

DIFFERENTIAL PRIVACY OF CROSS-ATTENTION WITH PROVABLE GUARANTEE

Anonymous authors

Paper under double-blind review

ABSTRACT

Cross-attention has become a fundamental module nowadays in many important artificial intelligence applications, e.g., retrieval-augmented generation (RAG), system prompt, guided stable diffusion, and many more. Ensuring cross-attention privacy is crucial and urgently needed because its key and value matrices may contain sensitive information about model providers and their users. In this work, we design a novel differential privacy (DP) data structure to address the privacy security of cross-attention with a theoretical guarantee. In detail, let n be the input token length of system prompt/RAG data, d be the feature dimension, $0 < \alpha \leq 1$ be the relative error parameter, R be the maximum value of the query and key matrices, R_w be the maximum value of the value matrix, and r, s, ϵ_s be parameters of polynomial kernel methods. Then, our data structure requires $\tilde{O}(ndr^2)$ memory consumption with $\tilde{O}(nr^2)$ initialization time complexity and $\tilde{O}(\alpha^{-1}r^2)$ query time complexity for a single token query. In addition, our data structure can guarantee that the process of answering user query satisfies (ϵ, δ) -DP with $\tilde{O}(n^{-1}\epsilon^{-1}\alpha^{-1/2}R^{2s}R_w r^2)$ additive error and $n^{-1}(\alpha + \epsilon_s)$ relative error between our output and the true answer. Furthermore, our result is robust to adaptive queries in which users can intentionally attack the cross-attention system. To our knowledge, this is the first work to provide DP for cross-attention and is promising to inspire more privacy algorithm design in large generative models (LGMs).

1 INTRODUCTION

The development of Artificial Intelligence (AI) has four stages: (1) prediction AI, e.g., ResNet (He et al., 2016) in image classification; (2) generation AI, e.g., ChatGPT (Achiam et al., 2023) in language generation; (3) autonomous agent AI, Voyager (Wang et al., 2023a) autonomously plays Minecraft game (Fan et al., 2022); (4) Artificial Generalization Intelligence (AGI). Humans have made rapid progress in generative AI, and we are excitingly heading to the third stage, the era of AI agent (Liu et al., 2023). One prevalent application of AI agents is customized large generative models (LGMs) agents (OpenAI, 2024a), e.g., AgentGPT (GitHub, 2024a), SuperAGI (GitHub, 2024d), MetaGPT (Hong et al., 2024b;a), GPT Researcher (GitHub, 2024c) and many so on. In particular, recently, Apple Inc. introduced Apple Intelligence (Apple, 2024), signaling the integration of LGMs into physical devices. This innovation allows devices to use personal information for real-life assistance, such as entering passport numbers when booking flights or informing users of their latest meetings. With increased AI capabilities, privacy concerns become significant, as the more personal information devices handle, the greater the potential privacy risks.

One fundamental technique used in LGMs is cross-attention (Vaswani et al., 2017), which is an essential module in retrieval-augmented generation (RAG) (Lewis et al., 2020), system prompt, guided stable diffusion, and many so on. In RAG, to be more professional, the LGMs answer user input queries by using a domain-specific database under cross-attention, which may contain specific privacy data and knowledge so that the LGMs gain additional power. For system prompts, based on cross-attention, some customized long prompts, e.g., user information or concrete rules, are concatenated before user input to follow human instructions better, which are commonly used in ChatGPT (GitHub, 2024b), Claude3 (Anthropic, 2024) and other commercial LGMs.

Consequently, protecting the privacy of domain-specific data in RAG or system prompts is crucial as they contain sensitive information about users and companies. These data and prompts are the

core assets of many start-ups. However, these data and prompts can be easily recovered (Li et al., 2023b), jailbroken (Jin et al., 2024), and released (Li et al., 2023a) by user adversarial attack (Yu et al., 2024), e.g., there are 1700 tokens in ChatGPT system prompts (Patel, 2024). These findings highlight the critical importance of robust privacy protections in LGMs, making privacy not just essential but an urgent issue that demands immediate attention.

To fundamentally preserve cross-attention privacy, we borrow the powerful tools from differential privacy (DP) (Dwork et al., 2006), which provides measurable privacy and combines with statistical machine learning seamlessly (Ponomareva et al., 2023). Thus, in this work, we would like to ask and answer the following question,

How can we use differential privacy to protect the security of cross-attention in LGMs?

Our work demonstrates that the Softmax cross-attention computation is equivalent to computing the weighted distance problem.

Definition 1.1 (Softmax cross-attention). *Let n and m be the token length of the data and input query, respectively. Let d be the feature dimension. Given fixed key matrix $K \in [0, R]^{n \times d}$ and fixed value matrix $V \in [-R_w, R_w]^{n \times d}$, for any input query matrix $Q \in [0, R]^{m \times d}$, the goal of the Softmax Cross-Attention Computation is to get the matrix $\text{Attn}(Q, K, V) \in \mathbb{R}^{m \times d}$, which is*

$$\text{Attn}(Q, K, V) := D^{-1}AV,$$

where $A \in \mathbb{R}^{m \times n}$ satisfies $A_{i,j} := \exp(\langle Q_i, K_j \rangle / d)$ for any $i \in [m], j \in [n]$ (Q_i and K_j denote the i -th and j -th rows of Q and K , respectively) and $D := \text{diag}(A\mathbf{1}_n) \in \mathbb{R}^{m \times m}$ is a diagonal matrix.

Note that $\text{Softmax}(QK^\top) = D^{-1}A \in \mathbb{R}^{m \times n}$ in Definition 1.1, which is the standard function used in transformers, and usually, we call it as attention matrix. Our main theorem, presented below, provides a robust solution of cross-attention, ensuring privacy and accuracy guarantees.

Theorem 1.2 (Main result; Informal version of Theorem 3.1). *Let Q, K, V, Attn be defined in Definition 1.1. Let $\alpha \in (0, 1)$ be the relative error parameter and p_f be the probability of failure parameter. Let r, s, ϵ_s be the parameters of the polynomial kernel methods (Lemma C.7). Let $l = \tilde{O}(r)$. Then, our Algorithm 1 requires $O(\ln dr)$ memory with $O(\ln r)$ initialization time and $\tilde{O}(\alpha^{-1}lr)$ query time, such that with probability $1 - p_f$, the output process of cross-attention satisfies (ϵ, δ) -DP and is robust to adaptive query with relative error $n^{-1}(\alpha + \epsilon_s)$ and additive error $\tilde{O}(n^{-1}\epsilon^{-1}\alpha^{-1/2}lR^{2s}R_w r)$.*

Our main technique in Theorem 1.2 ensures that cross-attention is differentially private by using the polynomial kernel approximation method and transforming it into a weighted distance problem. We then solve the problem by summing over weighted distances (depending on the value embedding) between the query embedding and the key embedding. We build a data structure for weighted Softmax queries in Section 4.3, and we extend this data structure to handle adaptive queries using the ϵ_0 -net/metric entropy argument in Section 4.4. Furthermore, our error decreases as the input token length grows, diminishing both the relative and additive errors to zero.

Our contributions are as follows:

- We demonstrate that cross-attention computations are equivalent to the weighted distance problem (Section 3).
- We design a novel algorithm (Algorithm 3) that privately answers weighted Softmax queries with high probability and a concrete accuracy bound.
- Our algorithm (Algorithm 1) handles multiple cross-attention queries and is robust against adaptive query attacks (Theorem 3.1), meaning that potential attackers cannot intentionally extract information of system prompts/RAG data.

To our knowledge, this is the first work to utilize DP to protect prompts in LGMs with theoretically provable guarantees. While some have explored protecting user/system prompts with DP (Edemacu & Wu, 2024; Mai et al., 2023), they are primarily empirical and lack theoretical guarantees. Additionally, many others are working on protecting private datasets by applying DP to the fine-tuning stage of LGMs (Behnia et al., 2022; Singh et al., 2024; Liu et al., 2024b; Yu et al., 2021; Li et al.,

2021; Shi et al., 2022a), which diverges from our work. The strength of DP lies in its strong, unambiguous, and concrete definition of privacy, enabling algorithm designs with provable privacy and accuracy analysis. Therefore, we believe that the theoretical aspects of DP applications in LGMs remain a highly impactful direction, and we aim to pave the way for further exploration in this area.

1.1 RELATED WORK

Differential Privacy in Data Structure and Attention. Differential privacy (DP) is a flourishing and powerful technique that has enormous applications in the topic of private machine learning. In the era of Large Generative Models (LGMs), there are three primary approaches to ensuring privacy: (1) during the pre-training stage: to protect training data (Abadi et al., 2016; Ponomareva et al., 2023), (2) during the adaptation stage: to protect target data (Behnia et al., 2022; Singh et al., 2024; Liu et al., 2024b; Yu et al., 2021; Li et al., 2021; Shi et al., 2022a; Huang et al., 2024), (3) during the inference stage: to protect user/system prompts (Edemacu & Wu, 2024) and RAG data (Lewis et al., 2020). To protect training data, DP-SGD (Abadi et al., 2016) uses DP optimizer to ensure data privacy, severing as the traditional baseline method. Recently, numerous works have aimed to improve this method by integrating DP in both the pre-training and fine-tuning stages of LGMs (Yu et al., 2021; Li et al., 2021; Golatkar et al., 2022; Behnia et al., 2022; Shi et al., 2022a; Mattern et al., 2022; Singh et al., 2024; Zheng et al., 2024; Liu et al., 2024b). **However, DP-SGD confines differential privacy to the optimizer. In contrast, we propose a novel approach that integrates DP directly into the attention mechanism, supported by strong theoretical analysis and guarantees. Given the resource-intensive nature of training LGMs, our technique offers a practical alternative for models trained with standard SGD, which lack inherent privacy guarantees. In such cases, applying DP-SGD would require retraining the models, which is computationally expensive, whereas our method avoids this additional cost.**

To protect user/system prompts, Edemacu & Wu (2024) provides a survey on both DP and non-DP methods. In the use of LGMs, prompting methods almost become a standard way for inference (Schulhoff et al., 2024). Given the billions of prompt interactions daily, ensuring privacy is essential (Mai et al., 2023). We refer readers to Appendix A for more related works.

Roadmap. In Section 2, we present the preliminary of differential privacy (DP) and cross-attention. In Section 3, we present the main result of our cross-attention theorem (Theorem 3.1). In Section 4, we outline the main results of our algorithms. In Section 5, we discuss DP-related topics and potential extensions. In Section 6, we conclude our paper.

2 PRELIMINARY

In this section, we give the preliminary of differential privacy (DP) and cross-attention. In Section 2.1, we describe the notations. In Section 2.2, we give definitions related to DP.

2.1 NOTATIONS

We use $\Pr[\cdot]$ to denote the probability. We use $\mathbb{E}[\cdot]$ to denote the expectation. We use $\text{Var}[\cdot]$ to denote the variance. For two vectors $x \in \mathbb{R}^d$ and $y \in \mathbb{R}^d$, we use $\langle x, y \rangle$ to denote the inner product between x, y , i.e., $\langle x, y \rangle = \sum_{i=1}^d x_i y_i$. We use $X \subset \mathbb{R}^d$ and $|X| = n$ to mean the same thing as $X \in \mathbb{R}^{n \times d}$. Also, we denote x_i^\top as the i -th row of X . We use $x_{i,j}$ to denote the j -th coordinate of $x_i \in \mathbb{R}^n$. We use $\mathbf{1}_n$ to denote a length- n vector where all the entries are ones. We use $\|x\|_p$ to denote the ℓ_p norm of a vector $x \in \mathbb{R}^n$, i.e., $\|x\|_1 := \sum_{i=1}^n |x_i|$, $\|x\|_2 := (\sum_{i=1}^n x_i^2)^{1/2}$, and $\|x\|_\infty := \max_{i \in [n]} |x_i|$. **We denote polynomial time complexity with respect to n as $\text{poly}(n)$. For a function f , we use $\tilde{O}(f)$ to represent f multiplied by a polylogarithmic factor, i.e., $f \cdot \text{poly}(\log f)$. This notation, known as soft- O or tilde notation, simplifies expressions by omitting logarithmic factors, focusing on the dominant term’s growth rate.**

2.2 DIFFERENTIAL PRIVACY DEFINITIONS

In this section, we give several definitions related to differential privacy (DP). We refer the reader to Dwork & Roth (2014) for more background and details on DP.

Definition 2.1 (Neighboring dataset). We define the two neighboring datasets as $X, X' \in \mathbb{R}^n$ such that $\|X - X'\|_1 \leq 1$, i.e., they differ on a single data point.

Definition 2.2 (Sensitivity). The sensitivity of a function $f : \mathbb{R}^n \rightarrow \mathbb{R}^d$ is defined by: $\Delta := \max_{X, X' \in \mathbb{R}^n, \|X - X'\|_1 = 1} \|f(X) - f(X')\|_1$.

Definition 2.3 ((ϵ, δ) -DP). For $\epsilon > 0, \delta \geq 0$, a randomized algorithm \mathcal{A} is (ϵ, δ) -DP, if for all $\mathcal{S} \subseteq \text{Range}(\mathcal{A})$ and for all X, X' such that $\|X - X'\|_1 \leq 1$:

$$\Pr[\mathcal{A}(X) \in \mathcal{S}] \leq \exp(\epsilon) \Pr[\mathcal{A}(X') \in \mathcal{S}] + \delta.$$

When $\delta = 0$, the algorithm is said to have pure differential privacy.

We mainly use the truncated Laplace mechanism, which has the following definitions.

Definition 2.4 (Truncated Laplace distribution). We use $\text{TLap}(\Delta, \epsilon, \delta)$ to denote the Truncated Laplace distribution with pdf proportional to $\exp(-\epsilon|z|/\Delta)$ on the region $[-B, B]$, where $B = \frac{\Delta}{\epsilon} \cdot \log(1 + \frac{\exp(\epsilon) - 1}{2\delta})$.

Fact 2.5 (Theorem 3 in Geng et al. (2020)). Let z denote a $\text{TLap}(\Delta, \epsilon, \delta)$ random variable. Then we have $\mathbb{E}[z] = 0$, and

$$\text{Var}[z] = \frac{2\Delta^2}{\epsilon^2} (1 - \delta \cdot \frac{\log^2(1 + \frac{\exp(\epsilon) - 1}{2\delta}) + 2\log(1 + \frac{\exp(\epsilon) - 1}{2\delta})}{e^\epsilon - 1}).$$

Furthermore, if $\delta = 0$, we have $\text{Var}[z] = 2\Delta^2/\epsilon^2$, meaning truncated Laplacian mechanism will be reduced to the standard Laplacian mechanism.

Lemma 2.6 (Laplace mechanism, (Dwork & Roth, 2014; Geng et al., 2020), see Lemma 2.2 in Andoni et al. (2023)). Given a numeric function f that takes a dataset X as the input, and has sensitivity Δ , the mechanism that outputs $f(X) + z$ where $z \sim \text{Lap}(\Delta/\epsilon)$ is $(\epsilon, 0)$ -DP. In addition, if $\epsilon, \delta \in (0, 0.5)$, $f(X) + z$, where $z \sim \text{TLap}(\Delta, \epsilon, \delta)$ is (ϵ, δ) -DP. Moreover, the truncated Laplace mechanism is always accuracy up to error B .

Algorithm 1 Adaptive query data structure

```

1: datastructure DPTREESOFTMAXADAPTIVE ▷ Theorem 4.4
2: members
3:    $\mathcal{D}_1, \dots, \mathcal{D}_{O(r \log(dR/(\epsilon_s p_f)))} : \text{DPTREESOFTMAX}$  ▷ Algorithm 3
4: end members
5: procedure INIT( $X \subset [0, R]^d, n \in \mathbb{N}_+, w \in [-R_w, R_w]^n, \epsilon \in (0, 1), \delta \in (0, 1), \delta' \in (0, 1), c \in (0, 0.1), \epsilon_s \in (0, 0.1), p_f \in (0, 0.01)$ )
6:    $l \leftarrow O(r \log(dR/(\epsilon_s p_f)))$ 
7:   for  $i = 1 \rightarrow l$  do
8:      $\mathcal{D}_i.\text{INIT}(X, n, w, \epsilon/l, \delta/l, \delta'/l, c, \epsilon_s)$ 
9:   end for
10: end procedure
11: procedure DISTANCEQUERY( $y \in [0, R]^d, \alpha \in (0, 1)$ )
12:    $l \leftarrow O(r \log(dR/(\epsilon_s p_f)))$ 
13:    $r \leftarrow 0^l$ 
14:   for  $i = 1 \rightarrow l$  do
15:      $r_i \leftarrow \mathcal{D}_i.\text{DISTANCEQUERY}(y, \alpha)$ 
16:   end for
17:   return Median of  $r$ 
18: end procedure
19: end datastructure

```

3 MAIN RESULTS: CROSS-ATTENTION

In this section, we show our main result for cross-attention. Theorem 3.1 states that we can ensure the entire cross-attention module satisfies DP and is robust to adaptive queries. Our high-level idea is based on the similarity between weighted distance problem and cross-attention. For a typical

weighted distance problem, we define the following: Let $w \in \mathbb{R}^n$ be the weights, $X \in \mathbb{R}^{n \times d}$ be the data matrix, where x_i^\top is the i -th row of X for $i \in [n]$, and let $y \in \mathbb{R}^d$ be the query. Suppose we need to answer ℓ_1 -distance query. We have

$$\sum_{i \in [n]} \underbrace{w_i}_{\text{weight}} \underbrace{\|y - x_i\|_1}_{\text{query data}}.$$

Now we introduce cross-attention. Let Q, K, V, Attn be defined in Definition 1.1. In a standard cross-attention process, K and V are accessible before inference, while the user input Q becomes available only when the user provides it. Here, K and V represent values stored in memory or disks and are considered private assets protected within the model, whereas Q is treated as public.

For the cross-attention mechanism Attn (Definition 1.1), we aim to ensure that the matrix AV satisfies DP guarantee. Let $A_{i,j} = \exp(\langle Q_i, K_j \rangle / d)$ for $i \in [m], j \in [n]$. Let $V_{j,k} \in \mathbb{R}$ be the (j, k) -th entry of V , for $j \in [n], k \in [d]$. Let $D = \text{diag}(A \mathbf{1}_n)$, acting as a normalizing factor that aggregates all the information. We store both K and its corresponding noises. For computing AV , we use the perturbed K , whereas for computing D , we rely on the original, unperturbed K . By post-processing property (Fact B.5), to ensure that the forward output $\text{Attn}(Q, K, V) = D^{-1}AV$ (Definition 1.1) satisfies DP, we only need to ensure the DP of its component AV .

The (i, k) -th entry of AV for each $i \in [m], k \in [d]$ is computed by

$$(AV)_{i,k} = \sum_{j=1}^n \underbrace{V_{j,k}}_{\text{weight}} \underbrace{\exp(\langle \underbrace{Q_i}_{\text{query}}, \underbrace{K_j}_{\text{data}} \rangle / d)}_{\text{weight}}, \quad (1)$$

which can be viewed as a weighted Softmax problem, where V provides the weights, Q is the query, and K is the dataset. Thus, we choose to add noise to K and V based on the similarity between the weighted distance problem and cross-attention. Furthermore, we find that we can only handle one column of V , i.e., $V_{*,k} \in \mathbb{R}^n$, in a single data structure. Therefore, we need to initialize a total of d different data structures, each with weights $V_{*,k}$ for $k \in [d]$.

Here, we present our main result below.

Theorem 3.1 (Softmax cross-attention, informal version of Theorem I.12). *Let Q, K, V, Attn be defined in Definition 1.1. Let $\alpha \in (0, 1)$ be the relative error parameter and p_f be the probability of failure parameter. Let r, s, ϵ_s be the parameters of the polynomial kernel methods (Lemma C.7). Let $\Gamma_{R,s} := \max_{j \in [s]} \frac{R_j^j}{\sqrt{j}}$ (Definition I.3). Let $l = O(r \log(dR/(\epsilon_s p_f)))$. There is a data structure $\text{DPTREESOFTMAXADAPTIVE}$ (Algorithm 1) that uses $O(\ln rd)$ spaces to ensure cross-attention satisfies DP and supports the following operations:*

- We initialize d data structures using $\text{INIT}(K, n, V_{*,k}, \epsilon \in (0, 1), \delta \in (0, 1), \delta' \in (0, 1), c \in (0, 0.1), \epsilon_s \in (0, 0.1), p_f \in (0, 0.01))$ (Algorithm 1), for $k \in [d]$. It takes $O(\ln r)$ time to initialize one data structure.
- At query time, for user input Q , we process one token at a time by passing the i -th row of Q , denoted $Q_i \in \mathbb{R}^d$, to $\text{DISTANCEQUERY}(Q_i, \alpha \in (0, 1))$ (Algorithm 1) for each $i \in [m]$. It takes $O(\alpha^{-1} l r \log^2 n)$ time to output an entry z in $\text{Attn}(Q, K, V)$ such that
 - the process of output z satisfies $(\epsilon, \delta + \delta')$ -DP,
 - the process of output z has relative error $n^{-1}(\alpha + \epsilon_s)$,
 - the process of output z has additive error $O(n^{-1} \epsilon^{-1} \alpha^{-1/2} l \Gamma_{R,s}^2 R_w r \sqrt{\log(l/\delta')} \cdot \log^{3/2} n)$,
 - it holds with probability $1 - p_f$ (where p_f is used in l),
 - it is robust to adaptive query.

In Theorem 3.1, we use our $\text{DPTREESOFTMAXADAPTIVE}$ (Algorithm 1) and guarantee that, for each query token of cross-attention, the output process satisfies (ϵ, δ) -DP with $n^{-1}(\alpha + \epsilon_s)$ relative error and $O(n^{-1} \epsilon^{-1} \alpha^{-1/2} l \Gamma_{R,s}^2 R_w r \sqrt{\log(l/\delta')} \cdot \log^{3/2} n)$ additive error, and $O(\alpha^{-1} l r \log^2 n)$ running time under adaptive query. More specifically, the algorithm creates d $\text{DPTREESOFTMAXADAPTIVE}$ data structures, each requiring $O(\ln r)$ memory consumption and $O(\ln r)$ initialization

time. Notably, our error is inversely proportional to n , meaning that as the input token length increases, both the relative and approximate errors approach zero. This is achieved by the normalizing matrix D (Definition 1.1). We refer the reader to Section I for proof details.

Thus, our algorithm theoretically protects system prompts/RAG data in cross-attention as discussed in Section 1. In Section 4, we provide a detailed technical overview, and in Section 5, we will present self-attention and DP-related discussion.

Algorithm 2 DPTree initialization and query

```

1: datastructure DPTREE ▷ Theorem C.1
2: members
3:    $b : \mathbb{R}^{2n-1}$ 
4:    $c : \mathbb{R}^{2n-1}$ 
5: end members
6: procedure INIT( $a \in \mathbb{R}^n, n \in \mathbb{N}_+, \Delta \in \mathbb{R}, \epsilon \in (0, 1), \delta \in (0, 1)$ ) ▷ Lemma D.3, Lemma C.3
7:    $b[n, 2n-1] \leftarrow a$ 
8:   for  $i = n \rightarrow 2n-1$  do
9:      $c[i] \leftarrow b[i] + \text{TLap}(\Delta, \epsilon/\log n, \delta/\log n)$ 
10:  end for
11:  for  $i = (\log n) \rightarrow 1$  do
12:    for  $j = 1 \rightarrow 2^{i-1}$  do
13:       $k \leftarrow 2^{i-1} + j - 1$ 
14:       $b[k] \leftarrow b[2k] + b[2k+1]$ 
15:       $c[k] \leftarrow b[k] + \text{TLap}(\Delta, \epsilon/\log n, \delta/\log n)$ 
16:    end for
17:  end for
18: end procedure
19: procedure QUERY( $x \in [n], y \in [n]$ ) ▷ Lemma D.4, D.5, D.6
20:   Trace from bottom nodes of  $x$  and  $y$  to find their lowest common ancestor, then we report
   the summation (based on  $c$ ) by using at most  $2 \log n$  nodes on the path. Let Value be the above
   summation.
21:   return Value
22: end procedure
23: procedure TRUEQUERY( $x \in [n], y \in [n]$ )
24:   Trace from bottom nodes of  $x$  and  $y$  to find their lowest common ancestor, then we report
   the summation (based on  $b$ ) by using at most  $2 \log n$  nodes on the path, where the height of the
   tree is  $\log n$ , and we need left and right boundary points. Let Value be the above summation.
25:   return Value
26: end procedure
27: end datastructure

```

4 KEY DATA STRUCTURE: DPTREE

This section provides our key data structures: DPTREE (Algorithm 2), DPTREEDISTANCE (Algorithm 5 and 6), DPTREEHIGHDIM (Algorithm 7), DPTREESOFTMAX (Algorithm 3), and DPTREESOFTMAXADAPTIVE (Algorithm 1).

In Section 4.1, we provide our high-level proof insights. In Section 4.2, we give our basic building block algorithms DPTREE, DPTREEDISTANCE and DPTREEHIGHDIM. In Section 4.3, we present our DPTREESOFTMAX algorithm that solves the weighted Softmax problem. In Section 4.4, we present our DPTREESOFTMAXADAPTIVE algorithm that enables DPTREESOFTMAX to handle adaptive query problem.

4.1 TECHNIQUE OVERVIEW

Notice that Eq. (1) is not a typical distance measure like ℓ_1 or ℓ_2 , but by using polynomial kernel method techniques, we transform it into a distance measure. Alman & Song (2023) states that the exponential inner product can be approximated by polynomial kernel function $P(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^r$,

i.e., $P(x)^\top P(y) \approx \exp(x^\top y/d)$ for two vector $x, y \in \mathbb{R}^d$, with a relative error. Then, by the Law of Cosines, we transform the inner product of polynomial kernel functions into a distance measure, i.e.,

$$2P(x)^\top P(y) = -\|P(x) - P(y)\|_2^2 + \|P(x)\|_2^2 + \|P(y)\|_2^2. \quad (2)$$

After transforming Eq. (1) into a distance measure, we design the DPTREE series data structures to provide cross-attention DP guarantee.

In summary, we first design the data structure DPTREE (Algorithm 2) that builds a binary segment tree with truncated Laplace noise added in the leaf nodes to ensure DP guarantee. Then, based on this data structure, we design DPTREEDISTANCE (Algorithm 5 and 6) to answer one dimensional weighted distance queries $\sum_{i=1}^n w_i \cdot |y - x_i|$, which utilizes DPTREE to store and return noised weights w_i multiplied with the approximated distances between the query y and data x_i . We further decompose high dimensional ℓ_p^p -distance problem into one dimensional ℓ_1 -distance problems using

$$\sum_{i=1}^n w_i \cdot \|y - x_i\|_p^p = \sum_{k=1}^d \sum_{i=1}^n w_i \cdot |y_k - x_{i,k}|^p. \quad (3)$$

Based on this decomposition, we design DPTREEHIGHDIM (Algorithm 7) which is capable of answering high dimension queries. Then, using Eq. (2) and DPTREEHIGHDIM, we design DPTREESOFTMAX (Algorithm 3) to answer Softmax queries. By building multiple copies of this data structure, we boost the success probability such that it can answer any query (including adaptive query) with an additive error, establishing the final data structure DPTREESOFTMAXADAPTIVE (Algorithm 1). See Section C for a more detailed outline of algorithms and proof techniques.

4.2 DPTREE, DPTREEDISTANCE, AND DPTREEHIGHDIM

The unweighted distance query has been explored in prior works (Huang & Roth, 2014; Backurs et al., 2024; Liu et al., 2024a). Specifically, Huang & Roth (2014) leverages online learning techniques to approximate the sum of distances, while Backurs et al. (2024) introduces a DP data structure based on a node-contaminated balanced binary tree. Furthermore, Liu et al. (2024a) presents a new data representation in tree nodes, where each node stores the sum of distances from one point to multiple points. In contrast, we focus on the weighted distance query, generalizing their results.

We design a basic data structure DPTREE (Algorithm 2) that answers summation queries by a summation segment tree with truncated Laplace noise (Definition 2.4). The algorithm first builds a binary summation tree in an array and then adds truncated Laplace noises to each node. In query time, we first trace from bottom nodes to find their lowest common ancestor, then report the summation by using at most $2 \log n$ nodes on the path (Algorithm 2). Based on the parallel composition rule of DP (Fact B.7), we find that if we have multiple disjoint interval queries, the error of the weighted sum of the intervals can be bounded independently of the number of queries (Lemma D.8). See more details in Section D.

We then design DPTREEDISTANCE, a one-dimensional weighted ℓ_1 distance data structure detailed in Algorithm 5 and 6. Initialization involves rounding each data point to the nearest multiple of a small interval and aggregating their weights into an array (illustrated in Figure 1), which is then input into our DPTREE. At query time, we retrieve aggregated weights within small intervals and multiply them by their distances to the query point. We introduce a relative error parameter α to reduce the number of iterations to $O(\log(n)/\alpha)$, improving efficiency. Guided by Eq.(3), we design DPTREEHIGHDIM (Algorithm 7), which extends DPTREEDISTANCE to higher dimension by constructing independent data structures for each coordinate. See details in Section F and G.

4.3 SOFTMAX ACTIVATION

In this section, we present DPTREESOFTMAX (Algorithm 3) that answers the weighted Softmax query (Definition 4.1) and is further used to design DP cross-attention. First, we introduce the definition of weighted Softmax query, an abstraction for the problem described in Eq. (1).

Definition 4.1 (Weighted Softmax query (without normalization)). *For the dataset $X \in [0, R]^{n \times d}$ where x_i^\top is the i -th row of X and query $y \in [0, R]^d$, we define the weighted exponential inner*

Algorithm 3 Softmax query

```

1: datastrucutre DPTREESOFTMAX ▷ Theorem 4.2
2: members
3:    $\mathcal{D}_1, \dots, \mathcal{D}_r : \text{DPTREEDISTANCE}$  ▷ Algorithm 5, Theorem I.7
4:    $P : [0, \Gamma_{R,s}]^{n \times r}$  ▷ Definition I.3 for  $\Gamma_{R,s}$ , Eq. (9) for  $s$ , Eq. (10) for  $r$ 
5:    $w : [-R_w, R_w]^n$ 
6:    $P_{wx}, s_w, \epsilon_s : \mathbb{R}$ 
7: end members
8: procedure INIT( $X \subset [0, R]^d, n \in \mathbb{N}_+, w \in [-R_w, R_w]^n, \epsilon \in (0, 1), \delta \in (0, 1), \delta' \in (0, 1),$   

    $c \in (0, 0.1), \epsilon_s \in (0, 0.1)$ ) ▷ Lemma C.7
9:    $\epsilon_s, w, P, P_{wx}, s_w \leftarrow \epsilon_s, w, 0^{n \times r}, 0, 0$ 
10:  for  $j = 1 \rightarrow n$  do
11:    Compute  $P(x_j)$  ▷ Polynomial kernel function  $P(\cdot)$ , Lemma I.5
12:    Compute  $w_j \|P(x_j)\|_2^2$ 
13:     $P_{wx} \leftarrow P_{wx} + w_j \|P(x_j)\|_2^2$ 
14:     $s_w \leftarrow s_w + w_j$ 
15:     $P_{j,:} \leftarrow P(x_j)$ 
16:  end for
17:  for  $i = 1 \rightarrow r$  do
18:     $\mathcal{D}_i.\text{INIT}(P_{:,i}, n, w, c\epsilon/\sqrt{r \log(1/\delta')}, \delta/r)$  ▷ Algorithm 5
19:  end for
20: end procedure
21: procedure DISTANCEQUERY( $y \in [0, R]^d, \alpha \in (0, 1)$ ) ▷ Lemma C.7
22:  Value  $\leftarrow 0$ 
23:  Compute  $P(y)$ 
24:  Compute  $\|P(y)\|_2^2$ 
25:  for  $i = 1 \rightarrow r$  do
26:    Value  $\leftarrow$  Value +  $\mathcal{D}_i.\text{DISTANCEQUERY}(P(y)_i, \alpha)$  ▷ Algorithm 6
27:  end for
28:  Value  $\leftarrow 0.5 \cdot (P_{wx} + s_w \|P(y)\|_2^2 - \text{Value})$ 
29:  return Value
30: end procedure
31: end datastrucutre

```

product/Softmax query to be:

$$\sum_{i \in [n]} w_i \exp(\langle x_i, y \rangle / d) = w^\top \exp(Xy/d).$$

Building on Definition 4.1, we develop a novel algorithm to answer differentially private weighted Softmax queries using the polynomial kernel method from Alman & Song (2023). Specifically, in Eq.(2), there is a term that computes the weighted ℓ_2^2 distance, which we calculate using DPTREEHIGHDIM. We then compute the exact term for the weighted ℓ_2^2 norms of the approximation kernel. By summing these terms with a controlled error, we extend DPTREEHIGHDIM to answer the Softmax query efficiently. More details can be found in Section I.

Theorem 4.2 (Softmax query, informal version of Theorem I.8). *Let $R \geq 1$. Let $r \leq \binom{2s+2d}{2s}$ and $s = O(\max\{\frac{\log(1/\epsilon_s)}{\log(\log(1/\epsilon_s)/R)}, R^2\})$. Let $\Gamma_{R,s} := \max_{j \in [s]} \frac{R^j}{\sqrt{j!}}$ (Definition I.3). Let the accuracy parameter be $\epsilon_s \in (0, 0.1)$. Our data structure DPTREESOFTMAX (Algorithm 3) uses $O(nr)$ spaces to solve Softmax query problem for dataset $X \subset [0, R]^d$ and support following operations:*

- INIT($X \subset [0, R]^d, n \in \mathbb{N}_+, w \in [-R_w, R_w]^n, \epsilon \in (0, 1), \delta \in (0, 1), \delta' \in (0, 1), c \in (0, 0.1), \epsilon_s \in (0, 0.1)$). (Algorithm 3) *It takes $O(nr)$ time to initialize the data structure.*
- DISTANCEQUERY($y \in [0, R]^d, \alpha \in (0, 1)$). (Algorithm 3) *It takes $O(\alpha^{-1} r \log^2 n)$ time to output a number z such that*
 - *the process of output z satisfies $(\epsilon, \delta + \delta')$ -DP private, which computes $w^\top \exp(Xy/d)$,*

- the error bound satisfies $|z - w^\top \exp(Xy/d)| \leq (\alpha + \epsilon_s) \cdot w^\top \exp(Xy/d) + O(\epsilon^{-1} \alpha^{-1/2} \Gamma_{R,s}^2 R_w r \sqrt{\log(1/\delta')} \cdot \log^{3/2} n)$,
- it holds with probability at least 0.99.

Remark 4.3. In Theorem 4.2, the parameter ϵ_s is the accuracy parameter for polynomial kernel approximation described in Section C.5. Besides, note that the error bound in Theorem 4.2 does not depend on δ but depends on δ' . The role of δ is to control a hidden constant term in the big O notation, i.e., increasing δ reduces the error by a small constant (Fact 2.5). In practice, we set δ as a small positive constant close to 0. Please refer to the Lemma D.6 for more details.

4.4 ADAPTIVE QUERY DATA STRUCTURE

We adapt our DPTREESOFTMAX to DPTREESOFTMAXADAPTIVE (Algorithm 1) to solve the adaptive query problem. By proving it can handle any query within the query space with a certain error, we ensure it effectively processes adaptive queries. We first boost the constant probability to high probability using the Chernoff bound (Lemma B.2). Employing an ϵ_0 -net argument and the union bound, we bound all query points within the net. Finally, we use the Lipschitz property of the weighted Softmax distance function with an additive error to bound all points in the query space. The corresponding proofs can be found in Section H and Section I.

Theorem 4.4 (Adaptive query Softmax data structure, informal version of Theorem I.11). *Let $R \geq 1$. Let $r \leq \binom{2s+2d}{2s}$ and $s = O(\max\{\frac{\log(1/\epsilon_s)}{\log(\log(1/\epsilon_s)/R)}, R^2\})$. Let $\Gamma_{R,s} := \max_{j \in [s]} \frac{R^j}{\sqrt{j!}}$ (Definition I.3). Let the accuracy parameter be $\epsilon_s \in (0, 0.1)$. Let $X \in [0, R]^{n \times d}$ be the dataset, $w \in [-R_w, R_w]^n$ be weights, $y \in [0, R]^d$ be the query, $\alpha \in (0, 1)$ be the relative error parameter and p_f be the failure probability parameter. Let $l = O(r \log(dR/(\epsilon_s p_f)))$. There is a data structure DPTREESOFTMAXADAPTIVE (Algorithm 1) that uses $O(\ln r)$ spaces to solve the weighted Softmax query problem for the dataset $X \subset [0, R]^d$ and supports the following operations:*

- **INIT**($X \subset [0, R]^d, n \in \mathbb{N}_+, w \in [-R_w, R_w]^n, \epsilon \in (0, 1), \delta \in (0, 1), \delta' \in (0, 1), c \in (0, 0.1), \epsilon_s \in (0, 0.1), p_f \in (0, 0.01)$). It takes $O(\ln r)$ time to initialize the data structure.
- **DISTANCEQUERY**($y \in [0, R]^d, \alpha \in (0, 1)$). It takes $O(\alpha^{-1} l r \log^2 n)$ time to output a number z such that
 - the process of output z satisfies $(\epsilon, \delta + \delta')$ -DP private, which computes $w^\top \exp(Xy/d)$,
 - the error bound satisfies $|z - w^\top \exp(Xy/d)| \leq (\alpha + \epsilon_s) \cdot w^\top \exp(Xy/d) + O(\epsilon^{-1} \alpha^{-1/2} l \Gamma_{R,s}^2 R_w r \sqrt{\log(l/\delta')} \cdot \log^{3/2} n)$,
 - it holds with probability at least $1 - p_f$ (where p_f is used in l),
 - it is robust to adaptive query.

Remark 4.5. We describe the parallelization of our algorithms. In the second for loop of **INIT** and the for loop of **DISTANCEQUERY** in Algorithm 3, the r DPTREEDISTANCE data structures instantiated for each coordinate are independent of each other. In addition, the for loops in Algorithm 1 are also parallelizable since the $l = O(r \log(dR/(\epsilon_s p_f)))$ copies are independent. After parallelization, we have the final time complexity of **INIT** to be $O(nr)$ and **DISTANCEQUERY** to be $O(\alpha^{-1} \log^2 n)$ in Algorithm 1 with $O(lr)$ GPU process.

5 DISCUSSION

How do we extend to self-attention and other data structures? As self-attention is a more fundamental module in LGMs, we would like to extend our data structure to this setting. However, the challenge we faced was the dynamic update in tree nodes for each query for self-attention, which our current analysis does not support. How we can solve this challenge is crucial, and we leave it as our future direction.

Moreover, we observe that Li et al. (2015) introduces the DP matrix mechanism, which offers an alternative to our currently used binary tree data structure. A preliminary idea for extending this is as follows: consider $A = \exp(QK^\top/d)$ as defined in Definition 1.1, where Q of size $m \times d$ represents the query matrix with m linear queries, and K serves as the database. Leveraging the

results from Li et al. (2015), we could design an alternative algorithm to enhance the current binary tree data structure, DPTREE. We leave this exploration for future work.

Why not add noise to some other places? Where and how to add DP noises is an important problem to ask during the DP algorithm design. In this paper, we consider the problem of $\sum_{i=1}^n w_i \exp(\langle x_i, y \rangle / d)$ where $y, x_i \in [0, R]^d$ and $w \in [-R_w, R_w]^n$ (Definition 4.1). Notice that the only place where we add noises is in the most basic building block data structure DPTREE (Algorithm 2). From Lemma C.3 and the way we initialize DPTREE in Algorithm 5, we see that the sensitivity Δ of this problem is $2R_w$.

A simple method for adding noise involves adding n noises to a length n array, with each item $w_i \exp(\langle x_i, y \rangle / d)$ for $i \in [n]$. However, this approach increases the error by a factor of n by basic composition (Fact B.6) and also makes the model dependent on the number of queries. Besides, it only supports a single query and requires rebuilding the tree for each new query, rendering it impractical. In contrast, our current noise-adding technique (Lines 9 and 15 of Algorithm 2) utilizes a summation tree such that the error only increases by a factor of $\text{poly} \log n$. This method also supports multiple queries, eliminating the need to rebuild the tree each time.

How to remove the relative error parameter α ? The relative error parameter α in Theorem 3.1 appears because of the $(1 + \alpha)$ -approximation introduced in Algorithm 5 (Remark F.3) to reduce the number of required iterations from naive $O(n)$ to $O(\log(n)/\alpha)$. However, we notice that a recent work (Liu et al., 2024a) does not utilize $(1 + \alpha)$ -approximation and still achieves $O(\log n)$ iteration number. They introduce a new tree node representation where each node stores the sum of distances from one point to multiple points, enabling the answer to be divided into only $\log n$ values, each combining two distance values, two count values, and y itself. Our DPTREE algorithms can be integrated with their method, thus removing parameter α .

6 CONCLUSION

To our knowledge, we are the first work to provide differential privacy for cross-attention. This paper presents the DPTREE data structures, which provide a differential privacy guarantee for the cross-attention module in large generative models. This is achieved by transforming the cross-attention mechanism into a weighted distance problem. Furthermore, our algorithm is robust to adaptive queries, allowing users to interact with the model arbitrarily without extracting sensitive information from the system prompts or RAG data. Our results may inspire more privacy algorithm design in large generative models.

REFERENCES

- Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC conference on computer and communications security*, pp. 308–318, 2016.
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Josh Alman and Zhao Song. Fast attention requires bounded entries. *Advances in Neural Information Processing Systems*, 36, 2023.
- Josh Alman and Zhao Song. The fine-grained complexity of gradient computation for training large language models. *arXiv preprint arXiv:2402.04497*, 2024a.
- Josh Alman and Zhao Song. How to capture higher-order correlations? generalizing matrix softmax attention to kronecker computation. In *The Twelfth International Conference on Learning Representations*, 2024b.
- Alexandr Andoni, Piotr Indyk, Sepideh Mahabadi, and Shyam Narayanan. Differentially private approximate near neighbor counting in high dimensions. In *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 43544–43562, 2023.

- Anthropic. System prompts, 2024. <https://docs.anthropic.com/en/docs/system-prompts>.
- Apple. Apple intelligence, 2024. <https://www.apple.com/apple-intelligence/>.
- Arturs Backurs, Zinan Lin, Sepideh Mahabadi, Sandeep Silwal, and Jakub Tarnawski. Efficiently computing similarities to private datasets. *arXiv preprint arXiv:2403.08917*, 2024.
- Rouzbeh Behnia, Mohammadreza Reza Ebrahimi, Jason Pacheco, and Balaji Padmanabhan. Ew-tune: A framework for privately fine-tuning large language models with differential privacy. In *2022 IEEE International Conference on Data Mining Workshops (ICDMW)*, pp. 560–566. IEEE, 2022.
- Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George Bm Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al. Improving language models by retrieving from trillions of tokens. In *International conference on machine learning*, pp. 2206–2240. PMLR, 2022.
- Chun-Fu Richard Chen, Quanfu Fan, and Rameswar Panda. Crossvit: Cross-attention multi-scale vision transformer for image classification. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 357–366, 2021.
- Justin Y Chen, Shyam Narayanan, and Yinzhan Xu. All-pairs shortest path distances with differential privacy: Improved algorithms for bounded and unbounded weights. *arXiv preprint arXiv:2204.02335*, 2022.
- Yeshwanth Cherapanamjeri, Sandeep Silwal, David P Woodruff, Fred Zhang, Qiuyi Zhang, and Samson Zhou. Robust algorithms on adaptive inputs from bounded adversaries. *arXiv preprint arXiv:2304.07413*, 2023.
- Herman Chernoff. A measure of asymptotic efficiency for tests of a hypothesis based on the sum of observations. *The Annals of Mathematical Statistics*, pp. 493–507, 1952.
- Vincent Cohen-Addad, Alessandro Epasto, Vahab Mirrokni, Shyam Narayanan, and Peilin Zhong. Near-optimal private and scalable k -clustering. *Advances in Neural Information Processing Systems*, 35:10462–10475, 2022a.
- Vincent Cohen-Addad, Chenglin Fan, Silvio Lattanzi, Slobodan Mitrovic, Ashkan Norouzi-Fard, Nikos Parotsidis, and Jakub M Tarnawski. Near-optimal correlation clustering with privacy. *Advances in Neural Information Processing Systems*, 35:33702–33715, 2022b.
- Itai Dinur, Uri Stemmer, David P Woodruff, and Samson Zhou. On differential privacy and adaptive data analysis with bounded space. In *Annual International Conference on the Theory and Applications of Cryptographic Techniques*, pp. 35–65. Springer, 2023.
- Wei Dong, Zijun Chen, Qiyao Luo, Elaine Shi, and Ke Yi. Continual observation of joins under differential privacy. *Proceedings of the ACM on Management of Data*, 2(3):1–27, 2024.
- Cynthia Dwork. Differential privacy: A survey of results. In *International conference on theory and applications of models of computation*, pp. 1–19. Springer, 2008.
- Cynthia Dwork and Aaron Roth. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science*, 9(3–4):211–407, 2014.
- Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. In *Theory of Cryptography: Third Theory of Cryptography Conference, TCC 2006, New York, NY, USA, March 4-7, 2006. Proceedings 3*, pp. 265–284. Springer, 2006.
- Kennedy Edemacu and Xintao Wu. Privacy preserving prompt engineering: A survey. *arXiv preprint arXiv:2404.06001*, 2024.
- Marek Eliáš, Michael Kapralov, Janardhan Kulkarni, and Yin Tat Lee. Differentially private release of synthetic graphs. In *Proceedings of the Fourteenth Annual ACM-SIAM Symposium on Discrete Algorithms*, pp. 560–578. SIAM, 2020.

- Alessandro Epasto, Vahab Mirrokni, Shyam Narayanan, and Peilin Zhong. k -means clustering with distance-based privacy. *Advances in Neural Information Processing Systems*, 36, 2024.
- Hossein Esfandiari, Vahab Mirrokni, and Shyam Narayanan. Tight and robust private mean estimation with few users. In *International Conference on Machine Learning*, pp. 16383–16412. PMLR, 2022.
- Chenglin Fan and Ping Li. Distances release with differential privacy in tree and grid graph. In *2022 IEEE International Symposium on Information Theory (ISIT)*, pp. 2190–2195. IEEE, 2022.
- Chenglin Fan, Ping Li, and Xiaoyun Li. k -median clustering via metric embedding: towards better initialization with differential privacy. *Advances in Neural Information Processing Systems*, 36, 2024.
- Linxi Fan, Guanzhi Wang, Yunfan Jiang, Ajay Mandlekar, Yuncong Yang, Haoyi Zhu, Andrew Tang, De-An Huang, Yuke Zhu, and Anima Anandkumar. Minedojo: Building open-ended embodied agents with internet-scale knowledge. *Advances in Neural Information Processing Systems*, 35: 18343–18362, 2022.
- Alireza Farhadi, MohammadTaghi Hajiaghayi, and Elaine Shi. Differentially private densest subgraph. In *International Conference on Artificial Intelligence and Statistics*, pp. 11581–11597. PMLR, 2022.
- Yeqi Gao, Zhao Song, Xin Yang, and Yufa Zhou. Differentially private attention computation. In *Neurips Safe Generative AI Workshop 2024*, 2024.
- Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen Wang. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997*, 2023.
- Quan Geng, Wei Ding, Ruiqi Guo, and Sanjiv Kumar. Tight analysis of privacy and utility tradeoff in approximate differential privacy. In *International Conference on Artificial Intelligence and Statistics*, pp. 89–99. PMLR, 2020.
- Badih Ghazi, Pritish Kamath, Ravi Kumar, Pasin Manurangsi, and Kewen Wu. On differentially private counting on trees. In *50th International Colloquium on Automata, Languages, and Programming (ICALP 2023)*, volume 261, pp. 66. Schloss Dagstuhl–Leibniz-Zentrum für Informatik, 2023.
- GitHub. Agentgpt, 2024a. <https://github.com/reworkd/AgentGPT>.
- GitHub. Chatgpt system prompt, 2024b. https://github.com/LouisShark/chatgpt_system_prompt.
- GitHub. Gpt researcher, 2024c. <https://github.com/assafelovic/gpt-researcher>.
- GitHub. Superagi, 2024d. <https://github.com/TransformerOptimus/SuperAGI>.
- Aditya Golatkar, Alessandro Achille, Yu-Xiang Wang, Aaron Roth, Michael Kearns, and Stefano Soatto. Mixed differential privacy in computer vision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8376–8386, 2022.
- Sivakanth Gopi, Yin Tat Lee, and Lukas Wutschitz. Numerical composition of differential privacy. *Advances in Neural Information Processing Systems*, 34:11631–11642, 2021.
- Sivakanth Gopi, Yin Tat Lee, and Daogao Liu. Private convex optimization via exponential mechanism. In *Conference on Learning Theory*, pp. 1948–1989. PMLR, 2022.
- Sivakanth Gopi, Yin Tat Lee, Daogao Liu, Ruqi Shen, and Kevin Tian. Private convex optimization in general norms. In *Proceedings of the 2023 Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pp. 5068–5089. SIAM, 2023.

- Adityanand Guntuboyina and Bodhisattva Sen. L1 covering numbers for uniformly bounded convex functions. In *Conference on Learning Theory*, pp. 12–1. JMLR Workshop and Conference Proceedings, 2012.
- Michael Hay, Vibhor Rastogi, Gerome Miklau, and Dan Suciu. Boosting the accuracy of differentially-private histograms through consistency. *arXiv preprint arXiv:0904.0942*, 2009.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*, 2022.
- Sirui Hong, Yizhang Lin, Bangbang Liu, Binhao Wu, Danyang Li, Jiaqi Chen, Jiayi Zhang, Jinlin Wang, Lingyao Zhang, Mingchen Zhuge, et al. Data interpreter: An llm agent for data science. *arXiv preprint arXiv:2402.18679*, 2024a.
- Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu Zheng, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, Chenyu Ran, Lingfeng Xiao, Chenglin Wu, and Jürgen Schmidhuber. MetaGPT: Meta programming for a multi-agent collaborative framework. In *The Twelfth International Conference on Learning Representations*, 2024b. URL <https://openreview.net/forum?id=VtmBAGCN7o>.
- Samuel B Hopkins, Gautam Kamath, Mahbod Majid, and Shyam Narayanan. Robustness implies privacy in statistical estimation. In *Proceedings of the 55th Annual ACM Symposium on Theory of Computing*, pp. 497–506, 2023.
- Jerry Yao-Chieh Hu, Donglin Yang, Dennis Wu, Chenwei Xu, Bo-Yu Chen, and Han Liu. On sparse modern hopfield model. In *Thirty-seventh Conference on Neural Information Processing Systems (NeurIPS)*, 2023.
- Jerry Yao-Chieh Hu, Pei-Hsuan Chang, Haozheng Luo, Hong-Yu Chen, Weijian Li, Wei-Po Wang, and Han Liu. Outlier-efficient hopfield layers for large transformer-based models. In *Forty-first International Conference on Machine Learning (ICML)*, 2024a.
- Jerry Yao-Chieh Hu, Bo-Yu Chen, Dennis Wu, Feng Ruan, and Han Liu. Nonparametric modern hopfield models. *arXiv preprint arXiv:2404.03900*, 2024b.
- Jerry Yao-Chieh Hu, Thomas Lin, Zhao Song, and Han Liu. On computational limits of modern hopfield models: A fine-grained complexity analysis. In *Forty-first International Conference on Machine Learning (ICML)*, 2024c.
- Tianhao Huang, Tao Yang, Ivan Habernal, Lijie Hu, and Di Wang. Private language models via truncated laplacian mechanism. *arXiv preprint arXiv:2410.08027*, 2024.
- Zhiyi Huang and Aaron Roth. Exploiting metric structure for efficient private query release. In *Proceedings of the twenty-fifth annual ACM-SIAM symposium on Discrete algorithms*, pp. 523–534. SIAM, 2014.
- Ziyue Huang and Ke Yi. Approximate range counting under differential privacy. In *37th International Symposium on Computational Geometry (SoCG 2021)*. Schloss-Dagstuhl-Leibniz Zentrum für Informatik, 2021.
- Haibo Jin, Leyang Hu, Xinuo Li, Peiyan Zhang, Chonghan Chen, Jun Zhuang, and Haohan Wang. Jailbreakzoo: Survey, landscapes, and horizons in jailbreaking large language and vision-language models. *arXiv preprint arXiv:2407.01599*, 2024.
- Christopher Jung, Katrina Ligett, Seth Neel, Aaron Roth, Saeed Sharifi-Malvajerdi, and Moshe Shenfeld. A new analysis of differential privacy’s generalization guarantees. *arXiv preprint arXiv:1909.03577*, 2019.

- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33: 9459–9474, 2020.
- Chao Li, Gerome Miklau, Michael Hay, Andrew McGregor, and Vibhor Rastogi. The matrix mechanism: optimizing linear counting queries under differential privacy. *The VLDB journal*, 24: 757–781, 2015.
- Chenyang Li, Yingyu Liang, Zhenmei Shi, Zhao Song, and Tianyi Zhou. Fourier circuits in neural networks: Unlocking the potential of large language models in mathematical reasoning and modular arithmetic. *arXiv preprint arXiv:2402.09469*, 2024a.
- Haoran Li, Yulin Chen, Jinglong Luo, Yan Kang, Xiaojin Zhang, Qi Hu, Chunkit Chan, and Yangqiu Song. Privacy in large language models: Attacks, defenses and future directions. *arXiv preprint arXiv:2310.10383*, 2023a.
- Haoran Li, Dadi Guo, Wei Fan, Mingshi Xu, Jie Huang, Fanpu Meng, and Yangqiu Song. Multi-step jailbreaking privacy attacks on chatgpt. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 4138–4153, 2023b.
- Ninghui Li, Min Lyu, Dong Su, and Weining Yang. *Differential privacy: From theory to practice*. Springer, 2017.
- Ping Li and Xiaoyun Li. Differential privacy with random projections and sign random projections. *arXiv preprint arXiv:2306.01751*, 2023a.
- Ping Li and Xiaoyun Li. Smooth flipping probability for differential private sign random projection methods. *Advances in Neural Information Processing Systems*, 36, 2024.
- Xiaoyu Li, Yingyu Liang, Zhenmei Shi, Zhao Song, and Junwei Yu. Fast john ellipsoid computation with differential privacy optimization. *arXiv preprint arXiv:2408.06395*, 2024b.
- Xiaoyun Li and Ping Li. Differentially private one permutation hashing and bin-wise consistent weighted sampling. *arXiv preprint arXiv:2306.07674*, 2023b.
- Xuechen Li, Florian Tramer, Percy Liang, and Tatsunori Hashimoto. Large language models can be strong differentially private learners. In *International Conference on Learning Representations*, 2021.
- Xuechen Li, Daogao Liu, Tatsunori B Hashimoto, Huseyin A Inan, Janardhan Kulkarni, Yin-Tat Lee, and Abhradeep Guha Thakurta. When does differentially private learning not suffer in high dimensions? *Advances in Neural Information Processing Systems*, 35:28616–28630, 2022.
- Yingyu Liang, Zhizhou Sha, Zhenmei Shi, and Zhao Song. Differential privacy mechanisms in neural tangent kernel regression. *arXiv preprint arXiv:2407.13621*, 2024a.
- Yingyu Liang, Zhizhou Sha, Zhenmei Shi, Zhao Song, and Yufa Zhou. Multi-layer transformers gradient can be approximated in almost linear time. *arXiv preprint arXiv:2408.13233*, 2024b.
- Yingyu Liang, Zhenmei Shi, Zhao Song, and Chiwun Yang. Toward infinite-long prefix in transformer. *arXiv preprint arXiv:2406.14036*, 2024c.
- Yingyu Liang, Zhenmei Shi, Zhao Song, and Yufa Zhou. Unraveling the smoothness properties of diffusion models: A gaussian mixture perspective. *arXiv preprint arXiv:2405.16418*, 2024d.
- Yingyu Liang, Zhenmei Shi, Zhao Song, and Yufa Zhou. Tensor attention training: Provably efficient learning of higher-order transformers. *arXiv preprint arXiv:2405.16411*, 2024e.
- Erzhi Liu, Jerry Yao-Chieh Hu, Alex Reneau, Zhao Song, and Han Liu. Differentially private kernel density estimation. *arXiv preprint arXiv:2409.01688*, 2024a.
- Zhihao Liu, Jian Lou, Wenjie Bao, Zhan Qin, and Kui Ren. Differentially private zeroth-order methods for scalable large language model finetuning. *arXiv preprint arXiv:2402.07818*, 2024b.

- Zhiwei Liu, Weiran Yao, Jianguo Zhang, Le Xue, Shelby Heinecke, Rithesh Murthy, Yihao Feng, Zeyuan Chen, Juan Carlos Niebles, Devansh Arpit, et al. Bolaa: Benchmarking and orchestrating llm-augmented autonomous agents. *arXiv preprint arXiv:2308.05960*, 2023.
- Peihua Mai, Ran Yan, Zhe Huang, Youjia Yang, and Yan Pang. Split-and-denoise: Protect large language model inference with local differential privacy. *arXiv preprint arXiv:2310.09130*, 2023.
- Justus Mattern, Zhijing Jin, Benjamin Weggenmann, Bernhard Schölkopf, and Mrinmaya Sachan. Differentially private language models for secure data sharing. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 4860–4873. Association for Computational Linguistics, 2022.
- Shyam Narayanan. Private high-dimensional hypothesis testing. In *Conference on Learning Theory*, pp. 3979–4027. PMLR, 2022.
- Shyam Narayanan. Better and simpler lower bounds for differentially private statistical estimation. *arXiv preprint arXiv:2310.06289*, 2023.
- OpenAI. Creating a gpt, 2024a. <https://help.openai.com/en/articles/8554397-creating-a-gpt>.
- OpenAI. Video generation models as world simulators, 2024b. <https://openai.com/research/video-generation-models-as-world-simulators>.
- Samet Oymak, Ankit Singh Rawat, Mahdi Soltanolkotabi, and Christos Thrampoulidis. On the role of attention in prompt-tuning. In *International Conference on Machine Learning*, pp. 26724–26768. PMLR, 2023.
- Dylan Patel. Chatgpt system prompt is 1700 tokens?!, 2024. <https://x.com/dylan522p/status/1755086111397863777>.
- William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4195–4205, 2023.
- Natalia Ponomareva, Hussein Hazimeh, Alex Kurakin, Zheng Xu, Carson Denison, H Brendan McMahan, Sergei Vassilvitskii, Steve Chien, and Abhradeep Guha Thakurta. How to dp-fy ml: A practical guide to machine learning with differential privacy. *Journal of Artificial Intelligence Research*, 77:1113–1201, 2023.
- Lianke Qin, Aravind Reddy, Zhao Song, Zhaozhuo Xu, and Danyang Zhuo. Adaptive and dynamic multi-resolution hashing for pairwise summations. In *2022 IEEE International Conference on Big Data (Big Data)*, pp. 115–120. IEEE, 2022.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in neural information processing systems*, 35:36479–36494, 2022.
- Sander Schulhoff, Michael Ilie, Nishant Balepur, Konstantine Kahadze, Amanda Liu, Chenglei Si, Yinheng Li, Aayush Gupta, Hyojung Han, Sevien Schulhoff, et al. The prompt report: A systematic survey of prompting techniques. *arXiv preprint arXiv:2406.06608*, 2024.
- Weiyang Shi, Ryan Shea, Si Chen, Chiyuan Zhang, Ruoxi Jia, and Zhou Yu. Just fine-tune twice: Selective differential privacy for large language models. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 6327–6340, 2022a.
- Zhenmei Shi, Jiefeng Chen, Kunyang Li, Jayaram Raghuram, Xi Wu, Yingyu Liang, and Somesh Jha. The trade-off between universality and label efficiency of representations from contrastive learning. In *The Eleventh International Conference on Learning Representations*, 2022b.

- Zhenmei Shi, Junyi Wei, Zhuoyan Xu, and Yingyu Liang. Why larger language models do in-context learning differently? *arXiv preprint arXiv:2405.19592*, 2024.
- Tanmay Singh, Harshvardhan Aditya, Vijay K Madiseti, and Arshdeep Bahga. Whispered tuning: Data privacy preservation in fine-tuning llms through differential privacy. *Journal of Software Engineering and Applications*, 17(1):1–22, 2024.
- Zhao Song, Yitan Wang, Zheng Yu, and Lichen Zhang. Sketching for first order method: efficient algorithm for low-bandwidth channel and vulnerability. In *International Conference on Machine Learning*, pp. 32365–32417. PMLR, 2023a.
- Zhao Song, Xin Yang, Yuanyuan Yang, and Lichen Zhang. Sketching meets differential privacy: fast algorithm for dynamic kronecker projection maintenance. In *International Conference on Machine Learning (ICML)*, pp. 32418–32462. PMLR, 2023b.
- 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.
- Roman Vershynin. An introduction with applications in data science. *Camb. Ser. Stat. Probab. Math.*, 47, 2017.
- Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. *arXiv preprint arXiv:2305.16291*, 2023a.
- Jiayu Wang, Yifei Ming, Zhenmei Shi, Vibhav Vineet, Xin Wang, and Neel Joshi. Is a picture worth a thousand words? delving into spatial reasoning for vision language models. *arXiv preprint arXiv:2406.14852*, 2024a.
- Yilin Wang, Zeyuan Chen, Liangjun Zhong, Zheng Ding, Zhizhou Sha, and Zhuowen Tu. Dolphin: Diffusion layout transformers without autoencoder. *arXiv preprint arXiv:2310.16305*, 2023b.
- Yilin Wang, Haiyang Xu, Xiang Zhang, Zeyuan Chen, Zhizhou Sha, Zirui Wang, and Zhuowen Tu. Omnicontrolnet: Dual-stage integration for conditional image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7436–7448, 2024b.
- Zirui Wang, Zhizhou Sha, Zheng Ding, Yilin Wang, and Zhuowen Tu. Tokencompose: Grounding diffusion with token-level supervision. *arXiv preprint arXiv:2312.03626*, 2023c.
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682*, 2022.
- David Woodruff, Fred Zhang, and Samson Zhou. On robust streaming for learning with experts: algorithms and lower bounds. *Advances in Neural Information Processing Systems*, 36:79518–79539, 2023.
- David P Woodruff. Sketching as a tool for numerical linear algebra. *Foundations and Trends® in Theoretical Computer Science*, 10(1–2):1–157, 2014.
- Dennis Wu, Jerry Yao-Chieh Hu, Teng-Yun Hsiao, and Han Liu. Uniform memory retrieval with larger capacity for modern hopfield models. In *Forty-first International Conference on Machine Learning (ICML)*, 2024a.
- Dennis Wu, Jerry Yao-Chieh Hu, Weijian Li, Bo-Yu Chen, and Han Liu. STanhop: Sparse tandem hopfield model for memory-enhanced time series prediction. In *The Twelfth International Conference on Learning Representations (ICLR)*, 2024b.
- Chenwei Xu, Yu-Chao Huang, Jerry Yao-Chieh Hu, Weijian Li, Ammar Gilani, Hsi-Sheng Goan, and Han Liu. Bishop: Bi-directional cellular learning for tabular data with generalized sparse modern hopfield model. In *Forty-first International Conference on Machine Learning (ICML)*, 2024a.

- Zhuoyan Xu, Zhenmei Shi, Junyi Wei, Fangzhou Mu, Yin Li, and Yingyu Liang. Towards few-shot adaptation of foundation models via multitask finetuning. In *The Twelfth International Conference on Learning Representations*, 2023.
- Zhuoyan Xu, Zhenmei Shi, and Yingyu Liang. Do large language models have compositional ability? an investigation into limitations and scalability. In *ICLR 2024 Workshop on Mathematical and Empirical Understanding of Foundation Models*, 2024b.
- Fei Yang, Shiqi Yang, Muhammad Atif Butt, Joost van de Weijer, et al. Dynamic prompt learning: Addressing cross-attention leakage for text-based image editing. *Advances in Neural Information Processing Systems*, 36, 2024.
- Mengmeng Yang, Taolin Guo, Tianqing Zhu, Ivan Tjuawinata, Jun Zhao, and Kwok-Yan Lam. Local differential privacy and its applications: A comprehensive survey. *Computer Standards & Interfaces*, pp. 103827, 2023.
- Da Yu, Saurabh Naik, Arturs Backurs, Sivakanth Gopi, Huseyin A Inan, Gautam Kamath, Janardhan Kulkarni, Yin Tat Lee, Andre Manoel, Lukas Wutschitz, et al. Differentially private fine-tuning of language models. In *International Conference on Learning Representations*, 2021.
- Jiahao Yu, Haozheng Luo, Jerry Yao-Chieh Hu, Wenbo Guo, Han Liu, and Xinyu Xing. Enhancing jailbreak attack against large language models through silent tokens. *arXiv preprint arXiv:2405.20653*, 2024.
- Ying Zhao and Jinjun Chen. A survey on differential privacy for unstructured data content. *ACM Computing Surveys (CSUR)*, 54(10s):1–28, 2022.
- Chunyan Zheng, Keke Sun, Wenhao Zhao, Haibo Zhou, Lixing Jiang, Shaoyang Song, and Chunlai Zhou. Locally differentially private in-context learning. In *LREC/COLING*, 2024.

Appendix

CONTENTS

| | | |
|----------|---|-----------|
| 1 | Introduction | 1 |
| 1.1 | Related Work | 3 |
| 2 | Preliminary | 3 |
| 2.1 | Notations | 3 |
| 2.2 | Differential Privacy Definitions | 3 |
| 3 | Main Results: Cross-Attention | 4 |
| 4 | Key Data Structure: DPTree | 6 |
| 4.1 | Technique Overview | 6 |
| 4.2 | DPTree, DPTreeDistance, and DPTreeHighDim | 7 |
| 4.3 | Softmax Activation | 7 |
| 4.4 | Adaptive Query Data Structure | 9 |
| 5 | Discussion | 9 |
| 6 | Conclusion | 10 |
| A | More Related Work | 20 |
| B | More Preliminary | 21 |
| B.1 | Probability Tools | 21 |
| B.2 | Algebraic Facts | 21 |
| B.3 | DP Facts | 21 |
| B.4 | Comparison of Truncated Laplace, Gaussian, and Laplace Mechanisms | 22 |
| C | Proof Outline | 22 |
| C.1 | Summation Segment Tree | 22 |
| C.2 | Sensitivity for Range Summation Problem | 23 |
| C.3 | Weighted ℓ_p^p Distance Problem | 23 |
| C.4 | One-Dimensional Weighted ℓ_1 Distance Data Structure | 24 |
| C.5 | Softmax Activation | 25 |
| C.6 | Adaptive Query | 25 |
| D | DPTree Algorithm | 26 |
| D.1 | Single Data Structure | 26 |
| D.2 | Boost the Constant Probability to High Probability | 26 |
| D.3 | Algorithm of Data Structure | 27 |

| | | |
|------|--|-----------|
| 972 | D.4 Disjoint Intervals | 28 |
| 973 | | |
| 974 | E Weighted ℓ_p^p Distance | 30 |
| 975 | | |
| 976 | E.1 One Dimensional Weighted Distance | 30 |
| 977 | E.2 High Dimensional Weighted Distance | 31 |
| 978 | | |
| 979 | F One-Dimensional Weighted ℓ_1 Distance Query | 32 |
| 980 | | |
| 981 | F.1 Runtime Analysis | 32 |
| 982 | F.2 Privacy and Accuracy Analysis | 33 |
| 983 | F.3 One Dimension Single Data Structure | 35 |
| 984 | | |
| 985 | | |
| 986 | G High-Dimensional Weighted ℓ_1 Query | 35 |
| 987 | | |
| 988 | G.1 Privacy and Accuracy Analysis for High Dimensional Weighted Distance | 35 |
| 989 | G.2 High Dimension Single Data Structure | 37 |
| 990 | | |
| 991 | H Adaptive Query | 37 |
| 992 | | |
| 993 | H.1 Boost the Constant Probability to High Probability | 37 |
| 994 | H.2 From Each Fixed Query Point to All On-net Points | 38 |
| 995 | H.3 From Net Points to All Points | 38 |
| 996 | H.4 Effect of Different Norms on the Result | 40 |
| 997 | | |
| 998 | | |
| 999 | I Softmax Activation | 42 |
| 1000 | | |
| 1001 | I.1 Exponential Inner Product | 42 |
| 1002 | I.2 Algorithm Modifications | 43 |
| 1003 | I.3 Adaptive Softmax | 44 |
| 1004 | I.4 Proof of Main Result | 46 |
| 1005 | | |
| 1006 | | |
| 1007 | | |
| 1008 | | |
| 1009 | | |
| 1010 | | |
| 1011 | | |
| 1012 | | |
| 1013 | | |
| 1014 | | |
| 1015 | | |
| 1016 | | |
| 1017 | | |
| 1018 | | |
| 1019 | | |
| 1020 | | |
| 1021 | | |
| 1022 | | |
| 1023 | | |
| 1024 | | |
| 1025 | | |

Roadmap. The appendix is organized as follows. In Section A, we provide more related works. In Section B, we give the preliminary of our paper. In Section C, we offer an outline of our proof techniques. In Section D, we give the analysis of the data structure DPTREE that can solve summation problem with DP and accuracy guarantee. In Section E, we show how to solve weighted distance problem. In Section F, we give our DPTREEDISTANCE data structure that can solve one dimensional ℓ_1 distance problem with DP and accuracy guarantee. In Section G, we present the analysis of our DPTREEHIGHDIM (Algorithm 7) data structure, which can address the high-dimensional ℓ_1 distance problem while ensuring differential privacy and accuracy guarantees. In Section H, we show how we can handle adaptive query. In Section I, we show how to extend our algorithm to Softmax activation and give the analysis of DPTREESOFTMAX (Algorithm 3) and DPTREESOFTMAXADAPTIVE (Algorithm 1).

A MORE RELATED WORK

Differential Privacy Guarantee Analysis. Ever since Dwork et al. (2006) proposes the notion of differential privacy (DP), it has become one of the most essential standards of privacy protection in both theoretical and empirical ways (Dwork, 2008; Li et al., 2017; Zhao & Chen, 2022; Ponomareva et al., 2023; Yang et al., 2023). DP provides a powerful, robust, and quantifiable privacy definition, allowing algorithm design with concrete privacy and accuracy guarantee (Hay et al., 2009; Esfandiari et al., 2022; Andoni et al., 2023; Li & Li, 2023b; Huang & Yi, 2021; Ghazi et al., 2023; Backurs et al., 2024; Cohen-Addad et al., 2022a; Epasto et al., 2024; Chen et al., 2022; Hopkins et al., 2023; Narayanan, 2022; 2023; Jung et al., 2019; Li & Li, 2024; Fan & Li, 2022; Fan et al., 2024; Li & Li, 2023a; Cherapanamjeri et al., 2023; Cohen-Addad et al., 2022b; Dong et al., 2024; Farhadi et al., 2022; Gopi et al., 2021; 2023; Li et al., 2022; Gopi et al., 2022; Eliáš et al., 2020; Song et al., 2023b; Dinur et al., 2023; Woodruff et al., 2023; Song et al., 2023a; Gao et al., 2024; Liang et al., 2024a; Li et al., 2024b). Additionally, new mechanisms have been proposed beyond the traditional Laplace, Gaussian, and Exponential mechanisms (Dwork & Roth, 2014). For example, truncated Laplace mechanism (Geng et al., 2020) is proved to be the current tightest the lower and upper bounds on the minimum noise amplitude and power cross all (ϵ, δ) -DP distributions.

Cross-Attention in System Prompt, RAG, Stable Diffusion and More. Cross-attention (Vaswani et al., 2017), first introduced in language translation, is a widely used technique in many advanced AI systems. For example, Stable Diffusion (Rombach et al., 2022; Liang et al., 2024d; Wang et al., 2023b;c; 2024b) and SORA (OpenAI, 2024b) employ cross-attention as a core module for a text-to-image conditional generation. This technique is also utilized by other multimodal models (Liang et al., 2024e), including Imagen (Saharia et al., 2022) and Diffusion Transformer (Peebles & Xie, 2023). In the realm of text-to-image editing, Hertz et al. (2022) analyzes and controls the cross-attention module to enable editing without requiring additional training. Furthermore, Yang et al. (2024) tackles the issue of inaccurate cross-attention maps, enhancing fine-grained control over edited regions while preventing unintended changes to other areas. In addition, Retrieval Augmented Generation (RAG) (Lewis et al., 2020; Borgeaud et al., 2022; Gao et al., 2023), a technique that improves model responses by retrieving information from a knowledge base or external documents, extensively uses cross-attention as its core design module. Cross-attention also has other applications. Oymak et al. (2023) demonstrates that the prompt-tuning (Liang et al., 2024c) task can be formulated as cross-attention, while Chen et al. (2021) uses cross-attention to fuse multi-scale features in vision transformers, thereby reducing computation. Moreover, attention-based Transformer architecture makes LGMs equipping many emergent ability (Wei et al., 2022), such as spatial reasoning (Wang et al., 2024a), mathematical reasoning (Li et al., 2024a), in-context learning ability (Shi et al., 2024), compositional ability (Xu et al., 2024b), few-shot adaptation ability (Shi et al., 2022b; Xu et al., 2023), and so on. There are some other works that used cross attention in Hopfield Models (Hu et al., 2023; Wu et al., 2024b; Hu et al., 2024c; Xu et al., 2024a; Wu et al., 2024a; Hu et al., 2024a;b).

B MORE PRELIMINARY

In Section B.1, we give the probability tools we use in the paper. In Section B.2, we provide the algebraic facts we use. In Section B.3, we give the DP facts we use in the paper. In Section B.4, we compare between popular DP mechanisms.

B.1 PROBABILITY TOOLS

In this section, we give several probability lemmas.

Lemma B.1 (Markov’s inequality). *If x is a nonnegative random variable and $t > 0$, we have*

$$\Pr[x \geq t] \leq \frac{\mathbb{E}[x]}{t}.$$

Lemma B.2 (Chernoff bound, (Chernoff, 1952)). *Let x_i be a Bernoulli random variable with probability p_i of being equal to 1 and $1 - p_i$ of being equal to 0, and all x_i for $i \in [n]$ are independent. Let $x = \sum_{i=1}^n x_i$. Let $\mu = \mathbb{E}[x] = \sum_{i=1}^n p_i$. Then, for all $\delta > 0$ we have*

$$\Pr[x \geq (1 + \delta)\mu] \leq \exp(-\delta^2\mu/3),$$

and for all $0 < \delta < 1$

$$\Pr[x \leq (1 - \delta)\mu] \leq \exp(-\delta^2\mu/2).$$

Lemma B.3 (Chebyshev’s inequality). *Let x (integrable) be a random variable with finite non-zero variance σ^2 (and thus finite expected value μ). Then for any real number $k > 0$,*

$$\Pr[|x - \mu| \geq k\sigma] \leq \frac{1}{k^2}.$$

B.2 ALGEBRAIC FACTS

Fact B.4 (Upper bound of exponential, Fact C.9 in Liang et al. (2024d)). *For $a \in \mathbb{R}$, $b \in \mathbb{R}$, $a, b \leq R$, where $R \geq 0$, we have*

$$|\exp(a) - \exp(b)| \leq \exp(R)|a - b|.$$

B.3 DP FACTS

In this section, we present several facts about differential privacy (DP). We first state the post-processing property, which means, in an algorithm, if one step is DP, all the following steps are DP.

Fact B.5 (Post-processing, see Fact 2.1 in Ghazi et al. (2023)). *Let \mathcal{A}_1 be an (ϵ, δ) -DP algorithm and \mathcal{A}_2 be a (randomized) post-processing algorithm. Then the algorithm $\mathcal{A}(X) = \mathcal{A}_2(\mathcal{A}_1(X))$ is still an (ϵ, δ) -DP algorithm.*

If we have many DP algorithms, we need a composition rule. The most straightforward composition is the basic/sequential composition rule.

Fact B.6 (Basic composition, see Fact 2.3 in Ghazi et al. (2023)). *Let \mathcal{A}_1 be an (ϵ_1, δ_1) -DP algorithm and \mathcal{A}_2 be an (ϵ_2, δ_2) -DP algorithm. Then $\mathcal{A}(X) = (\mathcal{A}_1(X), \mathcal{A}_2(\mathcal{A}_1(X), X))$ is an $(\epsilon_1 + \epsilon_2, \delta_1 + \delta_2)$ -DP algorithm.*

We can do much better if we know that the inputs are disjoint.

Fact B.7 (Parallel composition, see Fact 2.4 in Ghazi et al. (2023)). *Let \mathcal{A}_1 be an (ϵ_1, δ_1) -DP algorithm and \mathcal{A}_2 be an (ϵ_2, δ_2) -DP algorithm. Assume \mathcal{A}_1 and \mathcal{A}_2 depend on disjoint subsets of input coordinates. Then the algorithm $\mathcal{A}(X) = (\mathcal{A}_1(X), \mathcal{A}_2(\mathcal{A}_1(X), X))$ is a $(\max\{\epsilon_1, \epsilon_2\}, \max\{\delta_1, \delta_2\})$ -DP algorithm.*

In addition, we have the advanced composition, which improves the dependence of the number of DP algorithms to square root but compromises the term δ' .

Theorem B.8 (Advanced composition, see Theorem 3.20 in Dwork & Roth (2014)). *For all $\epsilon, \delta, \delta' \geq 0$, the class of (ϵ, δ) -differentially private mechanisms satisfies $(\epsilon', k\delta + \delta')$ -differential privacy under k -fold adaptive composition for:*

$$\epsilon' = k\epsilon(e^\epsilon - 1) + \epsilon\sqrt{2k \log(1/\delta')}.$$

B.4 COMPARISON OF TRUNCATED LAPLACE, GAUSSIAN, AND LAPLACE MECHANISMS

We first define the Laplace mechanism as below:

Definition B.9 (Laplace distribution). *We use $\text{Lap}(b)$ to denote the pdf: $p(z) = \frac{1}{2b} \exp(-\frac{|z|}{b})$.*

Fact B.10. *For $z \sim \text{Lap}(b)$, $\mathbb{E}[z] = 0$, and $\text{Var}[z] = 2b^2$. Furthermore, if $b = \Delta/\epsilon$, we have $\text{Var}[z] = 2\Delta^2/\epsilon^2$.*

In this paper, we use the Chebyshev inequality to bound the error, and from Geng et al. (2020), we know that the truncated Laplace mechanism has the current minimum variance across all (ϵ, δ) -DP distributions.

The variance of Gaussian mechanism in Theorem 3.22 in Dwork & Roth (2014):

$$\text{Var} = \frac{2\Delta^2 \log(1.25/\delta)}{\epsilon^2}.$$

The variance of Laplace mechanism in Fact B.10:

$$\text{Var} = \frac{2\Delta^2}{\epsilon^2}.$$

The variance of truncated Laplace mechanism in Fact 2.5, for $c \in (0, 1]$:

$$\text{Var} = \frac{2\Delta^2 c}{\epsilon^2}.$$

Thus, since it has the minimum variance, we choose the truncated Laplace mechanism to design our algorithms among these popular mechanisms.

C PROOF OUTLINE

This section provides the proof outline of our paper. In Section C.1, we analyze our DPTREE data structure. In Section C.2, we show the sensitivity of summation problem. In Section C.3, we explain the high-level idea behind the weighted ℓ_p^p distance query. In Section C.4, we show how to answer one-dimensional weighted ℓ_1 distance query. In Section C.5, we show how to answer Softmax distance query using previous algorithms. In Section C.6, we show how to handle adaptive query. By combining the results from these sections, we prove the main results in Section 4.

C.1 SUMMATION SEGMENT TREE

First, in order to solve the weighted distance problem, we need to have a basic DP algorithm (Algorithm 2) that can answer simple summation queries. After analyzing its DP and error in Section D, we state the data structure theorem.

Theorem C.1 (DPTREE data structure, informal version of Theorem D.1). *There is a data structure (see DPTREE in Algorithm 2) that uses $O(n)$ spaces to support the following operations:*

- **INIT**($a \in \mathbb{R}^n, n \in \mathbb{N}_+, \Delta \in \mathbb{N}_+, \epsilon \in (0, 1), \delta \in (0, 1)$). *It takes $O(n)$ time to initialize the data structure.*
- **QUERY**($x \in [n], y \in [n]$). *It takes $O(\log n)$ time to output a number z such that*
 - *the process of output z satisfies (ϵ, δ) -DP private, which computes $\sum_{i=x}^y a_i$,*
 - $|z - \sum_{i=x}^y a_i| \leq O(\epsilon^{-1} \Delta \log^{3/2} n)$,
 - *it holds with probability 0.99.*

During the design of the data structure, we found an interesting property based on the parallel composition rule of DP Fact B.7. We will now state the lemma, whose proof is provided in Section D.

Lemma C.2 (Weighted sum of disjoint interval queries, informal version of Lemma D.8). *If the following conditions hold that:*

- Let there be t disjoint intervals, i.e., S_j for $j \in [t]$, such that $S_j \cap S_k = \emptyset$ for all $j \neq k$.
- Let $\epsilon \in (0, 1)$ and $\delta \in (0, 1)$.
- Let a_j for $j \in [t]$ be a series that square converges to a , i.e., $\sum_{j=1}^t a_j^2 \leq a$.

Then, we have Alg. 2 is (ϵ, δ) -DP and output $\sum_{j=1}^t a_j \text{QUERY}(S_j)$ with the error upper bounded by

$$O(a^{1/2} \epsilon^{-1} \Delta \log^{3/2} n)$$

with probability 0.99.

From Lemma C.2, we can see that if we have multiple disjoint interval queries, the error of the weighted sum of the intervals can be bounded independently of the number of queries, as long as the sum of squared weights is finite.

C.2 SENSITIVITY FOR RANGE SUMMATION PROBLEM

Our DP summation tree data structure DPTREE (Algorithm 2) requires sensitivity parameter Δ . In this section, we show that for the summation problem, we have the sensitivity $\Delta = 2R$ if the input $X \in [-R, R]^n$.

Lemma C.3 (Sensitivity of summation). *Let $X \in [-R, R]^n$. We have the sensitivity $\Delta = 2R$ for DPTREE.INIT in Algorithm 2.*

Proof. Let's say two neighboring datasets X and X' differ in x_i and x'_i for some i in the array X . Then for a summation problem, i.e. $f(X) := \sum_{i=1}^n x_i$, we have

$$\Delta = \max_{X, X'} |f(X) - f(X')| = \max_{X, X'} |x_i - x'_i| = 2R.$$

where the first step follows from Definition 2.2, the second step follows from X, X' differ in x_i, x'_i , and the last step follows from each coordinate of the dataset is bounded in $[-R, R]$. \square

C.3 WEIGHTED ℓ_p^p DISTANCE PROBLEM

In this section, we introduce the intuition behind the method for handling the weighted ℓ_p^p distance problem. The formal lemmas and proofs can be found in Section E.

Given a dataset and a query point in d dimensions, we round each coordinate of the data points and the query point to the nearest multiple of a small interval. We then aggregate the weights of data points that have been rounded to the same position. Finally, we compute the sum of these aggregated weights multiplied by the distances between the query point and the data points over the rounded positions. This approach makes the computation more efficient while maintaining sufficient accuracy.

We provide an example of weighted ℓ_1 -distance of a one-dimensional dataset consisting of 10 data points, i.e., $X \in [0, 1]^{10}$ and a query $y = 0$ in Figure 1.

Lemma C.4 (Weighted ℓ_p^p -distance high dimension, informal version of Lemma E.2). *If the following conditions hold:*

- Let data $X \in [0, R]^{n \times d}$ and $x_i^\top \in [0, R]^d$ be the i -th row of x , weight $w \in \mathbb{R}^n$, query $y \in [0, R]^d$.
- We round each dimension of X and y to an integer multiple of R/n .
- Let $x_{i,k}, y_k$ denote the k -th coordinates of x_i, y for $k \in [d]$.
- Let $c_{j,k} := \sum_{j_0 \in S_{j,k}} w_{j_0}$ where the set $S_{j,k}$ is the set of index i such that the corresponding $x_{i,k}$ is rounded to jR/n for $j \in \{0, 1, 2, \dots, n\}$ for $k \in [d]$.
- After rounding, we assume that y_k is in the $l_k R/n$ position for $l_k \in \{0, 1, 2, \dots, n\}$ for $k \in [d]$.

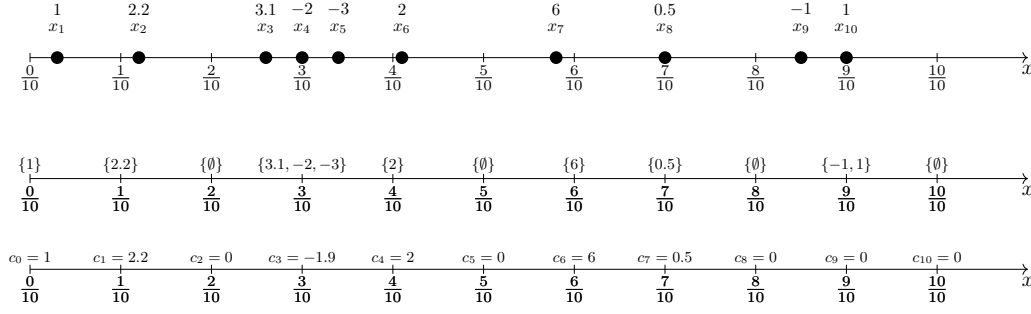


Figure 1: The visualization of how to build c_j of rounded dataset $X \in [0, 1]^{10}$ and compute the weighted ℓ_1 distance. The number above each x_i is w_i . See Algorithm 5 for details. Suppose $y = 0$. Then $\sum_{i=1}^n w_i |y - x_i| = 0.1 \cdot 2.2 + 0.3 \cdot 3.1 + 0.3 \cdot (-2) + 0.3 \cdot (-3) + 0.4 \cdot 2 + 0.6 \cdot 6 + 0.7 \cdot 0.5 + 0.9 \cdot (-1) + 0.9 \cdot 1 = 4.4$. And $\sum_{j=0}^n |k - j| c_j / n = 0.1 \cdot 2.2 + 0.3 \cdot (-1.9) + 0.4 \cdot 2 + 0.6 \cdot 6 + 0.7 \cdot 0.5 = 4.4$. See details in Lemma E.1.

For the weighted problem, we have

$$\sum_{i=1}^n w_i \cdot \|y - x_i\|_p^p = \sum_{k=1}^d \sum_{j=0}^n (|l_k - j| R/n + O(R/n))^p c_{j,k}.$$

where $O(R/n)$ is the rounding error for each data point.

Remark C.5. In Lemma C.4, we first round the dataset. This rounding simplifies the calculation by reducing the number of possible positions to consider, from real values in $[0, R]^d$ to the total $O(nd)$ spots. However, it also introduces an error $O(R/n)$ for one data point. Then, for one spot in the rounded dataset, we sum over the weights of that spot and multiply the corresponding distance raised to the power of p . Additionally, since we are dealing with ℓ_p^p distance, the rounding error is also raised to the power of p .

C.4 ONE-DIMENSIONAL WEIGHTED ℓ_1 DISTANCE DATA STRUCTURE

Based on previous discussions in Section C.1 and C.3, we can now describe our one-dimensional weighted ℓ_1 distance data structure, **DPTREEDISTANCE**, presented in Algorithm 5 and 6, which generalizes the results from Backurs et al. (2024). Drawing from the intuition in Section C.3, the initialization process is as follows: first, we round each data point in the dataset to the nearest multiple of a small interval and build an array that aggregates the corresponding weights. This array is then fed into our **DPTREE** data structure in Algorithm 2. At query time, we query the **DPTREE** to obtain the aggregated weights within a small interval and multiply these weights by the distance to the query point. Furthermore, we also introduce a relative error parameter α to reduce the total number of queries to $O(\log(n)/\alpha)$ instead of querying all n positions. We also analyze the DP and the error bound; see details in Section F.

Theorem C.6 (**DPTREEDISTANCE** data structure, informal version of Theorem F.6). *There is a data structure **DPTREEDISTANCE** (Algorithm 5 and 6) that uses $O(n)$ spaces to solve weighted ℓ_1 -distance query problem for dataset $X \subset [0, R]$ and support the following operations:*

- **INIT**($X \subset [0, R], n \in \mathbb{N}_+, w \in [-R_w, R_w]^n, \epsilon \in (0, 1), \delta \in (0, 1)$). (Algorithm 5) *It takes $O(n)$ time to initialize the data structure.*
- **DISTANCEQUERY**($y \in [0, R], \alpha \in (0, 1)$). (Algorithm 6) *It takes $O(\alpha^{-1} \log^2 n)$ time to output a number z such that*
 - *the process of output z satisfies (ϵ, δ) -DP, which computes $\sum_{i \in [n]} w_i |y - x_i|$,*
 - $|z - \sum_{i \in [n]} w_i |y - x_i|| \leq \alpha \sum_{i \in [n]} w_i |y - x_i| + O(\epsilon^{-1} \alpha^{-1/2} R R_w \log^{3/2} n),$
 - *it holds with probability 0.99.*

C.5 SOFTMAX ACTIVATION

We then describe how we extend the previous results to Softmax activation, i.e. exponential inner product function (Definition 4.1). From Alman & Song (2023), we know that Softmax activation can be approximated by polynomial kernel function $P(\cdot)$ with a certain error. The following lemma shows that we can transform weighted Softmax queries into polynomial kernels. More specifically, we have one term that computes the weighted ℓ_2^2 distance, which is the place where we add DP noises. Because of the decomposability of the ℓ_p^p distance, i.e. $\sum_{i \in [n]} w_i \|x_i - y\|_p^p = \sum_{j \in [d]} \sum_{i \in [n]} w_i |x_{i,j} - y_j|^p$, we can easily extend the results of Section C.4 to handle the ℓ_2^2 distance query. After that, we compute the term for the weighted ℓ_2^2 norms of approximation kernel exactly. Summing all these terms, with a certain error, we can answer the Softmax query. Related proofs can be found in Section I.

Lemma C.7 (Weighted Softmax approximation, informal version of Lemma I.6). *Let the accuracy parameter be $\epsilon_s \in (0, 0.1)$. Let $R \geq 1$. Let $r \leq \binom{2s+2d}{2s}$ and $s = O(\max\{\frac{\log(1/\epsilon_s)}{\log(\log(1/\epsilon_s)/R)}, R^2\})$. Let $\Gamma_{R,s} := \max_{j \in [s]} \frac{R^j}{\sqrt{j!}}$ (Definition I.3). Let $P(x) : [0, R]^d \rightarrow [0, \Gamma_{R,s}]^r$ be the s -th order polynomial kernel function defined in Lemma I.5. Then, we can approximate the exponential inner product using the polynomial kernel function:*

$$w^\top \exp(Xy/d) = -\frac{1}{2} \sum_{j \in [r]} \sum_{i \in [n]} w_i |P(x_i)_j - P(y)_j|^2 + \frac{1}{2} \sum_{i \in [n]} w_i (\|P(x_i)\|_2^2 + \|P(y)\|_2^2) + O(w^\top \exp(Xy/d) \cdot \epsilon_s).$$

Moreover, the vectors $P(\cdot)$ can be computed in $O(r)$ time.

C.6 ADAPTIVE QUERY

We introduce how we can modify our algorithm to solve the adaptive query problem using some tools in Qin et al. (2022). Our approach is based on proving that our algorithm can handle any query within the query space with a certain error. Since adaptive queries must lie within this space, our algorithm can effectively handle them. In Section C.5, we demonstrate our algorithm's capability to answer weighted Softmax distance queries with constant probability. We then use the Chernoff bound to boost the constant probability of our algorithm to a high probability. Next, we apply the notion of an ϵ_0 -net to bound all query points within the net using the union bound. Finally, we bound all points in the query space by utilizing the Lipschitz property of the weighted Softmax distance function and introducing an additive error. See the proofs in Sections H and I.

Lemma C.8 (Adaptive Softmax, informal version of Lemma I.10). *If the following conditions hold:*

- Let N be the ℓ_∞ ϵ_0 -net of \mathcal{B} , and let $|N|$ be the size of the net N .
- Let data set $X \in [0, R]^{n \times d}$, weights $w \in [-R_w, R_w]^n$, query $y \in [0, R]^d$.
- Let the relative error parameter $\alpha \in (0, 1)$, the failure probability $p_f \in (0, 0.01)$.
- We create $l = O(\log((R/\epsilon_0)^r/p_f))$ independent copies of the data structure $\{\text{DPTREESOFTMAX}_j\}_{j=1}^l$ (Algorithm 3) and take the median of the outputs with each data structure instantiated with $(\epsilon/l, (\delta + \delta')/l)$ -DP.
- Let $f(y) := \text{Median}(\{\text{DPTREESOFTMAX}_j.\text{DISTANCEQUERY}(y, \alpha)\}_{j=1}^l)$.
- Let $Z(y) := w^\top \exp(Xy/d)$.
- Let $B = O(\epsilon^{-1} \alpha^{-1/2} l \Gamma_{R,s}^2 R_w r \sqrt{\log(l/\delta')} \cdot \log^{3/2} n)$.

Then with probability $1 - p_f$, for all query points $q \in \mathcal{B}$, there exists a point $y \in N$ which is the closest to q , we can have the process of outputting median of l responses is $(\epsilon, \delta + \delta')$ -DP and the error satisfies

$$|f(y) - Z(q)| \leq (\alpha + \epsilon_s)Z(q) + B + 2n\sqrt{d}R_w \exp(R^2)\epsilon_0.$$

D DPTREE ALGORITHM

In this section, we give the analysis of privacy, accuracy and runtime of our DPTREE (Algorithm 2). In Section D.1, we give the theorem (Theorem D.1) of our data structure that can answer summation problem. In Section D.2, we improve our data structure from constant probability to high probability by applying Chernoff bound. In Section D.3, we give the analysis. In Section D.4, we show some results of our data structure if the input queries are disjoint.

D.1 SINGLE DATA STRUCTURE

We give the theorem of our DPTREE data structure that can answer the summation problem with DP, accuracy, runtime guarantee.

Theorem D.1 (DPTREE data structure, formal version of Theorem C.1). *There is a data structure (see DPTREE in Algorithm 2) that uses $O(n)$ spaces to support the following operations:*

- INIT($a \in \mathbb{R}^n, n \in \mathbb{N}_+, \Delta \in \mathbb{N}_+, \epsilon \in (0, 1), \delta \in (0, 1)$). *It takes $O(n)$ time to initialize the data structure.*
- QUERY($x \in [n], y \in [n]$). *It takes $O(\log n)$ time to output a number z such that*
 - *the process of output z satisfies (ϵ, δ) -DP private, which computes $\sum_{i=x}^y a_i$,*
 - $|z - \sum_{i=x}^y a_i| \leq O(\epsilon^{-1} \Delta \log^{3/2} n)$,
 - *it holds with probability 0.99.*

Proof. The proofs follow from combining Lemma D.3 (running time of initialization), Lemma D.4 (running time of query), Lemma D.5 (DP of query), and Lemma D.6 (error of query) together. \square

D.2 BOOST THE CONSTANT PROBABILITY TO HIGH PROBABILITY

We can use Chernoff bound to boost the high probability by repeating the data structure multiple times.

Theorem D.2 (High-probability). *There is a data structure (see DPTREEHIGHPROB in Algorithm 4) that uses $O(n \log(1/\delta_{\text{fail}}))$ spaces to support the following operations*

- INIT($a \in \mathbb{R}^n, n \in \mathbb{N}_+, \Delta \in \mathbb{N}_+, \epsilon \in (0, 1), \delta \in (0, 1), \delta_{\text{fail}} \in (0, 0.01)$). *It takes $O(n \log(1/\delta_{\text{fail}}))$ time to initialize the data structure.*
- QUERY($x \in [n], y \in [n]$). *It takes $O(\log(n) \cdot \log(1/\delta_{\text{fail}}))$ time to output a number z such that*
 - *the process of output z satisfies (ϵ, δ) -DP private, which computes $\sum_{i=x}^y a_i$,*
 - $|z - \sum_{i=x}^y a_i| \leq O(\epsilon^{-1} \Delta \log^{3/2}(n) \cdot \log(1/\delta_{\text{fail}}))$,
 - *it holds with probability $1 - \delta_{\text{fail}}$ for failure probability $\delta_{\text{fail}} \in (0, 0.01)$.*

Proof. Note that our data structure (Theorem D.1) succeeds with probability 0.99. The success of the algorithm (Theorem D.1) can be viewed as a Bernoulli random variable, to which we apply the Chernoff bound (Lemma B.2). By repeating the data structure $O(\log(1/\delta_{\text{fail}}))$ times and taking the median of the outputs, we boost the success probability. The details are following.

To boost the success probability, we assume the query is repeated l times. Let $i \in [l]$, and let z_i denote the indicator random variable for the success of the i -th instance of the data structure for a single query. Let $z = \sum_{i=1}^l z_i$ be the total success times. Since $p = \Pr[z_i = 1] = 0.99$, we can have $\mu = \mathbb{E}[z] = \sum_{i=1}^l p = lp$. Note that $p = 0.99$. By setting $\delta = 0.1$ and using Chernoff bound from Lemma B.2, we can show

$$\Pr[z \leq l/2] \leq \Pr[z \leq (1 - \delta)lp] \leq \exp(-\delta^2 lp/2).$$

Note that we want $z > l/2$ (since we want at least half to succeed so we could take the median),

$$\Pr[z > l/2] \geq 1 - \exp(-\delta^2 lp/2).$$

To ensure that failure probability is δ_{fail} , we have

$$\exp(-\delta^2 lp/2) = \delta_{\text{fail}}.$$

We can make this hold by choosing $l = O(\log(1/\delta_{\text{fail}}))$.

By the DP basic composition rule (Fact B.6), we need to choose $\epsilon = \epsilon'/O(\log(1/\delta_{\text{fail}}))$ and $\delta = \delta'/O(\log(1/\delta_{\text{fail}}))$ where ϵ', δ' are the ϵ, δ in Theorem D.1. \square

Algorithm 4 Boost constant probability

```

1: datastructure DPTREEHIGHPROB ▷ Theorem D.2
2: members
3:    $\mathcal{D}_1, \dots, \mathcal{D}_{O(\log(1/\delta_{\text{fail}}))} : \text{DPTREE}$  ▷ Alg. 2
4: end members
5: procedure INIT( $a \in \mathbb{R}^n, n \in \mathbb{N}_+, \Delta \in \mathbb{N}_+, \epsilon \in (0, 1), \delta \in (0, 1), \delta_{\text{fail}} \in (0, 0.01)$ )
6:   for  $i = 1 \rightarrow O(\log(1/\delta_{\text{fail}}))$  do
7:      $\mathcal{D}_i \leftarrow \text{INIT}(a, n, \Delta, \epsilon/O(\log(1/\delta_{\text{fail}})), \delta/O(\log(1/\delta_{\text{fail}})))$ 
8:   end for
9: end procedure
10: procedure QUERY( $x \in [n], y \in [n]$ )
11:    $r \leftarrow 0^{O(\log(1/\delta_{\text{fail}}))}$ 
12:   for  $i = 1 \rightarrow O(\log(1/\delta_{\text{fail}}))$  do
13:      $r_i \leftarrow \mathcal{D}_i \cdot \text{QUERY}(x, y)$ 
14:   end for
15:   return Median of  $r$ 
16: end procedure
17: end datastructure

```

D.3 ALGORITHM OF DATA STRUCTURE

In this section, we analyze the accuracy, DP, and runtime of Algorithm 2.

We first analyze the runtime.

Lemma D.3 (Runtime of initialization, Algorithm 2). *For the initialization, we have the time complexity of Algorithm 2 is $O(n)$.*

Proof. All the computations are dominated by $O(n)$ time. \square

Lemma D.4 (Runtime of query, Algorithm 2). *For each query, we have the time complexity of Algorithm 2 is $O(\log n)$.*

Proof. Due to the property of tree, we will use at most $2 \log n$ nodes in the tree, thus the running time is $O(\log n)$. \square

We now analyze the DP.

Lemma D.5 (Privacy of query, Algorithm 2). *The output process of QUERY (see Algorithm 2) is (ϵ, δ) -DP.*

Proof. Suppose that our dataset is $X \in [-R, R]^n$. Note that we only add noise in the pre-processing stage. There is no noise in the query stage. Since the problem we care about is summation, if we change one leaf node, the sensitivity $\Delta = 2R$ (see Lemma C.3). Since we add noise to each node in the tree, and each leaf node count will contribute to $\log n$ nodes, it is equivalent to our output function being in $\log n$ dimension. We will then blow up the DP parameter by $\log n$ factor. Thus, using the basic composition rule (Fact B.6), the DP guarantee for the whole tree data structure is $((\epsilon/\log n) \cdot \log n, (\delta/\log n) \cdot \log n)$ which is (ϵ, δ) -DP. \square

We now analyze the accuracy.

Lemma D.6 (Accuracy of query, Algorithm 2). *Let $\epsilon \in (0, 1)$ and $\delta \in (0, 1)$. Then, using Chebyshev's inequality and Fact 2.5, we have the error of `QUERY`(see Algorithm 2) output is upper bounded by:*

$$O(\epsilon^{-1} \Delta \log^{3/2} n).$$

with probability 0.99.

Proof. For an interval S , we define `TRUEQUERY`(S) to be the output of `DPTREE.TRUEQUERY` in Algorithm 2. Let `QUERY`(S) denote the noised interval query answer returned by `DPTREE.QUERY` in Algorithm 2. Let $z := \text{QUERY}(S) - \text{TRUEQUERY}(S)$, which from Algorithm 2 we can see this is the sum of $O(\log n)$ independent truncated Laplace random variables each with parameter $\text{TLap}(\Delta, \epsilon/\log n, \delta/\log n)$. Thus,

$$z = \sum_{i=1}^{O(\log n)} z_i$$

where $z_i \sim \text{TLap}(\Delta, \epsilon/\log n, \delta/\log n)$, and every z_i are independent to each other.

We know $\mu = \mathbb{E}[z] = 0$ since $\mathbb{E}[z_i] = 0$. From Fact 2.5, we know the variance for each z_i is $\text{Var}[z_i] = c\epsilon^{-2}\Delta^2 \log^2 n$ where $0 < c \leq 2$ and $c = 2$ when $\delta = 0$.

Therefore, we can show

$$\begin{aligned} \text{Var}[z] &= \text{Var}\left[\sum_{i=1}^{O(\log n)} z_i\right] \\ &= \sum_{i=1}^{O(\log n)} \text{Var}[z_i] \\ &= O(c\epsilon^{-2}\Delta^2 \log^3 n) \end{aligned} \tag{4}$$

where the first step follows from definition of z , the second step follows from every z_i are independent to each other, and the last step follows from $\text{Var}[z_i] = O(c\epsilon^{-2}\Delta^2 \log^2 n)$.

Note that we wish to bound $|z| = |\text{QUERY}(S) - \text{TRUEQUERY}(S)|$ as our error.

Using Lemma B.3, we can have

$$\Pr[|z| \geq k\sigma] \leq \frac{1}{k^2}.$$

We know that $\sigma = \sqrt{\text{Var}[z]} = O(c^{1/2}\epsilon^{-1}\Delta \log^{3/2} n)$. Picking $k = 10$, we have

$$\Pr[|z| < 10\sigma] \geq 0.99.$$

Thus, we conclude that error is bounded by $O(c^{1/2}\epsilon^{-1}\Delta \log^{3/2} n) = O(\epsilon^{-1}\Delta \log^{3/2} n)$ (since $c \in (0, 2]$) with probability 0.99. \square

D.4 DISJOINT INTERVALS

In this section, we show some interesting results for our `DPTREE` data structure if the input queries are disjoint.

Lemma D.7 (Sum of disjoint interval queries). *If the following conditions hold that:*

- *Let there be t disjoint intervals, i.e., S_j for $j \in [t]$, such that $S_j \cap S_k = \emptyset$ for all $j \neq k$.*
- *Let $\epsilon \in (0, 1)$ and $\delta \in (0, 1)$.*

Then, we have Algorithm 2 is (ϵ, δ) -DP and outputs $\sum_{j=1}^t \text{QUERY}(S_j)$ with the error upper bounded by

$$O(t^{1/2}\epsilon^{-1}\Delta \log^{3/2} n)$$

with probability 0.99.

Proof. From Lemma D.5, we know that DPTree.QUERY is (ϵ, δ) -DP. Then, from Fact B.7 and the disjoint intervals in Algorithm 6, we can conclude that the value returned is (ϵ, δ) -DP.

Let $\text{TRUEQUERY}(S_j)$ denote the true interval query answer returned by DPTREE.TRUEQUERY in Algorithm 2 for interval S_j . Let $\text{QUERY}(S_j)$ denote the noised interval query answer returned by DPTREE.QUERY in Algorithm 2 for interval S_j . Let $z_j := \text{QUERY}(S_j) - \text{TRUEQUERY}(S_j)$ and $z = \sum_{j=1}^t z_j$. From the proof of Lemma D.6, we know z_j is the sum of $O(\log n)$ independent truncated Laplace random variables each with parameter $\text{TLap}(\Delta, \epsilon/\log n, \delta/\log n)$ and the variance is bounded by

$$\text{Var}[z_j] = O(\epsilon^{-2} \Delta^2 \log^3 n)$$

Since the intervals S_j are disjoint, they are independent to each other. Then, we have

$$\begin{aligned} \text{Var}[z] &= \text{Var}\left[\sum_{j=1}^t z_j\right] \\ &= \sum_{j=1}^t \text{Var}[z_j] \\ &= O(t\epsilon^{-2} \Delta^2 \log^3 n) \end{aligned}$$

where the first step follows from definition of z , the second step follows from the intervals are disjoint, and the last step follows from $\text{Var}[z_j] = O(\epsilon^{-2} \Delta^2 \log^3 n)$.

Note that we wish to bound $|z|$ as our error.

Using Lemma B.3, we can have error bounded by

$$O(t^{1/2} \epsilon^{-1} \Delta \log^{3/2} n)$$

with probability 0.99. □

Moreover, this can be generalized to weighted sum of queries.

Lemma D.8 (Weighted sum of disjoint interval queries, formal version of Lemma C.2). *If the following conditions hold that:*

- Let there be t disjoint intervals, i.e., S_j for $j \in [t]$, such that $S_j \cap S_k = \emptyset$ for all $j \neq k$.
- Let $\epsilon \in (0, 1)$ and $\delta \in (0, 1)$.
- Let a_j for $j \in [t]$ be a series that square converges to a , i.e., $\sum_{j=1}^t a_j^2 \leq a$.

Then, we have Alg. 2 is (ϵ, δ) -DP and output $\sum_{j=1}^t a_j \text{QUERY}(S_j)$ with the error upper bounded by

$$O(a^{1/2} \epsilon^{-1} \Delta \log^{3/2} n)$$

with probability 0.99.

Proof. The DP proof is the same as in the proof of Lemma D.7.

Let $\text{TRUEQUERY}(S_j)$ and $\text{QUERY}(S_j)$ be same in the proof of Lemma D.7. Let $z_j := \text{QUERY}(S_j) - \text{TRUEQUERY}(S_j)$ and $z = \sum_{j=1}^t a_j z_j$. From the proof of Lemma D.7, we know the variance of z_j is bounded by

$$\text{Var}[z_j] = O(\epsilon^{-2} \Delta^2 \log^3 n)$$

Since the intervals S_j are disjoint, they are independent to each other. Then, we have

$$\text{Var}[z] = \text{Var}\left[\sum_{j=1}^t a_j z_j\right]$$

$$\begin{aligned}
&= \sum_{j=1}^t \text{Var}[a_j z_j] \\
&= \sum_{j=1}^t a_j^2 \text{Var}[z_j] \\
&= \sum_{j=1}^t a_j^2 \cdot O(\epsilon^{-2} \Delta^2 \log^3 n) \\
&= O(a \epsilon^{-2} \Delta^2 \log^3 n)
\end{aligned}$$

where the first step follows from the definition of z , the second step follows from the intervals are disjoint, the third step follows from the $\text{Var}[az] = a^2 \text{Var}[z]$ for a random variable z and a constant a , the fourth step follows from the $\text{Var}[z_j] = O(\epsilon^{-2} \Delta^2 \log^2 n)$, and the last step follows from $\sum_{j=1}^t a_j^2 \leq a$.

Note that we wish to bound $|z|$ as our error.

Using Lemma B.3, we can have error bounded by

$$O(a^{1/2} \epsilon^{-1} \Delta \log^{3/2} n)$$

with probability 0.99. \square

E WEIGHTED ℓ_p^p DISTANCE

In this section, we introduce how to handle weighted ℓ_p^p distance problem in the high level idea. In Section E.1, we show how to solve one dimensional weighted problem. In Section E.2, we show how to solve high dimensional weighted problem by decomposing each coordinate of the high dimensional dataset.

Suppose we have the original data $X \in [0, R]^n$ and weight $w \in \mathbb{R}^n$ and query $y \in [0, R]$. We want to compute the weighted ℓ_1 -distance, i.e.

$$\sum_{i=1}^n w_i \cdot |y - x_i|.$$

For data in d -dimension, due to the decomposability of ℓ_p^p distance, our problem will be: given $x_i \in [0, R]^d$ and $w_i \in \mathbb{R}$ for $i \in [n]$, and $y \in [0, R]^d$, we can compute

$$\sum_{i=1}^n w_i \cdot \|y - x_i\|_p^p = \sum_{j=1}^d \sum_{i=1}^n w_i \cdot |y_j - x_{i,j}|^p$$

where $x_{i,j}, y_j$ means the j -th coordinates of x_i, y for $j \in [d]$.

Therefore, we can solve one dimension problem first, and then the high dimension case can be solved automatically.

E.1 ONE DIMENSIONAL WEIGHTED DISTANCE

Now we can give the lemma for weighted distance of dataset.

Lemma E.1 (Weighted distance one dimension). *If the following conditions hold:*

- Let data $X \in [0, R]^n$, weight $w \in \mathbb{R}^n$, query $y \in [0, R]$.
- We round X and y to an integer multiple of R/n .
- Let $c_j = \sum_{j_0 \in S_j} w_{j_0}$ where set S_j is the set of index i such that the corresponding x_i is rounded to jR/n for $j \in \{0, 1, 2, \dots, n\}$.

- After rounding, we assume y is in the kR/n position for $k \in \{0, 1, 2, \dots, n\}$.

For the weighted problem, we have

$$\sum_{i=1}^n w_i \cdot |y - x_i| = \sum_{j=0}^n (|k - j|R/n + O(R/n))c_j.$$

Moreover, we have

$$\sum_{i=1}^n w_i \cdot |y - x_i|^p = \sum_{j=0}^n (|k - j|R/n + O(R/n))^p c_j$$

where $O(R/n)$ is the rounding error for each data point.

Proof. For each i , we have:

$$w_i \cdot |y - x_i| = w_i \cdot \left(\frac{|k - j|R}{n} + O\left(\frac{R}{n}\right) \right).$$

where $O(R/n)$ is the rounding error introduced by each data point, since each data point will be at most $O(R/n)$ away from its true position.

We can construct c_j by

$$c_j = \sum_{j_0 \in S_j} w_{j_0}$$

set S_j is the set of index i such that the corresponding x_i is rounded to jR/n . Moreover, c_j can be negative.

Summing over all i and grouping by j , we get:

$$\sum_{i=1}^n w_i \cdot |y - x_i| = \sum_{j=0}^n \left(\frac{|k - j|R}{n} + O\left(\frac{R}{n}\right) \right) c_j.$$

The total rounding error will be $O(R)$ because we have n data points, each with an error of at most $O(R/n)$.

Moreover, we have

$$\sum_{i=1}^n w_i \cdot |y - x_i|^p = \sum_{j=0}^n \left(\frac{|k - j|R}{n} + O\left(\frac{R}{n}\right) \right)^p c_j.$$

□

E.2 HIGH DIMENSIONAL WEIGHTED DISTANCE

Finally, we can solve the problem of weighted distance for d -dimensional dataset.

Lemma E.2 (Weighted ℓ_p^p -distance high dimension, formal version of Lemma C.4). *If the following conditions hold:*

- Let data $X \in [0, R]^{n \times d}$ and $x_i^\top \in [0, R]^d$ be the i -th row of x , weight $w \in \mathbb{R}^n$, query $y \in [0, R]^d$.
- We round each dimension of X and y to an integer multiple of R/n .
- Let $x_{i,k}, y_k$ denote the k -th coordinates of x_i, y for $k \in [d]$.
- Let $c_{j,k} := \sum_{j_0 \in S_{j,k}} w_{j_0}$ where set $S_{j,k}$ is the set of index i such that the corresponding $x_{i,k}$ is rounded to jR/n for $j \in \{0, 1, 2, \dots, n\}$ for $k \in [d]$.
- After rounding, we assume y_k is in the $l_k R/n$ position for $l_k \in \{0, 1, 2, \dots, n\}$ for $k \in [d]$.

For the weighted problem, we have

$$\sum_{i=1}^n w_i \cdot \|y - x_i\|_p^p = \sum_{k=1}^d \sum_{j=0}^n (|l_k - j|R/n + O(R/n))^p c_{j,k}$$

where $O(R/n)$ is the rounding error for each data point.

Proof. We can show

$$\begin{aligned} \sum_{i=1}^n w_i \cdot \|y - x_i\|_p^p &= \sum_{k=1}^d \sum_{i=1}^n w_i \cdot |y_k - x_{i,k}|^p \\ &= \sum_{k=1}^d \sum_{j=0}^n (|l_k - j|R/n + O(R/n))^p c_{j,k} \end{aligned}$$

where the first step follows from decomposability of ℓ_p^p -distance by dimension, the second step follows from Lemma E.1.

□

F ONE-DIMENSIONAL WEIGHTED ℓ_1 DISTANCE QUERY

In this section, we generalize the algorithms in Backurs et al. (2024) to weighted distance. Here, we compute the problem of one-dimensional weighted ℓ_1 distance query i.e. $\sum_{i \in [n]} w_i |y - x_i|$ for a given query $y \in [0, R]$, weights $w \in [-R_w, R_w]^n$ and dataset $X \subset [0, R]$ and $n = |X|$. In Section F.1, we analyze the runtime of our algorithm. In Section F.2, we analyze the DP and accuracy of our algorithm. In Section F.3, we give the theorem for our DPTREEDISTANCE data structure.

Algorithm 5 Pre-processing data structure

```

1: datastructure DPTREEDISTANCE ▷ Theorem F.6
2: members
3:    $\mathcal{D} : \text{DPTREE}$  ▷ Alg. 2
4:    $X : [0, R]^n$ 
5:    $w : [-R_w, R_w]^n$ 
6: end members
7: procedure INIT( $X \subset [0, R], n \in \mathbb{N}_+, w \in [-R_w, R_w]^n, \epsilon \in (0, 1), \delta \in (0, 1)$ ) ▷ Lemma E.1
8:    $X, w, a \leftarrow X, w, 0^{n+1}$ 
9:   for  $i = 1 \rightarrow n$  do
10:     $j \leftarrow \text{ROUND}(x_i, n)$  ▷  $x_i \in X$  for  $i \in [n]$ 
11:     $a_j \leftarrow a_j + w_i$ 
12:  end for
13:   $\mathcal{D}.\text{INIT}(a, n + 1, 2R_w, \epsilon, \delta)$  ▷ Alg. 2, Lemma C.3
14: end procedure
15: procedure ROUND( $x \in [0, R], n \in \mathbb{N}_+$ )
16:   Let  $j \in \{0, 1, 2, \dots, n - 1\}$  denote the integer such that  $jR/n \leq x < (j + 1)R/n$ 
17:   if  $|x - (j + 1)R/n| \leq |x - jR/n|$  then
18:      $j \leftarrow j + 1$ 
19:   end if
20:   return  $j$ 
21: end procedure
22: end datastructure

```

F.1 RUNTIME ANALYSIS

We first analyze the runtime.

Algorithm 6 One dimensional weighted ℓ_1 distance query

```

1: datastructure DPTREEDISTANCE ▷ Theorem F.6
2: procedure DISTANCEQUERY( $y \in [0, R], \alpha \in (0, 1)$ ) ▷ Lemma F.2, Lemma F.4, Lemma F.5
3:    $y \leftarrow \text{ROUND}(y, n) \cdot (R/n)$  ▷ Alg. 5
4:   Value  $\leftarrow 0$ 
5:   for  $j = 0, 1, \dots, O(\log(n)/\alpha)$  do
6:      $l_j \leftarrow \text{ROUND}(y + \frac{R}{(1+\alpha)^{j+1}}, n)$ 
7:      $r_j \leftarrow \text{ROUND}(y + \frac{R}{(1+\alpha)^j}, n)$  ▷ Consider the points to the right of  $y$ 
8:     Value  $\leftarrow \text{Value} + \mathcal{D}.\text{QUERY}(l_j, r_j) \cdot \frac{R}{(1+\alpha)^j}$  ▷ Alg. 2
9:   end for
10:  for  $j = 0, 1, \dots, O(\log(n)/\alpha)$  do
11:     $l_j \leftarrow \text{ROUND}(y - \frac{R}{(1+\alpha)^j}, n)$ 
12:     $r_j \leftarrow \text{ROUND}(y - \frac{R}{(1+\alpha)^{j+1}}, n)$  ▷ Consider the points to the left of  $y$ 
13:    Value  $\leftarrow \text{Value} + \mathcal{D}.\text{QUERY}(l_j, r_j) \cdot \frac{R}{(1+\alpha)^j}$  ▷ Alg. 2
14:  end for
15:  Return Value
16: end procedure
17: end datastructure

```

Lemma F.1 (Runtime of initialization, Algorithm 5). *For the initialization, we have the time complexity of INIT in Algorithm 5 is $O(n)$.*

Proof. In the initialization of INIT, the computations need $O(n)$ time to compute the count and $O(\log n)$ time to build the tree. Thus, total time is $O(n)$. \square

Lemma F.2 (Runtime of DISTANCEQUERY, Algorithm 6). *For the ℓ_1 distance query, we have the time complexity of DISTANCEQUERY in Algorithm 6 is $O(\alpha^{-1} \log^2 n)$.*

Proof. In DISTANCEQUERY, the computations need $O(\log n)$ time to compute one value from DPTREE.QUERY and this process need to be repeated $O(\alpha^{-1} \log n)$ times. \square

Remark F.3. In Line 8 and 13 of Algorithm 6, we use $R/(1+\alpha)^j$ to approximate the distance of each data point to the query in Lemma E.1, i.e. $|k-j|R/n$. This will introduce α relative error but also reduce the numbers of iteration from $O(n)$ to $O(\log(n)/\alpha)$.

F.2 PRIVACY AND ACCURACY ANALYSIS

We show the DP.

Lemma F.4 (Privacy of DISTANCEQUERY, Algorithm 6). *The output process of DISTANCEQUERY (Algorithm 6) is (ϵ, δ) -DP.*

Proof. From Lemma D.5, we know that DPTree.QUERY is (ϵ, δ) -DP output. We observe that intervals in Algorithm 6 are disjoint. Then, following the same logic in the proof of Lemma D.8, we can conclude that the value returned is (ϵ, δ) -DP. \square

We now analyze the accuracy of the algorithm.

Lemma F.5 (Accuracy of DISTANCEQUERY, Algorithm 6). *If the following conditions are satisfied:*

- Let $X \in [0, R]^n$ be a dataset consisting of n one-dimensional numbers, with weights $w \in [-R_w, R_w]^n$.
- Let $\alpha \in (0, 1)$ represent the relative error parameter utilized in Algorithm 6.
- Let \tilde{A} denote the output of the DISTANCEQUERY in Algorithm 6.

• Let $A_* := \sum_{i \in [n]} w_i |y - x_i|$ represent the true distance query value for a specific query y .

Then with probability 0.99, we have

$$|\tilde{A} - A_*| \leq \alpha A_* + O(\epsilon^{-1} \alpha^{-1/2} R R_w \log^{3/2} n).$$

Proof. To simplify the explanation, we consider only the distance query for the points in X located to the right of y . The proof can be symmetrically applied to the case of points to the left of y . For an interval $S_j := (l_j, r_j)$ where l_j, r_j are defined in Algorithm 6, we define $\text{TRUEQUERY}(S_j)$ to be the output of DPTREE.TRUEQUERY in Algorithm 2. Let

$$\hat{A} := \sum_{j=0}^{O(\log(n)/\alpha)} \frac{R}{(1+\alpha)^j} \cdot \text{TRUEQUERY}(S_j).$$

Since TRUEQUERY returns the sum of the corresponding weights, it aligns with the true answer $A_* := \sum_{i \in [n]} w_i |y - x_i|$. Thus, we have

$$|\hat{A} - A_*| \leq \alpha \cdot A_*,$$

because for all j , the distances between y and different points in X vary only by a multiplicative factor of $(1 + \alpha)$.

Next we show the additive error. Let $\text{QUERY}(S_j)$ denote the noised interval query answer returned by DPTREE.QUERY in Algorithm 2. Algorithm 6 outputs $\tilde{A} = \sum_{j=0}^{O(\log(n)/\alpha)} \frac{R}{(1+\alpha)^j} \cdot \text{QUERY}(S_j)$. We wish to bound

$$|\hat{A} - \tilde{A}| \leq \left| \sum_{j=0}^{O(\log(n)/\alpha)} \frac{R}{(1+\alpha)^j} \cdot (\text{TRUEQUERY}(S_j) - \text{QUERY}(S_j)) \right|.$$

Let $z_j := \text{QUERY}(S_j) - \text{TRUEQUERY}(S_j)$, which from Algorithm 2 we can see this is the sum of $O(\log n)$ independent truncated Laplace random variables.

From Lemma D.8, we only need to show that the series $\frac{1}{(1+\alpha)^j}$ for $j \in \{0, 1, \dots, O(\log(n)/\alpha)\}$ square converges to $1/\alpha$, since R is a constant.

We can show

$$\begin{aligned} \sum_{j=0}^{O(\log(n)/\alpha)} \frac{1}{(1+\alpha)^{2j}} &\leq \sum_{j=0}^{\infty} \frac{1}{(1+\alpha)^{2j}} \\ &\leq \sum_{j=0}^{\infty} \frac{1}{(1+\alpha)^j} \\ &= \frac{1}{1 - \frac{1}{1+\alpha}} \\ &= 1 + \frac{1}{\alpha} \\ &= O(1/\alpha) \end{aligned}$$

where the first step follows from we extend the finite sum to infinite sum, the second step follows from $\frac{1}{(1+\alpha)^{2j}} \leq \frac{1}{(1+\alpha)^j}$, the third step follows from the closed form of geometric sum, the fourth step follows from simple algebra, and the last step follows from $\alpha \in (0, 1)$.

Then from the proof of Lemma D.8, we can know that the variance is given by

$$O\left(\frac{R^2 R_w^2 \log^3 n}{\alpha \epsilon^2}\right) \quad (5)$$

since the sensitivity $\Delta = 2R_w$ from Lemma C.3.

Using Lemma B.3, we can have additive error bounded by

$$O\left(\frac{R \cdot R_w \log^{3/2} n}{\epsilon \sqrt{\alpha}}\right).$$

with probability 0.99. \square

F.3 ONE DIMENSION SINGLE DATA STRUCTURE

We therefore have the data structure that can solve weighted ℓ_1 -distance problem.

Theorem F.6 (DPTREEDISTANCE data structure, formal version of Theorem C.6). *There is a data structure DPTREEDISTANCE (Algorithm 5,6) that uses $O(n)$ spaces to solve weighted ℓ_1 -distance query problem for dataset $X \subset [0, R]$ and support the following operations:*

- INIT($X \subset [0, R], n \in \mathbb{N}_+, w \in [-R_w, R_w]^n, \epsilon \in (0, 1), \delta \in (0, 1)$). (Algorithm 5) It takes $O(n)$ time to initialize the data structure.
- DISTANCEQUERY($y \in [0, R], \alpha \in (0, 1)$). (Algorithm 6) It takes $O(\alpha^{-1} \log^2 n)$ time to output a number z such that
 - the process of output z satisfies (ϵ, δ) -DP private, which computes $\sum_{i \in [n]} w_i |y - x_i|$,
 - $|z - \sum_{i \in [n]} w_i |y - x_i|| \leq \alpha \sum_{i \in [n]} w_i |y - x_i| + O(\epsilon^{-1} \alpha^{-1/2} R R_w \log^{3/2} n)$,
 - it holds with probability 0.99.

Proof. The proofs follow from combining Lemma F.1 (running time of initialization), Lemma F.2 (running time of query), Lemma F.4 (DP of query), and Lemma F.5 (error of query) together. \square

G HIGH-DIMENSIONAL WEIGHTED ℓ_1 QUERY

In this section, we show how we can solve the high dimensional weighted ℓ_1 distance problem, generalizing results from Backurs et al. (2024). In Section G.1, we give the analysis of Algorithm 7. In Section G.2, we give the theorem of our DPTREEHIGHDIM data structure.

Algorithm 5,6 can be naturally extended to higher dimensions because of the decomposability of the ℓ_1 distance function. We construct d separate one-dimensional distance query data structures, each corresponding to a coordinate projection of the dataset.

G.1 PRIVACY AND ACCURACY ANALYSIS FOR HIGH DIMENSIONAL WEIGHTED DISTANCE

We now give the analysis of our Algorithm 7 for high dimensional weighted ℓ_1 -distance query.

Algorithm 7 High-dimensional weighted ℓ_1 distance query

```

1: datastrucutre DPTREEHIGHDIM ▷ Theorem G.3
2: members
3:    $\mathcal{D}_1, \dots, \mathcal{D}_d : \text{DPTREEDISTANCE}$  ▷ Alg. 5
4:    $X : [0, R]^{n \times d}$ 
5:    $w : [-R_w, R_w]^n$ 
6: end members
7: procedure INIT( $X \subset [0, R]^d, n \in \mathbb{N}_+, w \in [-R_w, R_w]^n, \epsilon \in (0, 1), \delta \in (0, 1), \delta' \in (0, 1), c \in (0, 0.1)$ ) ▷ Lemma E.2
8:    $X \leftarrow X$ 
9:    $w \leftarrow w$ 
10:  for  $i = 1 \rightarrow d$  do
11:     $\mathcal{D}_i.$ INIT( $X[:, i], n, w, c\epsilon/\sqrt{d \log(1/\delta')}$ ,  $\delta/d$ ) ▷ Alg. 5
12:  end for
13: end procedure
14: procedure DISTANCEQUERY( $y \in [0, R]^d, \alpha \in (0, 1)$ ) ▷ Lemma G.1, Lemma G.2
15:  Value  $\leftarrow 0$ 
16:  for  $i = 1 \rightarrow d$  do
17:    Value  $\leftarrow$  Value +  $\mathcal{D}_i.$ DISTANCEQUERY( $y_i, \alpha$ ) ▷ Alg. 6
18:  end for
19:  return Value
20: end procedure
21: end datastrucutre

```

Lemma G.1 (Privacy of DISTANCEQUERY, Algorithm 7). *If the following conditions hold*

- Let data set $X \in [0, R]^{n \times d}$, weights $w \in [-R_w, R_w]^n$, query $y \in [0, R]^d$.
- Let $\epsilon \in (0, 1)$, $\delta \in (0, 1)$, $\delta' \in (0, 1)$.
- Let $c \in (0, 0.1)$ be a small constant and A be the output of DISTANCEQUERY in Algorithm 7, where each one-dimensional algorithm is configured to be $(c\epsilon/\sqrt{d \log(1/\delta')}, \delta/d)$ -DP (see Line 11).
- Let $A_* = \sum_{i \in [n]} w_i \|y - x_i\|_1$ represent the true distance query value.
- Let $\epsilon = O(\log(1/\delta'))$.

Then, we have the output process of DISTANCEQUERY (Algorithm 7) is $(\epsilon, \delta + \delta')$ -DP.

Proof. The $(\epsilon, \delta + \delta')$ -DP guarantee follows from the approximate DP advanced composition result Theorem B.8. Our algorithm instantiate each one-dimensional data structure with $(c\epsilon/\sqrt{d \log(1/\delta')}, \delta/d)$ -DP total d times.

From advanced composition in Theorem B.8, for a sufficient small parameter ϵ and constant c , we have the final privacy loss parameter be:

$$O(c\epsilon\sqrt{2d \log(1/\delta')}/\sqrt{d \log(1/\delta')}) = O(\epsilon)$$

and the final failure probability parameter be:

$$d\delta/d + \delta' = \delta + \delta'.$$

□

Lemma G.2 (Accuracy of DISTANCEQUERY, Algorithm 7). *If the following conditions hold*

- Let data set $X \in [0, R]^{n \times d}$, weights $w \in [-R_w, R_w]^n$, query $y \in [0, R]^d$.
- Let $\epsilon \in (0, 1)$, $\delta \in (0, 1)$, $\delta' \in (0, 1)$.
- Let $c \in (0, 0.1)$ be a small constant and A be the output of DISTANCEQUERY in Algorithm 7, where each one-dimensional algorithm is configured to be $(c\epsilon/\sqrt{d \log(1/\delta')}, \delta/d)$ -DP (see Line 11).
- Let $A_* = \sum_{i \in [n]} w_i \|y - x_i\|_1$ represent the true distance query value.

With probability 0.99, we have

$$|A - A_*| \leq \alpha A_* + O(\epsilon^{-1} \alpha^{-1/2} R R_w d \sqrt{\log(1/\delta')} \cdot \log^3 n).$$

Proof. Let A_i be the i -th dimension output returned by \mathcal{D}_i in Algorithm 7. Let $A_{*,i}$ be the true distance query value in the i -th dimension. Observe that $A_* = \sum_{i=1}^d A_{*,i}$ and $A = \sum_{i=1}^d A_i$.

We follow the similar idea in the proof of Lemma F.5. Let $z_{j,i}$ be the random variables that represent z_j (used in the proof of Lemma F.5) for the i -th coordinate. We can observe that the overall error across d coordinates can be upper bounded by

$$\left| \sum_{i=1}^d \sum_{j=0}^{O(\log(n)/\alpha)} \frac{R z_{j,i}}{(1+\alpha)^j} \right|$$

where each $z_{j,i}$ is the sum of $O(\log n)$ truncated Laplace random variables independent to others. With ϵ scaled down by $c\epsilon/\sqrt{d \log(1/\delta')}$ and δ scaled down by δ/d , the variance of each individual dimension is given by (see Eq. (5))

$$O(\alpha^{-1} \epsilon^{-2} d R^2 R_w^2 \log(1/\delta') \log^3 n).$$

Thus, the total variance for d instantiated data structures is then

$$O(\alpha^{-1}\epsilon^{-2}d^2R^2R_w^2\log(1/\delta')\log^3n).$$

Finally, from Lemma B.3, we have the additive error given by

$$O(\alpha^{-1/2}\epsilon^{-1}dRR_w\sqrt{\log(1/\delta')}\cdot\log^{3/2}n).$$

□

G.2 HIGH DIMENSION SINGLE DATA STRUCTURE

We have the data structure that can solve weighted ℓ_1 -distance problem in d -dimensional data.

Theorem G.3 (DPTREEHIGHDIM data structure). *There is a data structure DPTREEHIGHDIM (Algorithm 7) that uses $O(nd)$ spaces to solve weighted ℓ_1 -distance query problem for dataset $X \subset [0, R]^d$ and support the following operations:*

- INIT($X \subset [0, R]^d, n \in \mathbb{N}_+, w \in [-R_w, R_w]^n, \epsilon \in (0, 1), \delta \in (0, 1), \delta' \in (0, 1), c \in (0, 0.1)$). (Algorithm 7) It takes $O(nd)$ time to initialize the data structure.
- DISTANCEQUERY($y \in [0, R]^d, \alpha \in (0, 1)$). (Algorithm 7) It takes $O(\alpha^{-1}d\log^2n)$ time to output a number z such that
 - the process of output z satisfies is $(\epsilon, \delta + \delta')$ -DP private, which computes $\sum_{i \in [n]} w_i \|y - x_i\|_1$,
 - $|z - \sum_{i \in [n]} w_i \|y - x_i\|_1| \leq \alpha \sum_{i \in [n]} w_i \|y - x_i\|_1 + O(\epsilon^{-1}\alpha^{-1/2}RR_wd\sqrt{\log(1/\delta')}\cdot\log^{3/2}n)$,
 - it holds with probability 0.99.

Proof. For the runtime analysis, since we loop data structure DPTREEDISTANCE d times, an additional d factor will appear for both initialization and query time complexity. The DP is proved by Lemma G.1. The accuracy is proved by Lemma G.2. □

H ADAPTIVE QUERY

In this section, we introduce how we can solve the adaptive query problem by our algorithm, using some tools from Qin et al. (2022). Our idea is that, if we can prove that our algorithm can solve any query in the query space with certain error. Then, since adaptive query must lie in this space, we can handle adaptive query. In Section H.1, we show how we can boost the constant probability of our algorithm to high probability. In Section H.2, we show how we can apply the notion of ϵ_0 -net and bound all query points in net. In Section H.3, we show how we can bound all points in the query space by introducing an additive error. In Section H.4, we examine the effects of different norms on our adaptive query proof.

First, from Theorem G.3, given query $y \in [0, R]^d, \alpha \in (0, 1)$ we have DISTANCEQUERY(y, α) that can solve d -dimension weighted ℓ_1 -distance problem with constant probability 0.99. Now we show how to improve it to solve adaptive query problem.

H.1 BOOST THE CONSTANT PROBABILITY TO HIGH PROBABILITY

We can repeat the data structure multiple times and take the median to boost the constant probability using Chernoff bound from Lemma B.2.

Lemma H.1 (Using Chernoff bound to boost the probability). *If the following conditions hold:*

- Let data set $X \in [0, R]^{n \times d}$, weights $w \in [-R_w, R_w]^n$, query $y \in [0, R]^d$.
- Let relative error parameter $\alpha \in (0, 1)$, the failure probability $p_f \in (0, 0.01)$.

- We create $l = O(\log(1/p_f))$ independent copies of data structure `DPTREEHIGHDIM` and take the median of the outputs with each data structure instantiated with $(\epsilon/l, (\delta + \delta')/l)$ -DP.
- Let $A_* = \sum_{i \in [n]} w_i \|y - x_i\|_1$ be the true answer.
- Let $B = O(\epsilon^{-1} \alpha^{-1/2} l R R_w d \sqrt{\log(l/\delta')} \cdot \log^{3/2} n)$.

Then for each fixed query point y , we can have the process of outputting the median of l responses is $(\epsilon, \delta + \delta')$ -DP and the error is upper bounded by $\alpha A_* + B$ with probability $1 - p_f$.

Proof. By basic composition Fact B.6, we prove the DP. Similar to the proof of Theorem D.2, we prove the error by Chernoff bound (Lemma B.2). \square

H.2 FROM EACH FIXED QUERY POINT TO ALL ON-NET POINTS

In this section, we build ϵ_0 -net and generalize from each fixed query point to all on-net points.

Definition H.2 (ℓ_p ϵ_0 -net, see Definition 4.2.1 in Vershynin (2017)). We define N be ℓ_p ϵ_0 -net of $\mathcal{B} := \{q \in [0, R]^d\}$ such that, for every point q in \mathcal{B} , there exists $y \in N$ satisfying $\|y - q\|_p \leq \epsilon_0$.

Fact H.3 (ℓ_∞ ϵ_0 -net). Let N be the ℓ_∞ ϵ_0 -net of \mathcal{B} , and $|N|$ be the size of net N . We have $|N| \leq (5R/\epsilon_0)^d$.

Fact H.4 (ℓ_2 ϵ_0 -net, see Lemma 5 in Woodruff (2014)). Let N be the ℓ_2 ϵ_0 -net of \mathcal{B} , and $|N|$ be the size of net N . We have $|N| \leq (5R/\epsilon_0)^d$.

Fact H.5 (ℓ_1 ϵ_0 -net, see Theorem 2 in Guntuboyina & Sen (2012)). Let N be the ℓ_1 ϵ_0 -net of \mathcal{B} , and $|N|$ be the size of net N . We have $|N| \leq (5R\sqrt{d}/\epsilon_0)^d$.

Lemma H.6 (From for each query point to for all points in net). If the following conditions hold:

- Let N be the ℓ_∞ ϵ_0 -net of \mathcal{B} , and $|N|$ be the size of net N .
- Let data set $X \in [0, R]^{n \times d}$, weights $w \in [-R_w, R_w]^n$, query $y \in [0, R]^d$.
- Let relative error parameter $\alpha \in (0, 1)$, the failure probability $p_f \in (0, 0.01)$.
- We create $l = O(\log(|N|/p_f))$ independent copies of data structure `DPTREEHIGHDIM` and take the median of the outputs with each data structure instantiated with $(\epsilon/l, (\delta + \delta')/l)$ -DP.
- Let $A_* = \sum_{i \in [n]} w_i \|y - x_i\|_1$ be the true answer.
- Let $B = O(\epsilon^{-1} \alpha^{-1/2} l R R_w d \sqrt{\log(l/\delta')} \cdot \log^{3/2} n)$.

Then with probability $1 - p_f$, for all query points $y \in N$, we can have the process of outputting the median of l responses is $(\epsilon, \delta + \delta')$ -DP and the error is upper bounded by $\alpha A_* + B$.

Proof. By basic composition Fact B.6, we prove the DP. From Lemma H.1, we know for each $y \in N$, the error is upper bounded by $\alpha A_* + B$ with probability $1 - p_f/|N|$.

Then, by union bound, with probability $1 - p_f$, the error of all $|N|$ query points in the net $y \in N$ is upper bounded by $\alpha A_* + B$. \square

H.3 FROM NET POINTS TO ALL POINTS

In this section, we show how to generalize points from net to all points in the query space. Since adaptive query must lie in this space, we complete the proof of adaptive query.

Lemma H.7 (Lipschitz of query function). If the following conditions hold:

- Let data set $X \in [0, R]^{n \times d}$, weights $w \in [-R_w, R_w]^n$, query $y \in [0, R]^d$.
- Let $Z(y) := \sum_{i \in [n]} w_i \|y - x_i\|_1$.

- Let $L = nR_w$.

Then, we have $Z(y)$ is L -Lipschitz (note that we have ℓ_1 Lipschitz here).

Proof. We can show

$$\begin{aligned}
 |Z(y) - Z(\tilde{y})| &= \left| \sum_{i \in [n]} w_i \|y - x_i\|_1 - \sum_{i \in [n]} w_i \|\tilde{y} - x_i\|_1 \right| \\
 &\leq \sum_{i \in [n]} |w_i| \cdot \left| \|y - x_i\|_1 - \|\tilde{y} - x_i\|_1 \right| \\
 &\leq \sum_{i \in [n]} |w_i| \cdot \|y - \tilde{y}\|_1 \\
 &= nR_w \cdot \|y - \tilde{y}\|_1
 \end{aligned}$$

where the first step follows from definition of $Z(y)$, the second step follows from triangular inequality, the third step follows from reverse triangular inequality, the fourth step follows from $w \in [-R_w, R_w]^n$. \square

Lemma H.8 (From points in net to all points in query space). *If the following conditions hold:*

- Let N be the ℓ_∞ ϵ_0 -net of \mathcal{B} , and $|N|$ be the size of net N .
- Let data set $X \in [0, R]^{n \times d}$, weights $w \in [-R_w, R_w]^n$, query $y \in [0, R]^d$.
- Let relative error parameter $\alpha \in (0, 1)$, the failure probability $p_f \in (0, 0.01)$.
- We create $l = O(\log((R/\epsilon_0)^d/p_f))$ independent copies of data structure $\{\text{DPTREEHIGHDIM}_j\}_{j=1}^l$ and take the median of the outputs with each data structure instantiated with $(\epsilon/l, (\delta + \delta')/l)$ -DP.
- Let $f(y) := \text{Median}(\{\text{DPTREEHIGHDIM}_j.\text{DISTANCEQUERY}(y, \alpha)\}_{j=1}^l)$.
- Let $Z(y) := \sum_{i \in [n]} w_i \|y - x_i\|_1$, where $Z(y)$ is L -Lipschitz with $L = nR_w$.
- Let $B = O(\epsilon^{-1} \alpha^{-1/2} l R R_w d \sqrt{\log(l/\delta')} \cdot \log^{3/2} n)$.

Then with probability $1 - p_f$, for all query points $q \in \mathcal{B}$, there exists a point $y \in N$ which is the closest to q , we can have the process of outputting the median of l responses is $(\epsilon, \delta + \delta')$ -DP and the error satisfy

$$|f(y) - Z(q)| \leq \alpha Z(q) + B + 2Ld\epsilon_0.$$

Proof. By basic composition Fact B.6, we prove the DP.

We define an event E such that:

$$\begin{aligned}
 &\forall y \in N \\
 &|f(y) - Z(y)| \leq \alpha Z(y) + B.
 \end{aligned}$$

From Lemma H.1, with $l = O(\log(|N|/p_f))$ we know

$$\Pr[\text{event } E \text{ holds}] \geq 1 - p_f$$

We can show

$$\begin{aligned}
 l &= O(\log(|N|/p_f)) \\
 &= O(\log((R/\epsilon_0)^d/p_f))
 \end{aligned}$$

where the first step follows from definition of l , the second step follows from Fact H.3.

We condition on event E to be held. Then, by definition of ℓ_∞ ϵ_0 -net (see Definition H.2), for each $q \notin N$, there exists $y \in N$ such that

$$\|y - q\|_\infty \leq \epsilon_0 \quad (6)$$

We know

$$\begin{aligned} |Z(y) - Z(q)| &\leq L \cdot \|y - q\|_1 \\ &\leq L \cdot d \|y - q\|_\infty \\ &\leq L \cdot d \epsilon_0 \end{aligned} \quad (7)$$

where the first step follows from Lemma H.7, the second step follows from $\|x\|_1 \leq d\|x\|_\infty$ for $x \in \mathbb{R}^d$, and the last step follows from Eq. (6).

Using the on-net query y to answer the off-net query q , for any $q \notin N$, we have

$$\begin{aligned} |f(y) - Z(q)| &\leq |f(y) - Z(y)| + |Z(y) - Z(q)| \\ &\leq |f(y) - Z(y)| + L \cdot d \cdot \epsilon_0 \\ &\leq \alpha Z(y) + B + L \cdot d \cdot \epsilon_0 \\ &\leq \alpha Z(q) + B + 2L \cdot d \cdot \epsilon_0 \end{aligned} \quad (8)$$

where the first step follows from triangular inequality, the second step follows from Eq. (7), the third step follows from Lemma H.6, and the last step follows from Eq. (7).

Thus, we complete the proof. \square

Therefore, even adaptive queries can be answered accurately, since any adaptive query can be assumed in \mathcal{B} .

H.4 EFFECT OF DIFFERENT NORMS ON THE RESULT

In the above proof, we have two different measure spaces, i.e. ℓ_∞ distance of ϵ_0 -net (Definition H.2) and ℓ_1 Lipschitz (Lemma H.7).

One might ask, will the norm we choose in two spaces have an impact on the final result? We can show that the norm we choose currently is sufficient to use.

For different norms, the only differences in the proofs will be Lipschitz smoothness in Eq. (7) and the cardinality of ϵ_0 -net, i.e. $|N|$ in Fact H.3.

Lemma H.9. *If we use ℓ_∞ ϵ_0 -net and use ℓ_1 Lipschitz in Lemma H.8, we have copies of data structure $l = O(d \log(nR/p_f))$.*

Proof. If we use ℓ_∞ to bound the distance to net, Eq. (7) is:

$$\begin{aligned} |Z(y) - Z(q)| &\leq nR_w \cdot \|y - q\|_1 \\ &\leq nR_w \cdot d \|y - q\|_\infty \\ &\leq nR_w \cdot d \epsilon_0 \end{aligned}$$

where the first step follows from Lemma H.7, the second step follows from $\|x\|_1 \leq d\|x\|_\infty$ for $x \in \mathbb{R}^d$, and the last step follows from ℓ_∞ ϵ_0 -net.

Then, Eq. (8) is

$$|f(y) - Z(q)| \leq \alpha Z(q) + B + 2nR_w \cdot d \cdot \epsilon_0$$

For ℓ_∞ distance, we have $|N| \leq (5R/\epsilon_0)^d$ in Fact H.3.

We can choose $\epsilon_0 = \Theta(1/n)$ to hide $nR_w \cdot d \cdot \epsilon_0$ term in B in Lemma H.8. Thus,

$$\begin{aligned} l &= O(\log(|N|/p_f)) \\ &= O(\log((R/\epsilon_0)^d/p_f)) \\ &= O(\log((nR)^d/p_f)) \\ &= O(d \log(nR/p_f)) \end{aligned}$$

where the last step follows from $\log(a^d/b) = O(d \log(a/b))$ for any $a > 1, 0 < b < 1, d > 1$. \square

Lemma H.10. *If we use ℓ_2 ϵ_0 -net and use ℓ_1 Lipschitz in Lemma H.8, we have copies of data structure $l = O(d \log(nR/p_f))$.*

Proof. If we use ℓ_2 to bound the distance to net, Eq. (7) changes to be:

$$\begin{aligned} |Z(y) - Z(q)| &\leq nR_w \cdot \|y - q\|_1 \\ &\leq nR_w \cdot \sqrt{d} \cdot \|y - q\|_2 \\ &\leq nR_w \cdot \epsilon_0 \sqrt{d} \end{aligned}$$

where the first step follows from Lemma H.7, the second step follows from $\|x\|_1 \leq \sqrt{d} \cdot \|x\|_2$ for $x \in \mathbb{R}^d$, and the last step follows from ℓ_2 ϵ_0 -net.

Then, Eq. (8) changes to be

$$|f(y) - Z(q)| \leq \alpha Z(q) + B + 2nR_w \cdot \epsilon_0 \sqrt{d}$$

For ℓ_2 distance, we also have $|N| \leq (5R/\epsilon_0)^d$ in Fact H.4.

We can choose $\epsilon_0 = \Theta(1/n)$ to hide $nR_w \cdot \sqrt{d} \cdot \epsilon_0$ term in B in Lemma H.8. Thus,

$$\begin{aligned} l &= O(\log(|N|/p_f)) \\ &= O(\log((R/\epsilon_0)^d/p_f)) \\ &= O(\log((nR)^d/p_f)) \\ &= O(d \log(nR/p_f)) \end{aligned}$$

where the last step follows from $\log(a^d/b) = O(d \log(a/b))$ for any $a > 1, 0 < b < 1, d > 1$. \square

Lemma H.11. *If we use ℓ_1 ϵ_0 -net and use ℓ_1 Lipschitz in Lemma H.8, we have copies of data structure $l = O(d \log(ndR/p_f))$.*

Proof. If we use ℓ_1 to bound the distance to net, Eq. (7) changes to be:

$$\begin{aligned} |Z(y) - Z(q)| &\leq nR_w \cdot \|y - q\|_1 \\ &\leq nR_w \cdot \epsilon_0 \end{aligned}$$

where the first step follows from Lemma H.7, and the last step follows from ℓ_1 ϵ_0 -net.

Then, Eq. (8) changes to be

$$|f(y) - Z(q)| \leq \alpha Z(q) + B + 2nR_w \cdot \epsilon_0$$

For ℓ_1 distance, we have $|N| \leq (5R\sqrt{d}/\epsilon_0)^d$.

We can choose $\epsilon_0 = \Theta(1/n)$ to hide $nR_w \cdot \epsilon_0$ term in B in Lemma H.8. Thus,

$$\begin{aligned} l &= O(\log(|N|/p_f)) \\ &= O(\log((R\sqrt{d}/\epsilon_0)^d/p_f)) \\ &= O(\log((nR\sqrt{d})^d/p_f)) \\ &= O(d \log(nRd/p_f)) \end{aligned}$$

where the last step follows from $\log(a^d/b) = O(d \log(a/b))$ for any $a > 1, 0 < b < 1, d > 1$. \square

From the above analysis, we can show that ℓ_∞ or ℓ_2 ϵ_0 -net is slightly better than ℓ_1 ϵ_0 -net.

- ℓ_∞ ϵ_0 -net, Lemma H.9: we have $l = O(d \log(nR/p_f))$.
- ℓ_2 ϵ_0 -net, Lemma H.10: we have $l = O(d \log(nR/p_f))$.
- ℓ_1 ϵ_0 -net, Lemma H.11: we have $l = O(d \log(nRd/p_f))$.

Thus, the norm we choose for ϵ_0 -net is sufficient good.

I SOFTMAX ACTIVATION

In this section, we introduce how we extend previous ℓ_1 distance results to the Softmax activation function, which is the most widely used distance measure in attention mechanism based models.

In Section I.1, we show how to extend to the Softmax distance function in Lemma I.6. In Section I.2, we show how to adjust our algorithms. In Section I.3, we extend our algorithm to be robust to adaptive query. In Section I.4, we give the proof of our main result Theorem 3.1.

I.1 EXPONENTIAL INNER PRODUCT

In this section, we show how we obtain the Softmax distance using ℓ_2^2 distance query. First, we provide some helpful results from Alman & Song (2023).

Definition I.1 (Definition 3.1 in Alman & Song (2023)). *Let $r \geq 1$ denote a positive integer. Let $\epsilon \in (0, 0.1)$ denote an accuracy parameter. Given a matrix $A \in \mathbb{R}_{\geq 0}^{n \times n}$, we say $\tilde{A} \in \mathbb{R}_{\geq 0}^{n \times n}$ is an (ϵ, r) -approximation of A if*

- $\tilde{A} = U_1 \cdot U_2^\top$ for some matrices $U_1, U_2 \in \mathbb{R}^{n \times r}$ (i.e., \tilde{A} has rank at most r), and
- $|\tilde{A}_{i,j} - A_{i,j}| \leq \epsilon \cdot A_{i,j}$ for all $(i, j) \in [n]^2$.

Lemma I.2 (Lemma 3.4 in Alman & Song (2023)). *Suppose $Q, K \in \mathbb{R}^{n \times d}$, with $\|Q\|_\infty \leq R$, and $\|K\|_\infty \leq R$. Let $A := \exp(QK^\top/d) \in \mathbb{R}^{n \times n}$. For accuracy parameter $\epsilon \in (0, 0.1)$, there is a positive integer s bounded above by*

$$s = O\left(\max\left\{\frac{\log(1/\epsilon)}{\log(\log(1/\epsilon)/R)}, R^2\right\}\right), \quad (9)$$

and a positive integer r bounded above by

$$r \leq \binom{2s + 2d}{2s} \quad (10)$$

such that: There is a matrix $\tilde{A} \in \mathbb{R}^{n \times n}$ that is an (ϵ, r) -approximation (Definition I.1) of $A \in \mathbb{R}^{n \times n}$. Furthermore, the matrices U_1 and U_2 defining \tilde{A} can be computed in $O(n \cdot r)$ time.

Here we consider the vector version of Lemma I.2.

Definition I.3. We define $\Gamma_{R,s} := \max_{j \in [s]} \frac{R^j}{\sqrt{j}!}$.

Then, we have $P(x) : [0, R]^d \rightarrow [0, \Gamma_{R,s}]^r$ where $P(\cdot)$ is polynomial kernel function defined in Alman & Song (2023).

Remark I.4. We use $\Gamma_{R,s}$ to denote the value range of our polynomial kernel methods function, i.e., $P(x) : [0, R]^d \rightarrow [0, \Gamma_{R,s}]^r$. The factorial term in $\Gamma_{R,s}$ comes from Taylor approximation coefficients. We take the maximum overall s order approximation terms to get the upper bound of our value range.

We use the polynomial approximation method, which has been applied to accelerate Transformer model extensively Alman & Song (2023; 2024a;b); Liang et al. (2024e;b).

Lemma I.5 (Polynomial approximation). *For any accuracy parameter $\epsilon_s \in (0, 0.1)$, let $R \geq 1$, and let $P(x) : [0, R]^d \rightarrow [0, \Gamma_{R,s}]^r$ be the s -th order polynomial kernel function defined in Alman & Song (2023) where $r \leq \binom{2s+2d}{2s}$ and $s = O(\max\{\frac{\log(1/\epsilon_s)}{\log(\log(1/\epsilon_s)/R)}, R^2\})$. Then, for any $x, y \in [0, R]^d$, we have*

$$|P(x)^\top P(y) - \exp(x^\top y/d)| \leq \epsilon_s \cdot \min\{\exp(x^\top y/d), P(x)^\top P(y)\}$$

Furthermore, the vectors $P(x)$ and $P(y)$ can be computed in $O(r)$ time.

Proof. Let $n = 1$. The proof follows from directly applying Lemma I.2. \square

Using the results from Alman & Song (2023) above, we can extend our results to Softmax activation.

Lemma I.6 (Weighted Softmax approximation, formal version of Lemma C.7). *Let accuracy parameter be $\epsilon_s \in (0, 0.1)$. Let $R \geq 1$. Let $r \leq \binom{2s+2d}{2s}$ and $s = O(\max\{\frac{\log(1/\epsilon_s)}{\log(\log(1/\epsilon_s)/R)}, R^2\})$. Let $P(x) : [0, R]^d \rightarrow [0, \Gamma_{R,s}]^r$ be the s -th order polynomial kernel function defined in Lemma I.5. Then we can approximate exponential inner product using polynomial kernel function:*

$$w^\top \exp(Xy/d) = -\frac{1}{2} \sum_{j \in [r]} \sum_{i \in [n]} w_i |P(x_i)_j - P(y)_j|^2 + \frac{1}{2} \sum_{i \in [n]} w_i (\|P(x_i)\|_2^2 + \|P(y)\|_2^2) + O(w^\top \exp(Xy/d) \cdot \epsilon_s)$$

Moreover, the vectors $P(\cdot)$ can be computed in $O(r)$ time.

Proof. From Lemma I.5, we can use polynomial kernel to approximate the Softmax function:

$$w^\top \exp(Xy/d) = \sum_{i \in [n]} w_i P(x_i)^\top P(y) + O(w^\top \exp(Xy/d) \cdot \epsilon_s).$$

The proof of approximation error and time complexity of constructing $P(\cdot)$ follows from Lemma I.5.

Then, we can show

$$\begin{aligned} 2 \sum_{i \in [n]} w_i P(x_i)^\top P(y) &= - \sum_{i \in [n]} w_i \|P(x_i) - P(y)\|_2^2 + \sum_{i \in [n]} w_i (\|P(x_i)\|_2^2 + \|P(y)\|_2^2) \\ &= - \sum_{j \in [r]} \sum_{i \in [n]} w_i |P(x_i)_j - P(y)_j|^2 + \sum_{i \in [n]} w_i (\|P(x_i)\|_2^2 + \|P(y)\|_2^2) \end{aligned}$$

where the first step follows from $\|x - y\|_2^2 = \|x\|_2^2 + \|y\|_2^2 - 2\langle x, y \rangle$, and the second step follows $\|x\|_2^2 = \sum_{j=1}^d |x_j|^2$ for $x \in \mathbb{R}^d$. \square

I.2 ALGORITHM MODIFICATIONS

Based on Lemma I.6, we can now extend our DP algorithms to handle Softmax activation. First, we need to construct $P(y)$ and $P(x_i)$ for $i \in [n]$, each costing $O(r)$ time. Then, for the second term in Lemma I.6, i.e. $\frac{1}{2} \sum_{i \in [n]} w_i (\|P(x_i)\|_2^2 + \|P(y)\|_2^2)$, we don't need to add DP noises in it; instead, we calculate this term exactly, preprocess it, and store the results in the algorithm. For the first term, $-\frac{1}{2} \sum_{j \in [r]} \sum_{i \in [n]} w_i |P(x_i)_j - P(y)_j|^2$, we can adjust our high dimensional DP distance query algorithm to solve it.

Due to the decomposability of ℓ_p^p norm, i.e.

$$\sum_{i \in [n]} w_i \|x_i - y\|_p^p = \sum_{j \in [d]} \sum_{i \in [n]} w_i |x_{i,j} - y_j|^p,$$

we can compute ℓ_2^2 norm easily (see details in Lemma E.2). We then show how to extend our one dimensional ℓ_1 distance algorithm (Algorithm 5 and 6) to ℓ_2^2 distance with minor modifications.

Theorem I.7 (DPTREEDISTANCE ℓ_2^2 distance). *With α scaled down by a factor of 2 and all QUERY instead multiplied by $R^2/(1 + \alpha/2)^{2j}$ in Lines 8 and 13 of Algorithm 6, i.e., from*

$$\bullet \text{ Lines 8 and 13: Value} \leftarrow \text{Value} + \mathcal{D}.\text{QUERY}(l_j, r_j) \cdot \frac{R}{(1+\alpha)^j}$$

to

$$\bullet \text{ Lines 8 and 13: Value} \leftarrow \text{Value} + \mathcal{D}.\text{QUERY}(l_j, r_j) \cdot \frac{R^2}{(1+\alpha/2)^{2j}}.$$

The data structure DPTREEDISTANCE (Algorithm 5,6) uses $O(n)$ spaces to solve weighted ℓ_2^2 -distance query problem for dataset $X \subset [0, R]$ and support the following operations:

- INIT($X \subset [0, R], n \in \mathbb{N}_+, w \in [-R_w, R_w]^n, \epsilon \in (0, 1), \delta \in (0, 1)$). (Algorithm 5) It takes $O(n)$ time to initialize the data structure.

- **DISTANCEQUERY**($y \in [0, R]$, $\alpha \in (0, 1)$). (Algorithm 6)

It takes $O(\alpha^{-1} \log^2 n)$ time to output a number z such that

- the process of output z satisfies (ϵ, δ) -DP private, which computes $\sum_{i \in [n]} w_i |y - x_i|^2$,
- $|z - \sum_{i \in [n]} w_i |y - x_i|^2| \leq \alpha \sum_{i \in [n]} w_i |y - x_i|^2 + O(\epsilon^{-1} \alpha^{-1/2} R^2 R_w \log^{3/2} n)$,
- it holds with probability 0.99.

Proof. The proof is similar to that of Theorem F.6, except that now our additive error includes R increased by a power of 2, i.e., from $O(\epsilon^{-1} \alpha^{-1/2} R R_w \log^{3/2} n)$ to $O(\epsilon^{-1} \alpha^{-1/2} R^2 R_w \log^{3/2} n)$. \square

Now we can give our result that can answer Softmax query.

Theorem I.8 (Softmax query, formal version of Theorem 4.2). *Let $R \geq 1$. Let $r \leq \binom{2s+2d}{2s}$ and $s = O(\max\{\frac{\log(1/\epsilon_s)}{\log(\log(1/\epsilon_s)/R)}, R^2\})$. Let $\Gamma_{R,s}$ be defined in Definition I.3. Let accuracy parameter be $\epsilon_s \in (0, 0.1)$. There is a data structure **DPTREESOFTMAX** (Algorithm 3) that uses $O(nr)$ spaces to solve Softmax query problem for dataset $X \subset [0, R]^d$ and support the following operations:*

- **INIT**($X \subset [0, R]^d$, $n \in \mathbb{N}_+$, $w \in [-R_w, R_w]^n$, $\epsilon \in (0, 1)$, $\delta \in (0, 1)$, $\delta' \in (0, 1)$, $c \in (0, 0.1)$, $\epsilon_s \in (0, 0.1)$). (Algorithm 3) It takes $O(nr)$ time to initialize the data structure.
- **DISTANCEQUERY**($y \in [0, R]^d$, $\alpha \in (0, 1)$). (Algorithm 3) It takes $O(\alpha^{-1} r \log^2 n)$ time to output a number z such that
 - the process of output z satisfies $(\epsilon, \delta + \delta')$ -DP private, which computes $w^\top \exp(Xy/d)$,
 - $|z - w^\top \exp(Xy/d)| \leq (\alpha + \epsilon_s) \cdot w^\top \exp(Xy/d) + O(\epsilon^{-1} \alpha^{-1/2} \Gamma_{R,s}^2 R_w r \sqrt{\log(1/\delta')} \cdot \log^{3/2} n)$,
 - it holds with probability 0.99.

Proof. The DP proof is the same as the proof of Lemma G.1.

We then show the time complexity. From Lemma I.6, we know that constructing $P(\cdot)$ requires $O(r)$ time. In the first for loop of **INIT**, the dominating time consumption is $O(nr)$. The second for loop also has a time complexity of $O(nr)$. Therefore, the total time complexity for **INIT** is $O(nr)$. In the **DISTANCEQUERY** function, constructing $P(y)$ takes $O(r)$ time. Within the for loop, it requires $O(\alpha^{-1} r \log^2 n)$. Thus, the total time complexity for **DISTANCEQUERY** is $O(\alpha^{-1} r \log^2 n)$.

The space complexity is $O(nr)$, since storing the $n \times r$ matrix P is the dominating factor.

The proof of the error follows from the triangle inequality by combining the errors in Lemma I.6 and Theorem I.7. \square

I.3 ADAPTIVE SOFTMAX

In this section, we show how to make Algorithm 3 robust to adaptive query. We follow the same idea from Section H. We notice that, in the Softmax activation, we have query function $Z(y) := w^\top \exp(Xy/d)$ different from the ℓ_1 -distance in Section H. Therefore, we need to recalculate Lipschitz constant first.

Lemma I.9 (Lipschitz of weighted Softmax). *If the following conditions hold:*

- Let data set $X \in [0, R]^{n \times d}$, weights $w \in [-R_w, R_w]^n$, query $y \in [0, R]^d$.
- Let $Z(y) := w^\top \exp(Xy/d)$.
- Let $L = nd^{-1/2} R R_w \exp(R^2)$.

Then, we have $Z(y)$ is L -Lipschitz (note that we have ℓ_1 Lipschitz here).

Proof. We can show

$$\begin{aligned}
|Z(y) - Z(\tilde{y})| &= \left| \sum_{i \in [n]} w_i \exp(x_i^\top y/d) - \sum_{i \in [n]} w_i \exp(x_i^\top \tilde{y}/d) \right| \\
&\leq \sum_{i \in [n]} |w_i| \cdot |\exp(x_i^\top y/d) - \exp(x_i^\top \tilde{y}/d)| \\
&\leq \sum_{i \in [n]} |w_i| \exp(R^2) |x_i^\top y/d - x_i^\top \tilde{y}/d| \\
&\leq \sum_{i \in [n]} |w_i| \exp(R^2) \|x_i\|_2 \cdot \|y - \tilde{y}\|_2/d \\
&\leq nR_w \exp(R^2) \sqrt{d}R \cdot \|y - \tilde{y}\|_2/d \\
&\leq nd^{-1/2} RR_w \exp(R^2) \|y - \tilde{y}\|_1
\end{aligned}$$

where the first step follows from definition of $Z(y)$, $Z(\tilde{y})$, the second step follows from triangular inequality, the third step follows from Fact B.4, the fourth step follows from Cauchy-Schwarz inequality $|u^\top v| \leq \|u\|_2 \cdot \|v\|_2$ for $u, v \in \mathbb{R}^d$, the fifth step follows from $w_i \in [-R_w, R_w]$ and $x_i \in [0, R]^d$, and the last step follows from $\|u\|_2 \leq \|u\|_1$ for $u \in \mathbb{R}^d$. \square

Then we can show how to extend our algorithm to be robust to adaptive query.

Lemma I.10 (Adaptive Softmax, formal version of Lemma C.8). *If the following conditions hold:*

- Let N be the ℓ_∞ ϵ_0 -net of \mathcal{B} , and $|N|$ be the size of net N .
- Let data set $X \in [0, R]^{n \times d}$, weights $w \in [-R_w, R_w]^n$, query $y \in [0, R]^d$.
- Let relative error parameter $\alpha \in (0, 1)$, the failure probability $p_f \in (0, 0.01)$.
- We create $l = O(\log((R/\epsilon_0)^r/p_f))$ independent copies of data structure $\{\text{DPTREESOFTMAX}_j\}_{j=1}^l$ (Algorithm 3) and take the median of the outputs with each data structure instantiated with $(\epsilon/l, (\delta + \delta')/l)$ -DP.
- Let $f(y) := \text{Median}(\{\text{DPTREESOFTMAX}_j.\text{DISTANCEQUERY}(y, \alpha)\}_{j=1}^l)$.
- Let $Z(y) := w^\top \exp(Xy/d)$, where $Z(y)$ is L -Lipschitz with $L = nd^{-1/2} RR_w \exp(R^2)$.
- Let $B = O(\epsilon^{-1} \alpha^{-1/2} l \Gamma_{R,s}^2 R_w r \sqrt{\log(l/\delta')} \cdot \log^3 n)$.

Then with probability $1 - p_f$, for all query points $q \in \mathcal{B}$, there exists a point $y \in N$ which is the closest to q , we can have the process of outputting the median of l responses is $(\epsilon, \delta + \delta')$ -DP and the error satisfies

$$|f(y) - Z(q)| \leq (\alpha + \epsilon_s)Z(q) + B + 2n\sqrt{d}RR_w \exp(R^2)\epsilon_0.$$

Proof. The proof follows from the same idea as the proof of Lemma H.8, except that we use Theorem I.8 and the Lipschitz in Lemma I.9. \square

Theorem I.11 (Adaptive query Softmax data structure, formal version of Theorem 4.4). *Let $R \geq 1$. Let $r \leq \binom{2s+2d}{2s}$ and $s = O(\max\{\frac{\log(1/\epsilon_s)}{\log(\log(1/\epsilon_s)/R)}, R^2\})$. Let $\Gamma_{R,s}$ be defined in Definition I.3. Let accuracy parameter be $\epsilon_s \in (0, 0.1)$. Let $X \in [0, R]^{n \times d}$ be the dataset, $w \in [-R_w, R_w]^n$ be weights, $y \in [0, R]^d$ be the query, $\alpha \in (0, 1)$ be the relative error parameter, and p_f be the failure probability parameter. Let $l = O(r \log(dR/(\epsilon_s p_f)))$. There is a data structure $\text{DPTREESOFTMAXADAPTIVE}$ (Algorithm 1) that uses $O(\ln r)$ spaces to solve weighted Softmax query problem for dataset $X \subset [0, R]^d$ and support the following operations:*

- $\text{INIT}(X \subset [0, R]^d, n \in \mathbb{N}_+, w \in [-R_w, R_w]^n, \epsilon \in (0, 1), \delta \in (0, 1), \delta' \in (0, 1), c \in (0, 0.1), \epsilon_s \in (0, 0.1), p_f \in (0, 0.01))$. (Algorithm 1) It takes $O(\ln d)$ time to initialize the data structure.

- **DISTANCEQUERY**($y \in [0, R]^d, \alpha \in (0, 1)$). (Algorithm 1) It takes $O(\alpha^{-1}ld \log^2 n)$ time to output a number z such that
 - the process of output z satisfies $(\epsilon, \delta + \delta')$ -DP private, which computes $w^\top \exp(Xy/d)$,
 - $|z - w^\top \exp(Xy/d)| \leq (\alpha + \epsilon_s) \cdot w^\top \exp(Xy/d) + O(\epsilon^{-1}\alpha^{-1/2}l\Gamma_{R,s}^2 R_w r \sqrt{\log(l/\delta')} \cdot \log^{3/2} n)$,
 - it holds with probability $1 - p_f$ (where p_f is used in l),
 - it is robust to adaptive query.

Proof. We only need to show how to pick ϵ_0 in the parameter l , because everything else is the same as Lemma I.10. We know the additive error introduced by adaptive query is $E_a := O(n\sqrt{d}RR_w \exp(R^2)\epsilon_0)$ and the relative error introduced by polynomial kernel approximation is $E_p := w^\top \exp(Xy/d) \cdot \epsilon_s$. It can be shown that:

$$\begin{aligned} E_p &:= w^\top \exp(Xy/d) \cdot \epsilon_s \\ &\leq \epsilon_s \|w\|_2 \cdot \|\exp(Xy/d)\|_2 \\ &= O(nR_w \epsilon_s \exp(R^2)) \end{aligned}$$

where the first step follows from definition of E_p , the second step follows from Cauchy–Schwarz inequality, and the last step follows from $w \in [-R_w, R_w]^n$, $X \in [0, R]^{n \times d}$, and $y \in [0, R]^d$.

Picking $\epsilon_0 = \Theta(\frac{\epsilon_s}{\sqrt{dR}})$, we can hide the error of adaptive query E_a in E_p . Thus, we have

$$\begin{aligned} l &= O(\log((R/\epsilon_0)^r / p_f)) \\ &= O(\log((\sqrt{d}R^2/\epsilon_s)^r / p_f)) \\ &= O(r \log(dR/(\epsilon_s p_f))) \end{aligned}$$

where the first step comes from the definition of l , the second step comes from picking $\epsilon_0 = \Theta(\frac{\epsilon_s}{\sqrt{dR}})$, and the last step follows from $\log(a^d/b) = O(d \log(a/b))$ for any $a > 1, 0 < b < 1, d > 1$. \square

I.4 PROOF OF MAIN RESULT

In this section, we give the proof of our main result of Theorem 3.1.

Theorem I.12 (Softmax cross-attention, formal version of Theorem 3.1). *Let Q, K, V, Attn be defined in Definition I.1. Let $\alpha \in (0, 1)$ be the relative error parameter and p_f be the probability of failure parameter. Let r, s, ϵ_s be parameters of polynomial kernel methods (Lemma C.7). Let $\Gamma_{R,s} := \max_{j \in [s]} \frac{R^j}{\sqrt{j!}}$ (Definition I.3). Let $l = O(r \log(dR/(\epsilon_s p_f)))$. There is a data structure **DPTREE****SOFTMAX****ADAPTIVE** (Algorithm 1) that uses $O(\ln rd)$ spaces to make cross-attention DP and supports the following operations:*

- We initialize d data structures using **INIT**($K, n, V_{*,k}, \epsilon \in (0, 1), \delta \in (0, 1), \delta' \in (0, 1), c \in (0, 0.1), \epsilon_s \in (0, 0.1), p_f \in (0, 0.01)$) (Algorithm 1), for $k \in [d]$. It takes $O(\ln r)$ time to initialize one data structure.
- At query time, for user input Q , we process one token at a time by passing the i -th row of Q , denoted $Q_i \in \mathbb{R}^d$, to **DISTANCEQUERY**($Q_i, \alpha \in (0, 1)$) (Algorithm 1) for each $i \in [m]$. It takes $O(\alpha^{-1}lr \log^2 n)$ time to output an entry z in $\text{Attn}(Q, K, V)$ such that
 - the process of output z satisfies $(\epsilon, \delta + \delta')$ -DP,
 - the process of output z has relative error $n^{-1}(\alpha + \epsilon_s)$,
 - the process of output z has additive error $O(n^{-1}\epsilon^{-1}\alpha^{-1/2}l\Gamma_{R,s}^2 R_w r \sqrt{\log(l/\delta')} \cdot \log^{3/2} n)$,
 - it holds with probability $1 - p_f$ (where p_f is used in l),
 - it is robust to adaptive query.

Proof. From Section 3, we know that we can ensure matrix AV in cross-attention computation satisfies DP. Next, from Theorem 4.4, for $i \in [m], j \in [n], k \in [d]$, we have $(AV)_{i,k}$ is $(\epsilon, \delta + \delta')$ -DP and also robust to adaptive query. For D , we compute the exact true value. By the post-processing property of DP (Fact B.5), our cross-attention process is DP.

Let $(AV)_{i,k}$ be the true value and $(\widetilde{AV})_{i,k}$ be the noisy value. Then, from Theorem 4.4, we have

$$|(AV)_{i,k} - (\widetilde{AV})_{i,k}| \leq (\alpha + \epsilon_s) \cdot (AV)_{i,k} + O(\epsilon^{-1} \alpha^{-1/2} l \Gamma_{R,s}^2 R_w r \sqrt{\log(l/\delta')} \cdot \log^{3/2} n). \quad (11)$$

For $D_{i,i}$, we can show

$$D_{i,i} = (A \cdot \mathbf{1}_n)_i = \sum_{j=1}^n \exp(\langle Q_i, K_j \rangle / d) \geq n \quad (12)$$

because $\langle Q_i, K_j \rangle \geq 0$ for bounded Q, K .

Finally, we can show the error of one entry is bounded by

$$\begin{aligned} |(D^{-1}AV)_{i,k} - (D^{-1}\widetilde{AV})_{i,k}| &= |D_{i,i}^{-1}((AV)_{i,k} - (\widetilde{AV})_{i,k})| \\ &= |D_{i,i}^{-1}| \cdot |((AV)_{i,k} - (\widetilde{AV})_{i,k})| \\ &\leq n^{-1}(\alpha + \epsilon_s) \cdot (AV)_{i,k} \\ &\quad + O(n^{-1} \epsilon^{-1} \alpha^{-1/2} l \Gamma_{R,s}^2 R_w r \sqrt{\log(l/\delta')} \cdot \log^{3/2}(n)) \end{aligned}$$

where the first step follows from definition, the second step follows from simple algebra, and the last step follows from Eq.(11) and (12). \square