)$ bits. In Section F, we provide a preliminary bit complexity analysis of our algorithm, in particular centering around the matrix inversion and pseudoinversion operations. To the best of our knowledge, there is no prior work on studying the bit complexity of numerical linear algebra operations in the QRAM model, and we leave a comprehensive analysis of bit complexity as a future direction.

4 RELATED WORK

Transformers and Attention Mechanism. Transformers (Vaswani et al., 2017) have been the driving force behind large language models (Devlin et al., 2019; Brown et al., 2020; Touvron et al., 2023; Bubeck et al., 2023; Team et al., 2023; Liu et al., 2024a). They are sequence-to-sequence generative models, where the sequence length is typically denoted by n . The key architectural component that distinguishes transformers from earlier models is the attention mechanism, which computes a softmax over the pairwise interactions of query-key vectors. However, computing the full softmax distribution requires $\Omega(n^2)$ time, due to the size of the attention matrix. This quadratic dependency renders transformers inefficient for long sequences, motivating a rich body of work aimed at approximating attention in subquadratic time. These approaches can be broadly categorized into three main classes: (1) *Pattern-based sparse attention*: only a subset of attention matrix entries are computed, with the subset determined by predefined patterns, such as sliding windows or graph-based sparsity structures (Daras et al., 2020; Kitaev et al., 2020; Roy et al., 2021; Sun et al., 2022; Child et al., 2019; Beltagy et al., 2020; Ainslie et al., 2020; Zaheer et al., 2020). (2) *Kernel-based linear attention*: these methods attempt to linearize the kernel by exploiting the identity $K(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle$ for a feature map $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$. When the kernel is exponential, exact computation requires $m = \infty$, so many heuristic approximations for ϕ have been proposed (Katharopoulos et al., 2020; Choromanski et al., 2021; Wang et al., 2020; Peng et al., 2021) with $m = O(d)$. (3) *Data structure-based attention*: these works design specialized data structures for approximating various components of attention. Examples include estimating the normalization factor via kernel density estimation (KDE) (Zandieh et al., 2023), using hashing to identify large entries (Han et al., 2024), applying polynomial approximation methods under bounded input conditions (Alman & Song, 2023), and other algorithmic innovations (Kacham et al., 2024; Zandieh et al., 2024; van den Brand et al., 2024; Song et al., 2024; Kannan et al., 2025; Chu et al., 2024; Chen et al., 2025b; Indyk et al., 2025). Our work falls into the third category, as we design quantum data structures to approximate each of the matrices involved in the attention computation.

Quantum Machine Learning. Given a machine learning problem, can we solve it faster on a quantum computer? The paradigm of using quantum mechanics to accelerate machine learning algorithms has sparked significant interest, leading to a wide array of results across diverse problem domains, including clustering (Kerenidis et al., 2019; Xue et al., 2023), classification (Li et al., 2019), regression (Chen & de Wolf, 2023), training neural networks (Chakrabarti et al., 2019; Kerenidis et al., 2020), convex optimization (Chakrabarti et al., 2020; van Apeldoorn et al., 2020a; Li & Zhang, 2022; Sidford & Zhang, 2023; Zhang et al., 2024b; Wang et al., 2024), mathematical programming (Brandão et al., 2019; van Apeldoorn et al., 2020b; van Apeldoorn & Gilyén, 2019; Kerenidis & Prakash, 2020; Kerenidis et al., 2021; van Apeldoorn et al., 2021; Apers & Gribling, 2023), graph sparsification (Apers & De Wolf, 2022), and recommender systems (Kerenidis & Prakash, 2017). Among the key quantum techniques, Grover search (Grover, 1996) plays a foundational role. It provides a quadratic speedup for database search problems: given a function $f : [n] \rightarrow \{0, 1\}$, the goal is to list up to m indices i such that $f(i) = 1$. The Grover search algorithm requires oracle access to f and can produce these m indices using only $O(\sqrt{mn})$ queries, in contrast to the $O(n)$ queries required classically. Several variants of Grover search have been developed to suit different computational settings. In this paper, we use the probabilistic version: given a list of n probabilities $p_1, \dots, p_n \in [0, 1]$, Grover search can be used to sample a list of indices where each i is selected independently with probability p_i . By the standard analysis of Grover search, this sampling requires $\tilde{O}(\sqrt{nP})$ queries to the probability values p_i where $P = \sum_{i=1}^n p_i$. Before our work, Gao et al. (2023) also applied Grover search to accelerate attention computation. However, their method requires a structural assumption: for each query $q_i \in \mathbb{R}^d$, the associated set $S_i = \{j \in [n] : \langle q_i, k_j \rangle \geq \tau\}$ must

have cardinality at most k . Under this assumption, their algorithm runs in time $\tilde{O}(n^{1.5}k^{0.5}d + nkd)$. Notably, if $k = n$, then their algorithm offers no speedup over the exact computation.

5 CONCLUSION

We consider the problem of approximating the attention module in the row query model, where the goal is to return individual rows of the approximate attention matrix. We design a quantum data structure that preprocesses Q , K , and V in $\tilde{O}(\epsilon^{-1}n^{0.5}\text{poly}(s_\lambda, d, \alpha))$ time, and answers any row query in $\tilde{O}(s_\lambda^2 + s_\lambda d)$ time. To the best of our knowledge, this is the first quantum algorithm to achieve sublinear dependence on n even in the row query model.

Our work also has several limitations, which raises interesting open questions. (1) The error guarantee we obtain is in Frobenius norm rather than spectral norm. While Frobenius norm bounds the sum of the squared ℓ_2 errors across all rows, the spectral norm provides a worst-case guarantee that each row is well approximated. Therefore, it would be desirable to strengthen the result to achieve a spectral norm guarantee. (2) The error bound we obtain contains an additive $\lambda\sqrt{n}$ term, which stems from bounding the infinity norm of the error matrix by \sqrt{n} times the spectral norm of it. This bound seems overly pessimistic, and it theoretically forces one to choose small value for λ , hindering the advantage of small statistical dimension. It would be interesting to remove the \sqrt{n} factor in the additive term. (3) While we provide a preliminary bit complexity analysis of our algorithm in Section F, we feel a more comprehensive study of bit complexity of numerical linear algebra in the QRAM model is needed. We leave this as a major future direction, as it will significantly broaden the practicality of these quantum speedups. (4) Our algorithm in its current form can only compute the full attention *without* the causal mask, as using the Nyström approximation implicitly assumes the complete interactions between queries and keys. To implement causal masking, one possibility is to design a quantum kernel density estimation data structure as shown in Zandieh et al. (2023) classically.

ETHICS STATEMENT

Our work is a theoretical quantum framework to approximate the attention module in sublinear time. We don't foresee any potential ethics concerns.

REPRODUCIBILITY STATEMENT

We include all the proofs in the appendix. For proofs of the exponential kernel, see Section A, for proofs of estimating the normalization factor, see Section B. For proofs of the leverage score approximate matrix multiplication, see Section C. The final conclusion is proved in Section D. We provide empirical verifications on the assumptions of the parameters in Section E, and a preliminary bit complexity analysis of our algorithm in Section F.

ACKNOWLEDGMENT

We would like to thank anonymous for very helpful discussions, and Ruizhe Zhang for answering our questions on the QRAM model. Lichen Zhang is supported by a Mathworks Fellowship and a Simons Dissertation Fellowship in Mathematics.

REFERENCES

- Muhammad Adnan, Akhil Arunkumar, Gaurav Jain, Prashant J. Nair, Ilya Soloveychik, and Purushotham Kamath. Keyformer: Kv cache reduction through key tokens selection for efficient generative inference. In *Proceedings of the 7th Conference on Machine Learning and Systems (MLSys)*, 2024.
- Joshua Ainslie, Santiago Ontanon, Chris Alberti, Vaclav Cvicek, Zachary Fisher, Philip Pham, Anirudh Ravula, Sumit Sanghai, Qifan Wang, and Li Yang. ETC: Encoding long and structured inputs in transformers. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 268–284. Association for Computational Linguistics, 2020.

- Ahmed El Alaoui and Michael W. Mahoney. Fast randomized kernel ridge regression with statistical guarantees. In *Advances in Neural Information Processing Systems 28 (NeurIPS 2015)*, pp. 775–783, 2015.
- Josh Alman and Zhao Song. Fast attention requires bounded entries. *Advances in Neural Information Processing Systems*, 36:63117–63135, 2023.
- Josh Alman, Ran Duan, Virginia Vassilevska Williams, Yinzhan Xu, Zixuan Xu, and Renfei Zhou. More asymmetry yields faster matrix multiplication. In *Proceedings of the 2025 Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pp. 2005–2039, 2025.
- Simon Apers and Ronald De Wolf. Quantum speedup for graph sparsification, cut approximation, and laplacian solving. *SIAM Journal on Computing*, 51(6):1703–1742, 2022.
- Simon Apers and Sander Gribling. Quantum speedups for linear programming via interior point methods. *arXiv preprint arXiv:2311.03215*, 2023.
- Payman Behnam, Yaosheng Fu, Ritchie Zhao, Po-An Tsai, Zhiding Yu, and Alexey Tumanov. Rocketkv: Accelerating long-context llm inference via two-stage kv cache compression. *arXiv preprint arXiv:2502.14051*, 2025.
- Iz Beltagy, Matthew E. Peters, and Arman Cohan. Longformer: The long-document transformer, 2020.
- Fernando G. S. L. Brandão, Amir Kalev, Tongyang Li, Cedric Yen-Yu Lin, Krysta M. Svore, and Xiaodi Wu. Quantum sdp solvers: Large speed-ups, optimality, and applications to quantum learning. In *46th International Colloquium on Automata, Languages, and Programming (ICALP 2019)*, volume 132 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pp. 27:1–27:14. Schloss Dagstuhl–Leibniz-Zentrum für Informatik, 2019.
- William Brandon, Mayank Mishra, Aniruddha Nrusimha, Rameswar Panda, and Jonathan Ragan-Kelley. Reducing transformer key-value cache size with cross-layer attention. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2024.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pp. 1877–1901. Curran Associates, Inc., 2020.
- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, Harsha Nori, Hamid Palangi, Marco Tulio Ribeiro, and Yi Zhang. Sparks of artificial general intelligence: Early experiments with gpt-4, 2023. URL <https://arxiv.org/abs/2303.12712>.
- Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 213–229. Springer, 2020.
- Shouvanik Chakrabarti, Yiming Huang, Tongyang Li, Soheil Feizi, and Xiaodi Wu. Quantum wasserstein generative adversarial networks. In *Advances in Neural Information Processing Systems 32 (NeurIPS 2019)*, pp. 6768–6779, 2019.
- Shouvanik Chakrabarti, Andrew M. Childs, Tongyang Li, and Xiaodi Wu. Quantum algorithms and lower bounds for convex optimization. *Quantum*, 4:221, 2020.
- Bo Chen, Xiaoyu Li, Yekun Ke, Yingyu Liang, Zhenmei Shi, and Zhao Song. Exploring the limits of kv cache compression in visual autoregressive transformers. *arXiv preprint arXiv:2503.14881*, 2025a.

- Bo Chen, Yingyu Liang, Zhizhou Sha, Zhenmei Shi, and Zhao Song. HSR-enhanced sparse attention acceleration. In *The Second Conference on Parsimony and Learning (Proceedings Track)*, 2025b. URL <https://openreview.net/forum?id=wsolgABiPZ>.
- Yanlin Chen and Ronald de Wolf. Quantum algorithms and lower bounds for linear regression with norm constraints. In *50th International Colloquium on Automata, Languages, and Programming (ICALP 2023)*, volume 261 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pp. 38:1–38:21. Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2023.
- Yifang Chen, Xiaoyu Li, Yingyu Liang, Zhenmei Shi, Zhao Song, and Yu Tian. Time and memory trade-off of kv-cache compression in tensor transformer decoding. *arXiv preprint arXiv:2503.11108*, 2025c.
- Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse transformers. *arXiv preprint arXiv:1904.10509*, 2019.
- Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, David Belanger, Lucy Colwell, and Adrian Weller. Rethinking attention with performers. In *International Conference on Learning Representations*, 2021.
- Jan K Chorowski, Dzmitry Bahdanau, Dmitriy Serdyuk, Kyunghyun Cho, and Yoshua Bengio. Attention-based models for speech recognition. *Advances in neural information processing systems*, 28, 2015.
- Timothy Chu, Josh Alman, Gary L Miller, Shyam Narayanan, Mark Sellke, and Zhao Song. Metric transforms and low rank representations of kernels for fast attention. *Advances in Neural Information Processing Systems*, 37:47014–47068, 2024.
- Kenneth L Clarkson and David P Woodruff. Low-rank approximation and regression in input sparsity time. *Journal of the ACM (JACM)*, 63(6):1–45, 2017.
- Arjan Cornelissen, Yassine Hamoudi, and Sofiene Jerbi. Near-optimal quantum algorithms for multivariate mean estimation. In *Proceedings of the 54th Annual ACM SIGACT Symposium on Theory of Computing*, STOC 2022, pp. 33–43, New York, NY, USA, 2022. Association for Computing Machinery.
- Giannis Daras, Nikita Kitaev, Augustus Odena, and Alexandros G. Dimakis. Smyrf: Efficient attention using asymmetric clustering. In *Advances in Neural Information Processing Systems*, volume 33, pp. 6470–6481, 2020.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, pp. 4171–4186. Association for Computational Linguistics, 2019.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021.
- Ran Duan, Hongxun Wu, and Renfei Zhou. Faster matrix multiplication via asymmetric hashing. In *FOCS*, 2023.
- Yuan Feng, Junlin Lv, Yukun Cao, Xike Xie, and S. Kevin Zhou. Identify critical kv cache in llm inference from an output perturbation perspective. *arXiv preprint arXiv:2502.03805*, 2025.
- Yeqi Gao, Zhao Song, Xin Yang, and Ruizhe Zhang. Fast quantum algorithm for attention computation, 2023.
- Lov K. Grover. A fast quantum mechanical algorithm for database search. In *Proceedings of the Twenty-Eighth Annual ACM Symposium on Theory of Computing*, STOC '96, pp. 212–219, New York, NY, USA, 1996. Association for Computing Machinery. ISBN 0897917855. doi: 10.1145/237814.237866. URL <https://doi.org/10.1145/237814.237866>.

- Yuzhou Gu, Zhao Song, Junze Yin, and Lichen Zhang. Low rank matrix completion via robust alternating minimization in nearly linear time. In *The Twelfth International Conference on Learning Representations (ICLR)*, 2024.
- Meng-Hao Guo, Tian-Xing Xu, Jiang-Jiang Liu, Zheng-Ning Liu, Peng-Tao Jiang, Tai-Jiang Mu, Song-Hai Zhang, Ralph R Martin, Ming-Ming Cheng, and Shi-Min Hu. Attention mechanisms in computer vision: A survey. *Computational Visual Media*, 8(3):331–368, 2022.
- Insu Han, Rajesh Jayaram, Amin Karbasi, Vahab Mirrokni, David P. Woodruff, and Amir Zandieh. Hyperattention: Long-context attention in near-linear time. In *International Conference on Learning Representations (ICLR)*, 2024.
- David Harvey and Joris van der Hoeven. Integer multiplication in time $O(n \log n)$. *Annals of Mathematics*, March 2021. doi: 10.4007/annals.2021.193.2.4. URL <https://hal.science/hal-02070778>.
- Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer, New York, 2nd edition, 2009.
- Nicholas J. Higham. *Accuracy and Stability of Numerical Algorithms*. Society for Industrial and Applied Mathematics, Philadelphia, PA, 2 edition, 2002. ISBN 0-89871-521-0. doi: 10.1137/1.9780898718027.
- Piotr Indyk, Michael Kapralov, Kshiteej Sheth, and Tal Wagner. Improved algorithms for kernel matrix-vector multiplication under sparsity assumptions. In *Proceedings of the International Conference on Learning Representations (ICLR)*, May 2025.
- Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu. TinyBERT: Distilling BERT for natural language understanding. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pp. 4163–4174. Association for Computational Linguistics, 2020.
- Praneeth Kacham, Vahab Mirrokni, and Peilin Zhong. Polysketchformer: Fast transformers via sketching polynomial kernels. In *Proceedings of the 41st International Conference on Machine Learning*, 2024.
- Ravindran Kannan, Chiranjib Bhattacharyya, Praneeth Kacham, and David P. Woodruff. Levattention: Time, space, and streaming efficient algorithm for heavy attentions. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2025.
- Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are rnns: Fast autoregressive transformers with linear attention. In *Proceedings of the 37th International Conference on Machine Learning*, pp. 5156–5165. PMLR, 2020.
- Iordanis Kerenidis and Anupam Prakash. Quantum recommendation systems. In *8th Innovations in Theoretical Computer Science Conference (ITCS 2017)*, volume 67 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pp. 49:1–49:21. Schloss Dagstuhl–Leibniz-Zentrum für Informatik, 2017.
- Iordanis Kerenidis and Anupam Prakash. A quantum interior point method for lps and sdps. *ACM Transactions on Quantum Computing*, 1(1):1–32, 2020.
- Iordanis Kerenidis, Jonas Landman, Alessandro Luongo, and Anupam Prakash. q-means: A quantum algorithm for unsupervised machine learning. In *Advances in Neural Information Processing Systems*, volume 32, pp. 4134–4144, 2019.
- Iordanis Kerenidis, Jonas Landman, and Anupam Prakash. Quantum algorithms for deep convolutional neural networks. In *Proceedings of the 8th International Conference on Learning Representations (ICLR)*, 2020.
- Iordanis Kerenidis, Anupam Prakash, and Dániel Szilágyi. Quantum algorithms for Second-Order Cone Programming and Support Vector Machines. *Quantum*, 5:427, 2021.

- Nikita Kitaev, Łukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. In *International Conference on Learning Representations*, 2020.
- Lilly Kumari, Shengjie Wang, Tianyi Zhou, Nikhil Sarda, Anthony Rowe, and Jeff Bilmes. Bumblebee: Dynamic kv-cache streaming submodular summarization for infinite-context transformers. In *Proceedings of the Conference on Learning for Molecules (COLM)*, 2024.
- Tongyang Li and Ruizhe Zhang. Quantum speedups of optimizing approximately convex functions with applications to logarithmic regret stochastic convex bandits. In *Advances in Neural Information Processing Systems*, volume 35, pp. 19565–19577, 2022.
- Tongyang Li, Shouvanik Chakrabarti, and Xiaodi Wu. Sublinear quantum algorithms for training linear and kernel-based classifiers. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 3815–3824. PMLR, June 2019.
- Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024a.
- Minghui Liu, Tahseen Rabbani, Tony O’Halloran, Ananth Sankaralingam, Mary-Anne Hartley, Brian Gravelle, Furong Huang, Cornelia Fermüller, and Yiannis Aloimonos. Hashevtic: A pre-attention kv cache eviction strategy using locality-sensitive hashing. *arXiv preprint arXiv:2412.16187*, 2024b.
- Xiaoyun Liu, Daniel Esser, Brandon Wagstaff, Anna Zavodni, Naomi Matsuura, Jonathan Kelly, and Eric Diller. Capsule robot pose and mechanism state detection in ultrasound using attention-based hierarchical deep learning. *Scientific Reports*, 12(1):21130, 2022.
- Cameron Musco and Christopher Musco. Recursive sampling for the nyström method. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS’17, pp. 3836–3848, Red Hook, NY, USA, 2017. Curran Associates Inc.
- Cameron Musco, Christopher Musco, and Aaron Sidford. Stability of the lanczos method for matrix function approximation. In *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms*, SODA ’18, pp. 1605–1624, USA, 2018. Society for Industrial and Applied Mathematics.
- Hao Peng, Nikolaos Pappas, Dani Yogatama, Roy Schwartz, Noah A. Smith, and Lingpeng Kong. Random feature attention. In *Proceedings of the 9th International Conference on Learning Representations (ICLR)*, 2021.
- Reiner Pope, Sholto Douglas, Aakanksha Chowdhery, Jacob Devlin, James Bradbury, Anselm Levskaya, Jonathan Heek, Kefan Xiao, Shivani Agrawal, and Jeff Dean. Efficiently scaling transformer inference. In *Proceedings of the 6th Conference on Machine Learning and Systems (MLSys)*, 2023.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21:1–67, 2020.
- Aurko Roy, Mohammad Saffar, Ashish Vaswani, and David Grangier. Efficient content-based sparse attention with routing transformers. *Transactions of the Association for Computational Linguistics*, 9:53–68, 2021.
- Rikhav Shah. Hermitian diagonalization in linear precision. In *Proceedings of the 2025 Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pp. 5599–5615, 2025.
- Aaron Sidford and Chenyi Zhang. Quantum speedups for stochastic optimization. In *Advances in Neural Information Processing Systems 36 (NeurIPS 2023)*, pp. 1–12, 2023.

- Aleksandros Sobczyk. Deterministic complexity analysis of hermitian eigenproblems. In Keren Censor-Hillel, Fabrizio Grandoni, Joël Ouaknine, and Gabriele Puppis (eds.), *52nd International Colloquium on Automata, Languages, and Programming, ICALP 2025, Aarhus, Denmark, July 8-11, 2025*, volume 334 of *LIPICs*, pp. 131:1–131:21. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2025.
- Zhao Song, Junze Yin, and Lichen Zhang. Solving attention kernel regression problem via preconditioner. In *International Conference on Artificial Intelligence and Statistics*, pp. 208–216. PMLR, 2024.
- Zhiqing Sun, Yiming Yang, and Shinjae Yoo. Sparse attention with learning to hash. In *International Conference on Learning Representations*, 2022.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- Joran van Apeldoorn and András Gilyén. Improvements in quantum sdp-solving with applications. In *46th International Colloquium on Automata, Languages, and Programming (ICALP 2019)*, volume 132 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pp. 99:1–99:15. Schloss Dagstuhl–Leibniz-Zentrum für Informatik, 2019.
- Joran van Apeldoorn, András Gilyén, Sander Gribling, and Ronald de Wolf. Convex optimization using quantum oracles. *Quantum*, 4:220, 2020a.
- Joran van Apeldoorn, András Gilyén, Sander Gribling, and Ronald de Wolf. Quantum sdp-solvers: Better upper and lower bounds. *Quantum*, 4:230, 2020b.
- Joran van Apeldoorn, Sander Gribling, Yinan Li, Harold Nieuwboer, Michael Walter, and Ronald de Wolf. Quantum algorithms for matrix scaling and matrix balancing. In *48th International Colloquium on Automata, Languages, and Programming (ICALP 2021)*, volume 198 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pp. 110:1–110:17. Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2021.
- Jan van den Brand, Zhao Song, and Tianyi Zhou. Algorithm and hardness for dynamic attention maintenance in large language models. In *International Conference on Machine Learning*, pp. 49008–49028. PMLR, 2024.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Evan Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, Nathan Lambert, Dustin Schwenk, Oyvind Tafjord, Taira Anderson, David Atkinson, Faeze Brahman, Christopher Clark, Pradeep Dasigi, Nouha Dziri, Allyson Ettinger, Michal Guerquin, David Heineman, Hamish Ivison, Pang Wei Koh, Jiacheng Liu, Saumya Malik, William Merrill, Lester James Validad Miranda, Jacob Morrison, Tyler Murray, Crystal Nam, Jake Poznanski, Valentina Pyatkin, Aman Rangapur, Michael Schmitz, Sam Skjonsberg, David Wadden, Christopher Wilhelm, Michael Wilson, Luke Zettlemoyer, Ali Farhadi, Noah A. Smith, and Hannaneh Hajishirzi. 2 OLMo 2 furious (COLM’s version). In *Second Conference on Language Modeling*, 2025.
- Hao Wang, Chenyi Zhang, and Tongyang Li. Near-optimal quantum algorithm for minimizing the maximal loss. In *Proceedings of the 12th International Conference on Learning Representations (ICLR)*, 2024.
- Sinong Wang, Belinda Z. Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-attention with linear complexity, 2020.

- Yuanchao Wang, Wenji Du, Chenghao Cai, and Yanyan Xu. Explaining the attention mechanism of end-to-end speech recognition using decision trees. *arXiv preprint arXiv:2110.03879*, 2021.
- Christopher K. I. Williams and Matthias Seeger. Using the nystrom method to speed up kernel machines. In *Proceedings of the 14th International Conference on Neural Information Processing Systems*, NIPS'00, pp. 661–667, Cambridge, MA, USA, 2000. MIT Press.
- Virginia Vassilevska Williams, Yinzhao Xu, Zixuan Xu, and Renfei Zhou. New bounds for matrix multiplication: from alpha to omega. In *Proceedings of the 2024 Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pp. 3792–3835. SIAM, 2024.
- Yunyang Xiong, Zhanpeng Zeng, Rudrasis Chakraborty, Mingxing Tan, Glenn Fung, Yin Li, and Vikas Singh. Nyströmformer: A nyström-based algorithm for approximating self-attention. In *AAAI*, pp. 14138–14148. AAAI Press, 2021.
- Yecheng Xue, Xiaoyu Chen, Tongyang Li, and Shaofeng H.-C. Jiang. Near-optimal quantum coresets construction algorithms for clustering. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 38881–38912. PMLR, 2023.
- Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Ruslan Salakhutdinov, and Quoc V. Le. Xlnet: Generalized autoregressive pretraining for language understanding. In *Advances in Neural Information Processing Systems*, volume 32, pp. 5754–5764. Curran Associates, Inc., 2019.
- Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontañón, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, and Amr Ahmed. Big bird: Transformers for longer sequences. In *Advances in Neural Information Processing Systems*, volume 33, pp. 17283–17297, 2020.
- Amir Zandieh, Insu Han, Majid Daliri, and Amin Karbasi. Kdeformer: Accelerating transformers via kernel density estimation. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 40605–40623. PMLR, 2023.
- Amir Zandieh, Insu Han, Vahab Mirrokni, and Amin Karbasi. Subgen: Token generation in sublinear time and memory. *arXiv preprint arXiv:2402.06082*, 2024.
- Rongzhi Zhang, Kuang Wang, Liyuan Liu, Shuohang Wang, Hao Cheng, Chao Zhang, and Yelong Shen. Lorc: Low-rank compression for llms kv cache with a progressive compression strategy. In *NeurIPS 2024 Workshop on Model Compression*, 2024a.
- Tong Zhang. Learning bounds for kernel regression using effective data dimensionality. *Neural Computation*, 17(9):2077–2098, 2005.
- Yexin Zhang, Chenyi Zhang, Cong Fang, Liwei Wang, and Tongyang Li. Quantum algorithms and lower bounds for finite-sum optimization. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 12345–12356. PMLR, 2024b.
- Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang. Informer: Beyond efficient transformer for long sequence time-series forecasting. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 11106–11115. AAAI Press, 2021.

Appendix

Roadmap. In Section A, we describe the quantum algorithm for exponential kernels. In Section B, we discuss how to estimate the normalization factor. In Section C, we show the details on approximating matrix multiplication via leverage scores. In Section D, we combine things together and obtain the main result. In Section E, we empirically verify the assumptions on the parameters. In Section F, we discuss the bit complexity of our algorithm.

A QUANTUM ALGORITHM FOR EXPONENTIAL KERNEL

In this section, we give a generic reduction from attention matrix to a kernel matrix. Given queries and keys $Q = \{q_1, \dots, q_n\}$, $K = \{k_1, \dots, k_n\}$, recall that we are interested in the matrix $\exp(QK^\top)$ where the (i, j) -th entry is $\exp(q_i^\top k_j)$, and this matrix is not a PSD kernel matrix. We show a reduction that first computes the exponential kernel $K(x, y) = \exp(\langle x, y \rangle)$ over the dataset $Q \cup K$, then we can effectively extract certain blocks of the kernel matrix E that approximates $\exp(QK^\top)$ well. We start with a lemma on block approximation.

Lemma A.1. *Let $E \in \mathbb{R}^{2n \times 2n}$ be a PSD matrix and $E = \begin{bmatrix} B & A \\ A^\top & C \end{bmatrix}$ where each block is of size $n \times n$. Suppose there exists a matrix $\tilde{E} \in \mathbb{R}^{2n \times 2n}$ such that $E \preceq \tilde{E} \preceq E + \lambda I$ for $\lambda > 0$ and let $\tilde{E} = \begin{bmatrix} \tilde{B} & \tilde{A} \\ \tilde{A}^\top & \tilde{C} \end{bmatrix}$, then we have*

$$\|A - \tilde{A}\| \leq \lambda \text{ and } \|A - \tilde{A}\|_F \leq \lambda\sqrt{n}.$$

Proof. We let $v \in \mathbb{R}^n$ be the vector that realizes the spectral norm $A - \tilde{A}$, consider the augmented vector $\begin{bmatrix} \mathbf{0}_n \\ v \end{bmatrix}$, then we see that

$$\begin{aligned} \left\| (E - \tilde{E}) \begin{bmatrix} \mathbf{0}_n \\ v \end{bmatrix} \right\|_2^2 &= \left\| \begin{bmatrix} (A - \tilde{A})v \\ (C - \tilde{C})v \end{bmatrix} \right\|_2^2 \\ &= \|(A - \tilde{A})v\|_2^2 + \|(C - \tilde{C})v\|_2^2 \\ &\leq \lambda^2, \end{aligned}$$

where the last step is by $\|E - \tilde{E}\| \leq \lambda$ and our test vector is unit norm. As $\|(C - \tilde{C})v\|_2^2$ is trivially non-negative, we conclude that $\|(A - \tilde{A})v\|_2 = \|A - \tilde{A}\| \leq \lambda$, as desired. To obtain a Frobenius norm bound, note that $\|A - \tilde{A}\|_F \leq \sqrt{n} \cdot \|A - \tilde{A}\| \leq \lambda\sqrt{n}$. \square

Our plan is to form the kernel matrix over the dataset $Q \cup K$ implicitly via Nyström approximation, then extract corresponding blocks to approximate $\exp(QK^\top)$.

Corollary A.2. *Let $Q, K \in \mathbb{R}^{n \times d}$ and let $E \in \mathbb{R}^{2n \times 2n}$ be the exponential kernel matrix over the dataset $Q \cup K$, suppose there exists an $\tilde{E} \in \mathbb{R}^{2n \times 2n}$ such that $E \preceq \tilde{E} \preceq E + \lambda I$ for some $\lambda > 0$, then there exists $\tilde{A} \in \mathbb{R}^{n \times n}$ such that*

$$\|\tilde{A} - \exp(QK^\top)\| \leq \lambda \text{ and } \|\tilde{A} - \exp(QK^\top)\|_F \leq \lambda\sqrt{n}.$$

Proof. The result is a consequence of Lemma A.1 by identifying that

$$E = \begin{bmatrix} \exp(QQ^\top) & \exp(QK^\top) \\ \exp(KQ^\top) & \exp(KK^\top) \end{bmatrix},$$

and \tilde{E} contains proper approximations for the desired blocks. \square

It remains to give an efficient algorithm to approximate the exponential kernel matrix E . A popular scheme is via Nyström approximation (Williams & Seeger, 2000): the algorithm selects a subset of

“landmark” points, and constructs \tilde{E} through these landmarks. Musco & Musco (2017) uses recursive ridge leverage score sampling to generate such an approximation efficiently. Musco & Musco (2017) presents an algorithm that uses $\tilde{O}(ns_\lambda \log(1/\delta))$ kernel function evaluations and $\tilde{O}(ns_\lambda^2 \log(1/\delta))$ additional runtime to compute an approximation \tilde{K} satisfying $K \preceq \tilde{K} \preceq K + \lambda I$ with probability at least $1 - \delta$. We restate their main result here for the sake of completeness.

Lemma A.3 (Theorem 7 of Musco & Musco (2017)). *Let $s = O(s_\lambda \log(s_\lambda/\delta))$, there exists a weighted sampling matrix $S \in \mathbb{R}^{n \times s}$, such that the Nyström approximation of E , $\tilde{E} = ES(S^\top ES)^\dagger S^\top E$ satisfies*

$$E \preceq \tilde{E} \preceq E + \lambda I,$$

holds with probability at least $1 - \delta$. Moreover, S can be computed using $O(ns)$ kernel evaluations and $O(ns^2)$ additional time.

Our main contribution is a quantum algorithm that generates the approximation in *sublinear time*. Before introducing the algorithm, we recall several key concepts.

Lemma A.3 relies on approximating the ridge leverage score on a sample, which can be captured by the notion of generalized ridge leverage score.

Definition A.4 (Generalized ridge leverage score, Musco & Musco (2017)). *Let $E \in \mathbb{R}^{n \times n}$ be a kernel matrix, let $\lambda > 0$, and let $S \in \mathbb{R}^{n \times s}$ be any weighted sampling matrix, the λ -generalized ridge leverage score with respect to S , $i \in [n]$,*

$$\tilde{\tau}_i^\lambda := \frac{1}{\lambda} (E - ES(S^\top ES + \lambda I)^{-1} S^\top E)_{i,i},$$

let $B \in \mathbb{R}^{n \times n}$ be any factorization of $E = BB^\top$, it can be equivalently defined as

$$\tilde{\tau}_i^\lambda = b_i^\top (B^\top S^\top SB + \lambda I)^{-1} b_i,$$

where b_i is the i -th row of B .

We also need a procedure introduced in Apers & Gribling (2023) that generates a spectral approximation of an $n \times d$ matrix, given only queries to its rows, using quantum leverage score sampling. We record it here.

Lemma A.5 (Theorem 3.1 of Apers & Gribling (2023)). *Let $U \in \mathbb{R}^{n \times d}$, $\epsilon, \delta \in (0, 1)$. There exists a quantum algorithm that computes a weighted sampling matrix $S \in \mathbb{R}^{n \times s}$ with $s = O(\epsilon^{-2} d \log(d/\delta))$ such that with probability at least $1 - \delta$,*

$$(1 - \epsilon)U^\top U \preceq U^\top SS^\top U \preceq (1 + \epsilon)U^\top U.$$

The quantum algorithm uses $\tilde{O}(\epsilon^{-1} n^{0.5} d^{0.5})$ row queries to U , and it takes time $\tilde{O}(\epsilon^{-1} n^{0.5} d^{1.5} + d^\omega)$. Moreover, if the leverage score sampling matrix contains $s \leq d$ rows, then the algorithm uses $\tilde{O}(n^{0.5} s^{0.5})$ row queries to U and it takes time $\tilde{O}(n^{0.5} s^{0.5} d + d^\omega)$. We use $\text{QLEVERAGE SCORE}(U, s)$ to denote this procedure that produces a leverage score sampling matrix $S \in \mathbb{R}^{n \times s}$.

We prove the key algorithmic result of this section.

Theorem A.6. *Let $\{x_1, \dots, x_n\} \subseteq \mathbb{R}^d$ be a dataset, $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^m$ be a kernel function, $\lambda > 0$ and $\delta \in (0, 1)$. Let E be the kernel matrix where $E_{i,j} = K(x_i, x_j)$. Suppose $s = O(s_\lambda \log(s_\lambda/\delta))$, then Algorithm 1 computes a weighted sampling matrix $S \in \mathbb{R}^{n \times s}$ such that with probability at least $1 - \delta$,*

$$E \preceq \tilde{E} \preceq E + \lambda I,$$

where $\tilde{E} = ES(S^\top ES)^\dagger S^\top E$. Moreover, S can be computed in $\tilde{O}(n^{0.5} s^{0.5})$ row queries to Q , K and in time $\tilde{O}(n^{0.5} s^{1.5} \cdot (\mathcal{T}_K + s) + s^\omega)$, where \mathcal{T}_K is the time to evaluate the kernel function.

Proof. We note that the major differences between Algorithm 1 and the algorithm in Musco & Musco (2017) are

- Musco & Musco (2017) algorithm is recursive, our algorithm unrolls the recursion and iteratively constructs the weighted sampling matrix;
- Musco & Musco (2017) computes all p_i 's classically, while we use QSAMPLE to generate samples.

Hence, the correctness is automatically satisfied. It remains to give a bound on the running time.

- Computing M_0 : $M_0 \in \mathbb{R}^{s \times s}$ contains the values of kernel functions over s^2 pairs, forming it takes $O(s^2) \cdot \mathcal{T}_K$ time;
- Computing \widehat{M} : we maintain the invariant that $M_t \in \mathbb{R}^{s \times s}$ for all $t \in [T]$, therefore computing \widehat{M} is inverting a $s \times s$ matrix, which takes $O(s^\omega)$ time;
- Computing $D_{t-1}^\top K_i$: this operation involves computing s weighted kernel function evaluations, given D_{t-1} stores a list of s indices together with weights, it can be done in $O(s) \cdot \mathcal{T}_K$ time;
- Oracle for q_i : for any fixed i , note that we need to form $D_{t-1}^\top K_i$ using $O(s) \cdot \mathcal{T}_K$ time, and computing the quadratic form takes $O(s^2)$ time. Thus each oracle call takes $O(s) \cdot \mathcal{T}_K + O(s^2)$ time;
- Computing \widetilde{D}_t : this step requires to compute at most n probabilities, and each probability can be computed via an oracle call in $O(s) \cdot \mathcal{T}_K + O(s^2)$ time, so it remains to give a bound on the sum of probabilities. By the definition of p_i ,

$$\sum_{i=1}^n p_i \leq 16 \log(2s/\delta) \sum_{i=1}^n q_i,$$

and the sum of q_i 's is

$$\begin{aligned} \sum_{i=1}^n q_i &= \frac{5}{\lambda} \cdot (\mathbf{K}(x_i, x_i) - (D_{t-1}^\top K_i)^\top \widehat{M} (D_{t-1}^\top K_i)) \\ &= \frac{5}{\lambda} \cdot (E - ED_{t-1}(D_{t-1}^\top ED_{t-1} + \lambda I)^{-1} D_{t-1}^\top E)_{i,i} \\ &= 5 \cdot \sum_{i=1}^n \widetilde{\tau}_i^\lambda, \end{aligned}$$

by Theorem 8 of Musco & Musco (2017), the sum of λ -generalized ridge leverage score with sampling matrix D_{t-1} is at most $O(s_\lambda \log(s_\lambda/\delta)) = s$, thus the runtime is $\widetilde{O}(n^{0.5} s^{1.5} \cdot (\mathcal{T}_K + s))$.

Finally, note that the loop is dominated by the last iteration, and at each iteration, the number of points to consider is divided by half, we conclude the overall runtime of Algorithm 1 is

$$\widetilde{O}(n^{0.5} s^{1.5} \cdot (\mathcal{T}_K + s) + s^\omega),$$

as desired. \square

We can then apply Theorem A.6 to exponential kernel function and the dataset $Q \cup K$ to compute a Nyström sampling matrix S .

Corollary A.7. *Let $Q, K \in \mathbb{R}^{n \times d}$, $\lambda > 0$ and $\delta \in (0, 1)$. Define the dataset $X = \{x_1, x_2, \dots, x_{2n}\} \subseteq \mathbb{R}^d$ where for $i \in [n]$, $x_i = q_i$ and for $i \in \{n+1, \dots, 2n\}$, $x_i = k_i$. Let E be the kernel matrix where $E_{i,j} = \exp(\langle x_i, x_j \rangle)$. Suppose $s = O(s_\lambda \log(s_\lambda/\delta))$, then there exists an algorithm that computes a weighted sampling matrix $S \in \mathbb{R}^{2n \times s}$ such that, let $\widetilde{E} = ES(S^\top ES)^\dagger S^\top E$, then with probability at least $1 - \delta$, $E \preceq \widetilde{E} \preceq E + \lambda I$. Moreover, S can be computed in $\widetilde{O}(n^{0.5} s^{1.5} \cdot (d + s) + s^\omega)$ time.*

Proof. Apply Theorem A.6 to the kernel function $\mathbf{K}(x_i, x_j) = \exp(\langle x_i, x_j \rangle)$ and note that the kernel function can be computed in $O(d)$ time. \square

B ESTIMATING THE NORMALIZATION FACTOR

Given a sublinear quantum algorithm to approximate the matrix $\exp(QK^\top)$, our next step is to estimate the normalization factor $\exp(QK^\top)\mathbf{1}_n$ to compute the softmax matrix. We first show that given a Nyström approximation to the $2n \times 2n$ kernel matrix E , how to compute the normalization factor and the approximate guarantees.

Lemma B.1. *Let $M \in \mathbb{R}^{n \times n}$ be a symmetric matrix, then we have*

$$\|M\|_\infty \leq \sqrt{n} \cdot \|M\|.$$

Proof. Fix any $i \in [n]$, we examine the row $M_{i,*}$, set the test vector x to be $x_j = \begin{cases} +1, & \text{if } M_{i,j} \geq 0, \\ -1, & \text{otherwise.} \end{cases}$, then

$$\begin{aligned} \|M_{i,*}\|_1 &= M_{i,*}^\top x \\ &= \langle M e_i, x \rangle \\ &\leq \|M e_i\|_2 \cdot \|x\|_2 \\ &\leq \|M\| \cdot \|x\|_2 \\ &= \sqrt{n} \cdot \|M\|. \end{aligned}$$

The conclusion can be achieved by noting that this bound works for any row i . \square

There are two major issues for estimating the normalization factor:

- Corollary A.7 only allows us to compute the sampling matrix in sublinear time, explicitly forming the Nyström approximation \tilde{E} however, would require $\Omega(n)$ time since the matrix is of size $n \times n$;
- Even though we are given the explicit factorization $\tilde{E} = UU^\top$ where $U \in \mathbb{R}^{2n \times s}$, we would have to compute n normalization factors, which would require $\Omega(n)$ time.

In other words, because the output has size $\Omega(n)$, one cannot expect any quantum algorithm to run in $o(n)$ time. Instead, we design a quantum data structure with preprocessing time $o(n)$ time, and can support query to compute the normalization factor to any row efficiently.

In particular, we are interested in the following algorithmic task: given query access to the rows of a matrix $U \in \mathbb{R}^{n \times s}$ and a vector $v \in \mathbb{R}^n$, output a vector $\tilde{\mu} \in \mathbb{R}^s$ such that $\|\tilde{\mu} - U^\top v\|_{(U^\top U)^{-1}} \leq \epsilon$, which can be solved via Lemma 2.6. For our application, $\|v\|_\infty = 1$. However, we are interested in the quantity $UU^\top v$ so we need to measure the error $\|U(\tilde{\mu} - U^\top v)\|_2$. How would a bound on the $\|\cdot\|_{(U^\top U)^{-1}}$ be useful? We prove a structural lemma below.

Lemma B.2. *Let $x \in \mathbb{R}^s$ and $U_2 \in \mathbb{R}^{n \times s}$ satisfying $\|x\|_{(U_2^\top U_2)^{-1}} \leq \epsilon$ for some $\epsilon \in (0, 1)$, let $U_1 \in \mathbb{R}^{n \times s}$, we have*

$$\|U_1 x\|_2 \leq \epsilon \cdot \|U_1 U_2^\top\|.$$

Proof. We define the vector $y = (U_2^\top U_2)^{-1} x$ and $z = U_2 y$, then

$$\begin{aligned} \|z\|_2^2 &= y^\top U_2^\top U_2 y \\ &= x^\top (U_2^\top U_2)^{-1} x \\ &= \|x\|_{(U_2^\top U_2)^{-1}}^2 \\ &\leq \epsilon^2, \end{aligned}$$

moreover, the vector of interest is $U_1 x$ which is

$$\begin{aligned} U_1 x &= U_1 (U_2^\top U_2)^{-1} y \\ &= (U_1 U_2^\top)^{-1} U_2 y \end{aligned}$$

$$= (U_1 U_2^\top) z,$$

subsequently its ℓ_2 norm can be bounded as

$$\begin{aligned} \|U_1 x\|_2 &\leq \|U_1 U_2^\top\| \cdot \|z\|_2 \\ &\leq \epsilon \cdot \|U_1 U_2^\top\|, \end{aligned}$$

as desired. \square

Corollary B.3. *Let $\epsilon \in (0, 1)$, $U_1, U_2 \in \mathbb{R}^{n \times s}$ where $\tilde{A} = U_1 U_2^\top$, $v \in \mathbb{R}^n$, suppose there exists a vector $\tilde{\mu} \in \mathbb{R}^s$ with $\|\tilde{\mu} - U_2^\top v\|_{(U_2^\top U_2)^{-1}} \leq \epsilon$, then we have*

$$\|\tilde{A}v - U_1 \tilde{\mu}\|_2 \leq \epsilon \cdot \|U_1 U_2^\top\|.$$

Proof. We will apply Lemma B.2 by setting $x = \tilde{\mu} - U_2^\top v$ and by noting that $U_1 x = U_1 \tilde{\mu} - U_1 U_2^\top v = U_1 \tilde{\mu} - \tilde{A}v$. \square

We are now in the position to state our formal theorem, which provides an end-to-end guarantee on estimating the normalization factor. For simplicity, we will prove the statement with high probability guarantee, i.e., the success probability is $1 - 1/\text{poly}(n)$.

Theorem B.4. *Let $Q, K \in \mathbb{R}^{n \times d}$, $\lambda > 0$ and $\epsilon \in (0, 1)$. Let $s = \tilde{O}(s_\lambda)$ where s_λ is the statistical dimension of the exponential kernel on $Q \cup K$. There exists a data structure (Algorithm 2) with the following specification:*

- Preprocessing in $\tilde{O}(n^{0.5} s^{0.5} / \epsilon)$ row queries to Q, K and time $\tilde{O}(n^{0.5} s^{1.5} (s + d) / \epsilon + s^\omega)$;
- For any $i \in [n]$, it outputs an approximate normalization factor for row i in time $O(s(s + d))$.

Moreover, with probability at least $1 - 1/\text{poly}(n)$, it holds that for any $i \in [n]$, the output b_i satisfies

$$|b_i - \exp(q_i K^\top) \mathbf{1}_n| \leq O(\epsilon \|A\| + \lambda \sqrt{n}),$$

if $\frac{\lambda \sqrt{n}}{\|A\|} \leq 1$, then the bound can be further simplified to

$$|b_i - \exp(q_i K^\top) \mathbf{1}_n| \leq O(\lambda \sqrt{n}),$$

and the preprocessing time simplifies to

$$\tilde{O}(s^{1.5} (s + d) \|A\| / \lambda + s^\omega).$$

Proof. Given Q, K , let $E \in \mathbb{R}^{2n \times 2n}$ be the associated exponential kernel matrix. We will first invoke Corollary A.7 to compute a sampling matrix $S \in \mathbb{R}^{2n \times s}$ where $s = \tilde{O}(s_\lambda)$ such that $\tilde{E} = ES(S^\top ES)^\dagger S^\top E$ approximates E , in time $\tilde{O}(n^{0.5} s^{1.5} (s + d) + s^\omega)$. Set $U = ES(S^\top ES)^\dagger / 2$, we have that $\tilde{E} = UU^\top$ for $U = \begin{bmatrix} U_1 \\ U_2 \end{bmatrix}$ with $U_1, U_2 \in \mathbb{R}^{n \times s}$, and our desired approximate block for A is $\tilde{A} = U_1 U_2^\top$. Note that forming U explicitly would take $\Omega(n)$ time, so we instead implement a row oracle for U_2 . Since $U_2 \in \mathbb{R}^{n \times s}$, we only need to compute s entries for each row, and let $N = (S^\top ES)^\dagger / 2$, we see that $(U_2)_{j,*} = N(ES)_{j+n,*}$ and $(ES)_{j+n,*}$ contains values in the form of $S_k \cdot \exp(\langle x_{j+n}, x_k \rangle)$ for $k \in S$. N can be computed in $O(s^2 d + s^\omega)$ time, and row oracle for any $j \in [n]$ can be implemented in $O(sd + s^2)$ time. By Lemma 2.6, $\tilde{\mu}$ can be computed in $\tilde{O}(n^{0.5} s^{1.5} (s + d) / \epsilon)$ time. To query the normalization factor for row i , note that it can be computed via $(U_1 \tilde{\mu})_i = \langle (U_1)_{i,*}, \tilde{\mu} \rangle$, which can be computed using row oracle, in $O(s(s + d))$ time. Thus, the overall runtime of our procedure can be summarized as

- Preprocessing time $\tilde{O}(n^{0.5} s^{1.5} (s + d) / \epsilon + s^\omega)$;
- Query time $O(s(s + d))$.

It remains to give an approximation guarantee. With probability at least $1 - 1/\text{poly}(n)$, we have $\|A - \tilde{A}\| \leq \lambda$, and observe that

$$\begin{aligned} |\tilde{a}_i^\top \mathbf{1}_n - \exp(q_i K^\top) \mathbf{1}_n| &\leq \|(\tilde{A} - A)v\|_\infty \\ &\leq \|\tilde{E} - E\|_\infty \cdot \|v\|_\infty \\ &\leq \lambda\sqrt{n}, \end{aligned}$$

where the second step is by the matrix infinity norm is the induced norm of vector ℓ_∞ norm, and the last step is by Lemma B.1. On the other hand, our final output b_i is an approximation to $\tilde{a}_i^\top \mathbf{1}_n$. Let $\tilde{y} := U_1 \tilde{\mu}$, by Corollary B.3, we have

$$\|\tilde{A}v - \tilde{y}\|_2 \leq \epsilon \cdot \|\tilde{A}\|,$$

this holds with probability at least $1 - \delta$, conditioning on this event, note that by Lemma A.1, we have that $\|\tilde{A}\| \leq \|A\| + \lambda$. Thus, we conclude our final result by

$$\begin{aligned} |b_i - \exp(q_i K^\top) \mathbf{1}_n| &\leq |b_i - \tilde{a}_i \mathbf{1}_n| + |\tilde{a}_i^\top \mathbf{1}_n - \exp(q_i K^\top) \mathbf{1}_n| \\ &\leq \|\tilde{A}v - \tilde{y}\|_2 + \lambda\sqrt{n} \\ &\leq \epsilon \cdot (\lambda + \|A\|) + \lambda\sqrt{n}. \end{aligned}$$

Now, suppose $\lambda\sqrt{n} \leq \|A\|$, then we could set $\epsilon = \frac{\lambda\sqrt{n}}{\|A\|}$, the error bound simplifies to $O(\lambda\sqrt{n})$. \square

C APPROXIMATE MATRIX MULTIPLICATION VIA LEVERAGE SCORE

It remains to handle the value matrix, and we will do so via a machinery called approximate matrix multiplication.

Definition C.1 (Approximate matrix multiplication, Clarkson & Woodruff (2017)). *Let $A \in \mathbb{R}^{n \times d}$, $B \in \mathbb{R}^{n \times m}$ and let $C = A^\top B \in \mathbb{R}^{d \times m}$. The approximate matrix multiplication problem asks to design a random matrix $S \in \mathbb{R}^{n \times s}$, such that*

$$\Pr[\|A^\top S S^\top B - C\|_F \leq \epsilon \|A\|_F \|B\|_F] \geq 1 - \delta,$$

where $\epsilon, \delta \in (0, 1)$. We call such S satisfying (ϵ, δ) -AMM.

To generate the random matrix S , our strategy will be performing leverage score sampling over V . However, standard proof (see, e.g., Clarkson & Woodruff (2017)) requires V to have orthonormal columns. We provide a proof for the case where V does not have orthonormal columns (albeit it requires extra factors in blowups). Before doing so, we define a parameter that quantifies this blowup which we call *row distortion*.

Definition C.2 (Row distortion). *Let $A \in \mathbb{R}^{n \times d}$ for $n \geq d$, we define the row distortion of A , denoted by $\alpha(A)$, as*

$$\alpha(A) := \frac{d}{\|A\|_F^2} \cdot \max_{i \in [n]} \frac{\|a_i\|_2^2}{\tau_i},$$

where a_i is the i -th row of A and τ_i is the i -th leverage score of A (Definition 2.1). When A is clear from context, we use α as an abbreviation.

Lemma C.3. *Let $A \in \mathbb{R}^{n \times d}$ with $n \geq d$, then the row distortion of A satisfies*

$$\alpha(A) \leq \frac{d}{\text{srank}(A)},$$

where $\text{srank}(A) = \frac{\|A\|_F^2}{\|A\|_2^2}$ is the stable rank of A .

Proof. We derive an upper bound on $\|a_i\|_2^2$, let $A = U\Sigma V^\top$ be its SVD, then

$$\|a_i\|_2^2 = \|e_i^\top U \Sigma V^\top\|_2^2$$

$$\begin{aligned}
&\leq \|u_i\|_2^2 \cdot \|\Sigma V^\top\|^2 \\
&= \tau_i \cdot \|U \Sigma V^\top\|^2 \\
&= \tau_i \cdot \|A\|^2,
\end{aligned}$$

where the third step is by the definition of leverage score and spectral norm is unitary invariant. We thus obtain the following bound on $\alpha(A)$:

$$\begin{aligned}
\alpha(A) &= \frac{d}{\|A\|_F^2} \cdot \max_{i \in [n]} \frac{\|a_i\|_2^2}{\tau_i} \\
&\leq \frac{d}{\|A\|_F^2} \cdot \max_{i \in [n]} \frac{\tau_i \cdot \|A\|^2}{\tau_i} \\
&= d \cdot \frac{\|A\|^2}{\|A\|_F^2} \\
&= \frac{d}{\text{srank}(A)},
\end{aligned}$$

where we recall that $\text{srank}(A) = \frac{\|A\|_F^2}{\|A\|^2}$. \square

We are now ready to prove a generalized approximate matrix multiplication based on leverage score sampling, when the matrix does not have orthonormal columns.

Lemma C.4. *Let $A \in \mathbb{R}^{n \times d}$, $B \in \mathbb{R}^{n \times m}$, let $S \in \mathbb{R}^{n \times s}$ be the leverage score sampling matrix of A with $s = (\epsilon^{-2} \alpha \log(1/\delta))$ for $\epsilon, \delta \in (0, 1)$ and α is the row distortion of A (Definition C.2). Then, S is an (ϵ, δ) -AMM.*

Proof. For the sampling matrix S , it is a scaled submatrix of the permutation matrix, where for any $m \in [s]$, $S_{m, z_m} = \frac{1}{\sqrt{s p_m}}$ where $p_m \geq \frac{\tau_m}{d}$ and $z_m = i$ with probability p_i . Let a_i, b_j denote the i -th and j -th row of A and B , respectively. We can write

$$A^\top S S^\top B - A^\top B = \frac{1}{s} \sum_{i \in [n], m \in [s]} a_i b_i^\top \left(\frac{\mathbb{I}[z_m = i]}{p_i} - 1 \right),$$

taking expectation, we obtain

$$\begin{aligned}
\mathbb{E}[A^\top S S^\top B - A^\top B] &= \frac{1}{s} \sum_{i=1}^n a_i b_i^\top \left(\frac{p_i}{p_i} - 1 \right) \\
&= 0,
\end{aligned}$$

to bound the second moment of $\|A^\top S S^\top B - A^\top B\|_F$, we first expand the definition of Frobenius norm square:

$$\begin{aligned}
&\mathbb{E} \text{tr}[(A^\top S S^\top B - A^\top B)(A^\top S S^\top B - A^\top B)] \\
&= \mathbb{E} \frac{1}{s^2} \text{tr} \left[\sum_{i, j \in [n], m \in [s]} b_j a_j^\top a_i b_i^\top \left(\frac{\mathbb{I}[z_m = j]}{p_j} - 1 \right) \left(\frac{\mathbb{I}[z_m = i]}{p_i} - 1 \right) \right] \\
&= \frac{1}{s^2} \sum_{m=1}^s \text{tr} \left[\sum_{i=1}^n \frac{1}{p_i} \cdot b_i a_i^\top a_i b_i^\top - B^\top A A^\top B \right] \\
&= \frac{1}{s} \text{tr} \left[\sum_{i=1}^n \frac{1}{p_i} \cdot b_i a_i^\top a_i b_i^\top - B^\top A A^\top B \right] \\
&\leq \frac{1}{s} \left(\sum_{i=1}^n \frac{1}{p_i} \|a_i\|_2^2 \|b_i\|_2^2 - \text{tr}[B^\top A A^\top B] \right) \\
&\leq \frac{1}{s} (\alpha \|A\|_F^2 \|B\|_F^2 - \|A^\top B\|_F^2)
\end{aligned}$$

$$\leq \frac{\alpha}{s} \|A\|_F^2 \|B\|_F^2,$$

where the first step is by definition of S , the second step is by applying expectation and use $\mathbb{E}[A^\top S S^\top B - A^\top B] = 0$, the fourth step is by $\text{tr}[b_i a_i^\top a_i b_i^\top] = \|a_i b_i^\top\|_F^2 \leq \|a_i\|_2 \|b_i\|_2$, the fifth step is by $p_i \geq \frac{\tau_i}{d}$, therefore

$$\begin{aligned} \frac{1}{p_i} &\leq \frac{d}{\tau_i} \\ &= \frac{\|A\|_F^2}{\|a_i\|_2^2} \cdot \frac{d}{\|A\|_F^2} \cdot \frac{\|a_i\|_2^2}{\tau_i} \\ &\leq \alpha \cdot \frac{\|A\|_F^2}{\|a_i\|_2^2}, \end{aligned}$$

where the last step is by the definition of α . By Chebyshev's inequality, we can choose $s = O(\alpha/\epsilon^2)$ so that the approximate matrix multiplication holds with constant probability, and one could boost the success probability to $1 - \delta$ by either taking $\log(1/\delta)$ independent copies via a Chernoff bound, or directly through Bernstein inequality. \square

We are ready to state our final result on approximating the value matrix V .

Theorem C.5. *Let $V \in \mathbb{R}^{n \times d}$, $\epsilon \in (0, 1)$ and α be the row distortion of V . There exists a quantum algorithm that computes a weighted sampling matrix $S \in \mathbb{R}^{n \times s}$ with $s = \tilde{O}(\epsilon^{-2}\alpha)$ such that for any fixed matrix $B \in \mathbb{R}^{n \times m}$, S is an $(\epsilon, 1/\text{poly}(n))$ -AMM. Moreover, S can be computed using $\tilde{O}(\epsilon^{-1}n^{0.5}\alpha^{0.5})$ row queries to V and $\tilde{O}(\epsilon^{-1}n^{0.5}\alpha^{0.5}d + d^\omega)$ time.*

Proof. The proof is by composing Lemma A.5 and Lemma C.4, and note that for $\tilde{O}(\epsilon^{-2}\alpha)$ rows, the sum of leverage scores is at most $\tilde{O}(\epsilon^{-2}\alpha)$. \square

D PUT THINGS TOGETHER

We are now ready to state our final algorithm and its guarantee. Recall that, we define $D = \exp(QK^\top)\mathbf{1}_n$ and $D' = \exp(KQ^\top)\mathbf{1}_n$. We use \tilde{D}, \tilde{D}' to denote their approximations.

We prove a simple inequality that quantifies the perturbation on the inverse.

Lemma D.1. *Let $C, D \in \mathbb{R}^{n \times n}$, if D is nonsingular and $\|C - D\| \leq \epsilon$, and $\|D^{-1}\| < 1/\epsilon$, then C is also nonsingular and $\|C^{-1}\| \leq \frac{\|D^{-1}\|}{1 - \epsilon \cdot \|D^{-1}\|}$.*

Proof. We will make use of Neumann series, which states that for $\|A\| < 1$, $(I - A)^{-1}$ admits the expansion

$$(I - A)^{-1} = \sum_{k=0}^{\infty} A^k,$$

this leads to a bound on the norm:

$$\begin{aligned} \|(I - A)^{-1}\| &= \left\| \sum_{k=0}^{\infty} A^k \right\| \\ &\leq \sum_{k=0}^{\infty} \|A^k\| \\ &\leq \sum_{k=0}^{\infty} \|A\|^k \\ &= \frac{1}{1 - \|A\|}, \end{aligned} \tag{1}$$

now, to prove our desired bound, we write $C = D + E$ where E is the perturbation, then $C = D + E = D(I + D^{-1}E)$, and we will apply Eq. (1) to $-D^{-1}E$:

$$\begin{aligned}\|D^{-1}E\| &\leq \|D^{-1}\| \cdot \|E\| \\ &= \|D^{-1}\| \cdot \|C - D\| \\ &< 1/\epsilon \cdot \epsilon \\ &= 1,\end{aligned}$$

therefore

$$\begin{aligned}\|C^{-1}\| &= \|(I + D^{-1}E)^{-1}D^{-1}\| \\ &\leq \|D\| \cdot \|(I - D^{-1}E)^{-1}\| \\ &\leq \frac{\|D^{-1}\|}{1 - \|D^{-1}E\|} \\ &\leq \frac{\|D^{-1}\|}{1 - \|E\| \cdot \|D^{-1}\|} \\ &\leq \frac{\|D^{-1}\|}{1 - \epsilon \cdot \|D^{-1}\|},\end{aligned}$$

this completes the proof. \square

Theorem D.2 (Formal version of Theorem 3.1). *Let $Q, K, V \in \mathbb{R}^{n \times d}$ be the query, key and value matrices for attention, let $\epsilon, \lambda > 0$. Let $E \in \mathbb{R}^{2n \times 2n}$ be the exponential kernel matrix with the dataset $Q \cup K$, and let s_λ be the statistical dimension of E (Definition 2.2), α be the row distortion of V (Definition C.2). There exists a quantum data structure (Algorithm 3) that preprocesses Q, K, V only through row queries to these matrices and with probability at least $1 - 1/\text{poly}(n)$, for any $i \in [n]$, it outputs a vector $\tilde{r}_i \in \mathbb{R}^d$ where*

$$\tilde{r}_i = e_i^\top \tilde{D}^{-1} \tilde{A} \tilde{V}.$$

If in addition, we have $\|D^{-1}\| < \frac{1}{\epsilon\|A\| + \lambda\sqrt{n}}$, then the approximations \tilde{D}, \tilde{A} and \tilde{V} satisfy that

$$\|\tilde{D}^{-1} \tilde{A} \tilde{V} - D^{-1} A V\|_F \leq \epsilon \cdot (\beta \cdot \|D^{-1}\|) \cdot (\|A\|_F + \lambda\sqrt{n}) \cdot \|V\|_F,$$

where $\beta = \frac{1}{1 - (\epsilon\|A\| + \lambda\sqrt{n})\|D^{-1}\|}$. Moreover, the algorithm has the following runtime specification:

- Preprocesses in $\tilde{O}(\epsilon^{-1}n^{0.5}s_\lambda^{0.5})$ row queries to Q, K and $\tilde{O}(\epsilon^{-1}n^{0.5}\alpha^{0.5})$ row queries to V , and $\tilde{O}(\epsilon^{-1}n^{0.5}(s_\lambda^{2.5} + s_\lambda^{1.5}d + \alpha^{0.5}d) + d^\omega + s_\lambda^\omega + \epsilon^{-2}s_\lambda\alpha d)$ time;
- For any $i \in [n]$, it outputs \tilde{r}_i in $\tilde{O}(s_\lambda^2 + s_\lambda d)$ time.

Proof. By Theorem C.5, we know that with probability at least $1 - 1/\text{poly}(n)$, the following bound holds:

$$\begin{aligned}\|\tilde{D}^{-1} \tilde{A} S_V S_V^\top V\|_F &\leq \epsilon \cdot \|\tilde{D}^{-1} \tilde{A}\|_F \cdot \|V\|_F \\ &\leq \epsilon \cdot \|\tilde{D}^{-1}\| \cdot \|\tilde{A}\|_F \cdot \|V\|_F,\end{aligned}$$

where the second step is by $\|\tilde{D}^{-1} \tilde{A}\|_F \leq \|\tilde{D}^{-1}\| \cdot \|\tilde{A}\|_F$. By Theorem B.4, we know that

$$\|\tilde{D} - D\| \leq \epsilon\|A\| + \lambda\sqrt{n},$$

note that as long as the error satisfies that $\|D^{-1}\| < \frac{1}{\epsilon\|A\| + \lambda\sqrt{n}}$, then by Lemma D.1, we obtain a bound on $\|\tilde{D}^{-1}\|$:

$$\|\tilde{D}^{-1}\| \leq \frac{\|D^{-1}\|}{1 - (\epsilon\|A\| + \lambda\sqrt{n})\|D^{-1}\|}.$$

Finally, by Corollary A.2, we have

$$\|\tilde{A}\|_F \leq \|A\|_F + \lambda\sqrt{n}.$$

For the runtime, it suffices to combine Corollary A.7, Theorem B.4 and Theorem C.5, and the only additional runtime term is the $\epsilon^{-2}s_\lambda\alpha d$, which is the time to form matrix \tilde{R} and \tilde{L} . \square

E EMPIRICAL VERIFICATIONS ON PARAMETERS

In this section, we empirically verify the assumptions on the parameters. In particular, we focus on the following metrics:

- $\|D^{-1}\| \leq \frac{1}{\epsilon\|A\|+\lambda\sqrt{n}}$, we specifically check that what is the maximum possible ϵ so that $\|D^{-1}\| \leq \frac{1}{\epsilon\|A\|}$.
- $\frac{\|A\|_F}{\|A\|}$, this is important as our error guarantee is in terms of Frobenius norm rather than the more typical spectral norm (Zandieh et al., 2023; Han et al., 2024), we verify that this ratio is small.
- $\frac{\|V\|_F}{\|V\|}$, this is similar to the above test, we verify that this ratio is close to \sqrt{d} .
- $\frac{d}{\text{srank}(V)}$, this quantity serves as an upper bound of $\alpha(V)$, we verify that this quantity is a small constant rather than the upper bound d .
- $\frac{\|A\|_\infty}{\|A\|}$, in our error analysis, we have to pay an extra \sqrt{n} factor when converting the spectral norm to matrix infinity norm, we empirically show that this ratio is a small constant rather than the \sqrt{n} scaling.

To conduct our experiment, we use the OLMo2-1B and OLMo2-7B models, in particular their stage1 pretraining checkpoints (Walsh et al., 2025). We list the model architecture in the following.

	Sequence length n	Value dimension d	Number of layers L	Number of heads H
OLMo2-1B	4096	128	16	16
OLMo2-7B	4096	128	32	32

Table 2: Model architecture for OLMo2-1B and OLMo2-7B.

We compute the corresponding attention modules D, A, V using the pretraining datasets for these models, with batch size 2 and 16 batches. We then compute the statistics for each head and each layer, then aggregate the statistics over all layers. We report the mean of these statistics.

	ϵ_{\max}	$\frac{\ A\ _F}{\ A\ }$	$\frac{\ V\ _F}{\ V\ }$	$\frac{d}{\text{srank}(V)}$	$\frac{\ A\ _\infty}{\ A\ }$
OLMo2-1B	0.1708	1.3769	11.3137	1.7126	2.7439
OLMo2-7B	0.1685	1.3586	11.3137	2.1345	2.7439

Table 3: Mean statistics across all heads and all layers. ϵ_{\max} is the maximum ϵ such that $\|D^{-1}\| \leq \frac{1}{\epsilon\|A\|}$.

Through these verifications, we make the following preliminary observations:

- To satisfy the $\|D^{-1}\| \leq \frac{1}{\epsilon\|A\|}$ assumption, it is enough to pick $\epsilon \leq 0.17$, which is larger than common choice of $\epsilon \approx 0.1$. This gives us enough room to tune the parameter ϵ to achieve a good balance between efficiency and accuracy.
- The ratio $\frac{\|A\|_F}{\|A\|}$ is a constant smaller than 2, much smaller than the worst case \sqrt{n} predicted by the theory (recall that $n = 4096$ and $\sqrt{n} = 64$). This suggests that we don’t need to scale down ϵ by a factor of \sqrt{n} to recover the spectral norm error guarantee.
- The ratio $\frac{\|V\|_F}{\|V\|} \approx 11$ is roughly \sqrt{d} as $d = 128$, this confirms the theory, but it does not impair the sublinear scaling in n of our algorithm: we could simply scale down ϵ by a factor of \sqrt{d} to absorb this blowup and increase the runtime by a factor of \sqrt{d} .
- The ratio $\frac{d}{\text{srank}(V)}$ is a small constant, recall that $\text{srank}(V)$ can be as small as 1, causing the ratio to be d , our experiment shows that $\alpha(V)$ is close to a small constant rather than d .
- The ratio $\frac{\|A\|_\infty}{\|A\|}$ is a constant smaller than 3, we check this quantity as in proving the approximation guarantee for \tilde{D} , we make use of the fact that $\|A\|_\infty \leq \sqrt{n} \cdot \|A\|$, this again

shows that instead of the worst case \sqrt{n} scaling, this distortion is only by a constant factor, implying the $\lambda\sqrt{n}$ additive error term is more likely $O(\lambda)$ in practice. This greatly enlarges the range of choice for λ to achieve better speedup.

F BIT COMPLEXITY OF OUR ALGORITHM

In this section, we give a preliminary analysis on the bit complexity of our algorithm, in particular the bit complexity of matrix inversion operation. We will make use of the following standard algorithm for backward stable matrix inversion.

Lemma F.1 (Higham (2002); Harvey & van der Hoeven (2021)). *Let $A \in \mathbb{R}^{s \times s}$ be nonsingular, there exists an algorithm that computes B^{-1} such that*

$$\|B - A\| \leq \delta \cdot s^c \cdot \kappa(A)^{C \log s} \cdot \|A^{-1}\|,$$

for absolute constant $c, C > 0$ with bit complexity $O(s^3 \cdot M(b))$ where $b = O(\log(\kappa(A)) + \log(1/\delta))$ and $M(b) = O(b \log b)$.

We note that the backward stable error guarantee is exactly what has been analyzed in Gu et al. (2024):

Lemma F.2 (Lemma G.3 in Gu et al. (2024)). *Let $A, B \in \mathbb{R}^{s \times s}$ be matrices such that $\|A - B\| \leq \delta$, then*

$$|\tau_i(A) - \tau_i(B)| \leq \delta \cdot \kappa^{2.5}(A).$$

This means that by setting $\delta = 1/\text{poly}(\kappa(A))$, we can approximate the leverage scores well, and the number of bits $b = O(\log(\kappa(A)))$. Note that we apply matrix inversions for two type of matrices:

- $S^\top ES + \lambda I$, where S is the ridge leverage score sampling matrix for E ;
- $V^\top V$, where V is the value matrix.

In the latter case, we only need to pay the $O(\log(\kappa(V)))$ factor, which in practice, is very small: in our experiments, we see that on average, the log of the condition number is smaller than 4 and the largest log of the condition number is smaller than 20. The interesting part is the former case.

To analyze $\kappa(S^\top ES + \lambda I)$, we upper bound the spectral norm and lower bound the smallest eigenvalue. First, observe that $S^\top ES \succeq 0$, so trivially we have $S^\top ES + \lambda I \succeq \lambda I$, thus the smallest eigenvalue is at least λ . To bound $\|S^\top ES + \lambda I\|$, we note that

$$\begin{aligned} \|S^\top ES + \lambda I\| &\leq \|S^\top ES\| + \lambda \\ &\leq \|S\|^2 \cdot \|E\| + \lambda, \end{aligned}$$

we bound the two spectral norms respectively. For $\|S\|^2$, we bound it probabilistically: let $c_i = \begin{cases} 1, & \text{if } i \text{ is sampled with probability } p_i \\ 0, & \text{otherwise} \end{cases}$, and consider the matrix SS^\top , note that by definition, SS^\top is a diagonal matrix with

$$(SS^\top)_{i,i} = \frac{c_i}{p_i},$$

note that as c_i is a Bernoulli random variable with probability p_i , we have $\mathbb{E}[c_i] = p_i$ hence $\mathbb{E}\left[\frac{c_i}{p_i}\right] = 1$, and

$$\begin{aligned} \mathbb{E}[\|S\|^2] &= \mathbb{E}\left[\max_{i \in [n]} \frac{c_i}{p_i}\right] \\ &\leq \mathbb{E}\left[\sum_{i=1}^n \frac{c_i}{p_i}\right] \\ &= \sum_{i=1}^n \mathbb{E}\left[\frac{c_i}{p_i}\right] \end{aligned}$$

$$= n,$$

hence by Markov’s inequality, with constant probability (say 0.99), we have that $\|S\|^2 \leq O(n)$. Condition on this event, we analyze $\|E\|$: let $R = \max\{\max_i \|q_i\|_2^2, \max_i \|k_i\|_2^2\}$, then

$$\begin{aligned} \|E\| &\leq \text{tr}[E] \\ &= \sum_{i=1}^n \exp(\|q_i\|_2^2/\sqrt{d}) + \exp(\|k_i\|_2^2/\sqrt{d}) \\ &\leq 2n \exp(R/\sqrt{d}), \end{aligned}$$

combining the above, we obtain a final (probabilistic) upper bound on the condition number of $S^\top ES + \lambda I$:

$$\begin{aligned} \kappa(S^\top ES + \lambda I) &\leq \frac{\|S\|^2 \cdot \|E\| + \lambda}{\lambda} \\ &\leq 1 + \frac{Cn^2 \exp(R/\sqrt{d})}{\lambda}, \end{aligned}$$

this gives the final bound on $\log(\kappa(S^\top ES + \lambda I))$:

$$\log(\kappa(S^\top ES + \lambda I)) \leq R/\sqrt{d} + \log(n/\lambda).$$

In practice, the data-dependent parameter R/\sqrt{d} is small: for both OLMo2-1B and OLMo2-7B models, these values are 20.1147 and 21.2227 respectively. Hence, the final bit complexity is $\tilde{O}(s^3(d^{-0.5}R + \log(\kappa(V)) + \log(n/\lambda)))$.

When performing leverage score sampling over V , we need to compute the inverse $(V^\top V)^{-1}$, thus it is mandatory to obtain an upper bound on the condition number of V . To compute such an upper bound, we note that the algorithm computes a leverage score sampling matrix S with $O(\epsilon^{-2}d \log d)$ rows, and the matrix $SV \in \mathbb{R}^{\epsilon^{-2}d \log d \times d}$. Computing the condition number and spectral norm of SV can be done classically, in $\text{poly}(d)$ time.

To establish a relation between the conditioning of V and SV , observe that S provides a subspace embedding property: $(1 - \epsilon)V^\top V \preceq V^\top S^\top SV \preceq (1 + \epsilon)V^\top V$, this implies that $\kappa(SV) \leq \sqrt{\frac{1+\epsilon}{1-\epsilon}} \cdot \kappa(V) \leq (1 + O(\epsilon)) \cdot \kappa(V)$. This ensures that the bit complexity b depends only on $O(\log(\kappa(V)))$.

Finally, to compute the spectral norms and condition numbers required by the algorithms, we could use the algorithms in Musco et al. (2018); Shah (2025); Sobczyk (2025).

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.