

000 001 002 003 004 005 006 007 008 009 010 011 012 013 014 015 016 017 018 019 020 021 022 023 024 025 026 027 028 029 030 031 032 033 034 035 036 037 038 039 040 041 042 043 044 045 046 047 048 049 050 051 052 053 ON DIFFERENTIALLY PRIVATE STRING DISTANCES

Anonymous authors

Paper under double-blind review

ABSTRACT

Given a database of bit strings $A_1, \dots, A_m \in \{0, 1\}^n$, a fundamental data structure task is to estimate the distances between a given query $B \in \{0, 1\}^n$ with all the strings in the database. In addition, one might further want to ensure the integrity of the database by releasing these distance statistics in a secure manner. In this work, we propose differentially private (DP) data structures for this type of tasks, with a focus on Hamming and edit distance. On top of the strong privacy guarantees, our data structures are also time- and space-efficient. In particular, our data structure is ϵ -DP against any sequence of queries of arbitrary length, and for any query B such that the maximum distance to any string in the database is at most k , we output m distance estimates. Moreover,

- For Hamming distance, our data structure answers any query in $\tilde{O}(mk + n)$ time and each estimate deviates from the true distance by at most $\tilde{O}(k/e^{\epsilon/\log k})$;
- For edit distance, our data structure answers any query in $\tilde{O}(mk^2 + n)$ time and each estimate deviates from the true distance by at most $\tilde{O}(k/e^{\epsilon/(\log k \log n)})$.

For moderate k , both data structures support sublinear query operations in the combined size of the query and its output. We obtain these results via a novel adaptation of the randomized response technique as a bit flipping procedure, applied to the sketched strings.

1 INTRODUCTION

Estimating string distances is one of the most fundamental problems in computer science and information theory, with rich applications in high-dimensional geometry, computational biology and machine learning. The problem could be generically formulated as follows: given a collection of strings $A_1, \dots, A_m \in \Sigma^n$ where Σ is the alphabet, the goal is to design a data structure to preprocess these strings such that when a query $B \in \Sigma^n$ is given, the data structure needs to quickly output estimates of $\|A_i - B\|$ for all $i \in [m]$, where $\|\cdot\|$ is the distance of interest. Assuming the symbols in Σ can allow constant time access and operations, a naive implementation would be to simply compute all the distances between A_i 's and B , which would require $O(mn)$ time. Designing data structures with $o(mn)$ query time has been the driving research direction in string distance estimations. To make the discussion concrete, in this work we will focus on binary alphabet ($\Sigma = \{0, 1\}$) and for distance, we will study Hamming and edit distance. Hamming distance (Hamming, 1950) is one of the most natural distance measurements for binary strings, with its deep root in error detecting and correction for codes. It finds large array of applications in database similarity searches (Indyk & Motwani, 1998; Charikar, 2002; Norouzi et al., 2012) and clustering algorithms (Huang, 1997; Huang & Ng, 1999).

Compared to Hamming distance, edit distance or the Levenshtein distance (Levenshtein, 1966) could be viewed as a more robust distance measurement for strings: it counts the minimum number of operations (including insertion, deletion and substitution) to transform from A_i to B . To see the robustness compared to Hamming distance, consider $A_i = (01)^n$ and $B = (10)^n$, the Hamming distance between these two strings is n , but A_i could be easily transformed into B by deleting the first bit and adding a 0 to the end, yielding an edit distance of 2. Due to its flexibility, edit distance is particularly useful for sequence alignment in computational biology (Wang et al., 2015; Young et al., 2021; Berger et al., 2021), measuring text similarity (Navarro, 2001; Sidorov et al., 2015) and

054 natural language processing, speech recognition (Fiscus et al., 2006; Droppo & Acero, 2010) and
 055 time series analysis (Martea, 2009; Gold & Sharir, 2018).

056 In addition to data structures with fast query times, another important consideration is to ensure the
 057 database is *secure*. Consider the scenario where the database consists of private medical data of m
 058 patients, where each of the A_i is the characteristic vector of n different symptoms. A malicious
 059 adversary might attempt to count the number of symptoms each patient has by querying $\mathbf{0}_n$, or
 060 detecting whether patient i has symptom j by querying e_j and $\mathbf{0}_n$ where e_j is the j -th standard basis
 061 in \mathbb{R}^n . It is hence crucial to curate a private scheme so that the adversary cannot distinguish the case
 062 whether the patient has symptom j or not. This notion of privacy has been precisely captured by
 063 *differential privacy* (Dwork, 2006; Dwork et al., 2006), which states that for neighboring databases,
 064 the output distribution of the data structure query should be close with high probability, hence any
 065 adversary cannot distinguish between the two cases.

066 Motivated by both privacy and efficiency concerns, we ask the following natural question:
 067

068 *Is it possible to design a data structure to estimate Hamming and edit distance, that are both
 069 differentially private, and time/space-efficient?*

070 We provide an affirmative answer to the above question, with the main results summarized in the
 071 following two theorems. We will use $D_{\text{ham}}(A, B)$ to denote the Hamming distance between A and
 072 B , and $D_{\text{edit}}(A, B)$ to denote the edit distance between A and B . We also say a data structure is
 073 ϵ -DP if it provides ϵ -DP outputs against any sequence of queries, of arbitrary length.

074 **Theorem 1.1.** *Let $A_1, \dots, A_m \in \{0, 1\}^n$ be a database, $k \in [n]$ and $\epsilon > 0, \beta \in (0, 1)$, then there
 075 exists a randomized algorithm with the following guarantees:*

- 077 • The data structure is ϵ -DP;
- 078
- 079 • It preprocesses A_1, \dots, A_m in $\tilde{O}(mn)$ time¹;
- 080
- 081 • It consumes $\tilde{O}(mk)$ space;
- 082
- 083 • Given any query $B \in \{0, 1\}^n$ such that $\max_{i \in [m]} D_{\text{ham}}(A_i, B) \leq k$, it outputs m estimates
 084 z_1, \dots, z_m with $|z_i - D_{\text{ham}}(A_i, B)| = \tilde{O}(k/e^{\epsilon/\log k})$ for all $i \in [m]$ in $\tilde{O}(mk + n)$
 085 time, and it succeeds with probability at least $1 - \beta$.

086 **Theorem 1.2.** *Let $A_1, \dots, A_m \in \{0, 1\}^n$ be a database, $k \in [n]$ and $\epsilon > 0, \beta \in (0, 1)$, then there
 087 exists a randomized algorithm with the following guarantees:*

- 088 • The data structure is ϵ -DP;
- 089
- 090 • It preprocesses A_1, \dots, A_m in $\tilde{O}(mn)$ time;
- 091
- 092 • It consumes $\tilde{O}(mn)$ space;
- 093
- 094 • Given any query $B \in \{0, 1\}^n$ such that $\max_{i \in [m]} D_{\text{edit}}(A_i, B) \leq k$, it outputs m estimates
 095 z_1, \dots, z_m with $|z_i - D_{\text{edit}}(A_i, B)| = \tilde{O}(k/e^{\epsilon/(\log k \log n)})$ for all $i \in [m]$ in $\tilde{O}(mk^2 + n)$
 096 time, and it succeeds with probability at least $1 - \beta$.

097 Before diving into the details, we would like to make several remarks regarding our data structure
 098 results. Note that instead of solving the general Hamming distance and edit distance problem, we
 099 impose the assumption that the query B has the property that for any $i \in [m]$, $\|A_i - B\| \leq k$.
 100 Such an assumption might seem restrictive at its first glance, but under the standard complexity
 101 assumption Strong Exponential Time Hypothesis (SETH) (Impagliazzo & Paturi, 2001; Impagliazzo
 102 et al., 2001), it is known that there is no $O(n^{2-o(1)})$ time algorithm for exact or even approximate
 103 edit distance (Belazzougui & Zhang, 2016; Chakraborty et al., 2016a;b; Naumovitz et al., 2017;
 104 Rubinstein et al., 2019; Rubinstein & Song, 2020; Goldenberg et al., 2020; Jin et al., 2021; Boroujeni
 105 et al., 2021; Kociumaka et al., 2021; Bhattacharya & Koucký, 2023; Koucký & Saks, 2024). It
 106 is therefore natural to impose assumptions that the query is “near” to the database in pursuit of
 107

¹Throughout the paper, we will use $\tilde{O}(\cdot)$ to suppress polylogarithmic factors in m, n, k and $1/\beta$.

faster algorithms (Ukkonen, 1985; Myers, 1986; Landau & Vishkin, 1988; Goldenberg et al., 2019; Kociumaka & Saha, 2020; Goldenberg et al., 2023). In fact, assuming SETH, $O(n + k^2)$ runtime for edit distance when $m = 1$ is optimal up to sub-polynomial factors (Goldenberg et al., 2023). Thus, in this paper, we consider the setting where $\max_{i \in [m]} \|A_i - B\| \leq k$ for both Hamming and edit distance and show how to craft private and efficient mechanisms for this class of distance problems.

Regarding privacy guarantees, one might consider the following simple augmentation to any fast data structure for Hamming distance: compute the distance estimate via the data structure, and add Laplace noise to it. Since changing one coordinate of the database would lead to the Hamming distance change by at most 1, Laplace mechanism would properly handle this case. However, our goal is to release a *differentially private data structure* that is robust against potentially infinitely many queries, and a simple output perturbation won't be sufficient as an adversary could simply query with the same B , average them to reduce the variance and obtain a relatively accurate estimate of the de-noised output. To address this issue, we consider the *differentially private function release communication model* (Hall et al., 2013), where the curator releases an ϵ -DP description of a function $\hat{e}(\cdot)$ that is ϵ -DP without seeing any query in advance. The client can then use $\hat{e}(\cdot)$ to compute $\hat{e}(B)$ for any query B . This strong guarantee ensures that the client could feed infinitely many queries to $\hat{e}(\cdot)$ without compromising the privacy of the database.

There are several prior works study related problems, such as (Kim et al., 2021) for n -grams extraction and (Steiner, 2024) for pattern matching. These algorithms do not release a private data structure, thus the quality of queries would degrade as more and more queries are served. (Bernardini et al., 2025) provides a private data structure for substring and document counting, which is an orthogonal direction to our setting.

Organization. Section 2 presents preliminaries. Section 3 introduces our new DP Hamming distance data structure and its theoretical analysis. Section 4 introduces our new DP edit distance data structure and its theoretical analysis. Lastly, Section 5 provides concluding remarks. Section A discusses related work. We discuss limitations of our work in Section B and future directions in C. In Section D we provide the proofs for Hamming distance data structure and E for edit distance.

2 PRELIMINARY

Let E be an event, we use $\mathbf{1}[E]$ to denote the indicator variable if E is true. Given two length- n bit strings A and B , we use $D_{\text{ham}}(A, B)$ to denote $\sum_{i=1}^n \mathbf{1}[A_i \neq B_i]$. We use $D_{\text{edit}}(A, B)$ to denote the edit distance between A and B , i.e., the minimum number of operations to transform A to B where the allowed operations are insertion, deletion and substitution. We use \oplus to denote the XOR operation. For any positive integer n , we use $[n]$ to denote the set $\{1, 2, \dots, n\}$. We use $\Pr[\cdot]$, $\mathbb{E}[\cdot]$ and $\text{Var}[\cdot]$ to denote probability, expectation and variance respectively. We say two database \mathcal{D}_1 and \mathcal{D}_2 are neighboring if there exists one $i \in [n]$ such that $\mathcal{D}_1(A_i)$ and $\mathcal{D}_2(A_i)$ differ by one bit.

2.1 CONCENTRATION BOUNDS

We will mainly use two concentration inequalities in this paper.

Lemma 2.1 (Chebyshev's Inequality). *Let X be a random variable with $0 < \text{Var}[X] < \infty$. For any real number $t > 0$,*

$$\Pr[|X - \mathbb{E}[X]| > t] \leq \frac{\text{Var}[X]}{t^2}.$$

Lemma 2.2 (Hoeffding's Inequality). *Let X_1, \dots, X_n with $a_i \leq X_i \leq b_i$ almost surely. Let $S_n = \sum_{i=1}^n X_i$, then for any real number $t > 0$,*

$$\Pr[|S_n - \mathbb{E}[S_n]| > t] \leq 2 \exp\left(-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right).$$

2.2 DIFFERENTIAL PRIVACY

Differential privacy (DP) is the key privacy measure we will be trying to craft our algorithm to possess. In this paper, we will solely focus on pure DP (ϵ -DP).

Definition 2.3 (ϵ -Differential Privacy). *We say an algorithm \mathcal{A} is ϵ -differentially private (ϵ -DP) if for any two neighboring databases \mathcal{D}_1 and \mathcal{D}_2 and any subsets of possible outputs S , we have*

$$\Pr[\mathcal{A}(\mathcal{D}_1) \in S] \leq e^\epsilon \cdot \Pr[\mathcal{A}(\mathcal{D}_2) \in S],$$

where the probability is taken over the randomness of \mathcal{A} .

162 Since we will be designing data structures, we will work with the *function release communication*
 163 *model* (Hall et al., 2013) where the goal is to release a function that is ϵ -DP against any sequence of
 164 queries of arbitrary length.

165 **Definition 2.4** (ϵ -DP Data Structure). *We say a data structure \mathcal{A} is ϵ -DP, if \mathcal{A} is ϵ -DP against any*
 166 *sequence of queries of arbitrary length. In other words, the curator will release an ϵ -DP description*
 167 *of a function $\hat{e}(\cdot)$ without seeing any query in advance.*

169 Finally, we will be utilizing the *post-processing* property of ϵ -DP.

170 **Lemma 2.5** (Post-Processing). *Let \mathcal{A} be ϵ -DP, then for any deterministic or randomized function g*
 171 *that only depends on the output of \mathcal{A} , $g \circ \mathcal{A}$ is also ϵ -DP.*

173 3 DIFFERENTIALLY PRIVATE HAMMING DISTANCE DATA STRUCTURE

174 To start off, we introduce our data structure for differentially private Hamming distance. In particular,
 175 we will adapt a data structure due to (Porat & Lipsky, 2007): this data structure computes a
 176 sketch of length $\tilde{O}(k)$ bit string to both the database and query, then with high probability, one could
 177 retrieve the Hamming distance from these sketches. Since the resulting sketch is also a bit string, a
 178 natural idea is to inject Laplace noise on each coordinate of the sketch. Since for two neighboring
 179 databases, only one coordinate would change, we could add Laplace noise of scale $1/\epsilon$ to achieve
 180 ϵ -DP. However, this approach has a critical issue: one could show that with high probability, the
 181 magnitude of each noise is roughly $O(\epsilon^{-1} \log k)$, aggregating the k coordinates of the sketch, this
 182 leads to a total error of $O(\epsilon^{-1} k \log k)$. To decrease this error to $O(1)$, one would have to choose
 183 $\epsilon = k \log k$, which is too large for most applications.

184 Instead of Laplace noise, we present a novel scheme that flips each bit of the sketch with certain
 185 probability. Our main contribution is to show that this simple scheme only produces an error of
 186 $O(e^{-\epsilon} \log k)$. To appreciate the significance, by denoting $t := (\log k)/\epsilon$, we see that the Laplace
 187 mechanism has an error of $O(t^{-1} k)$, whereas our error is only $O(e^{-1/t} k)$, which is exponentially
 188 small! In what follows, we will describe a data structure when the database is only one string A and
 189 with constant success probability, and we will discuss how to extend it to m bit strings, and how to
 190 boost the success probability to $1 - \beta$ for any $\beta > 0$. We summarize the main result below.

191 **Theorem 3.1.** *Given a string A of length n . There exists an ϵ -DP data structure DPHAMMINGDIS-
 192 TANCE (Algorithm 1), with the following operations*

- 193 • $\text{INIT}(A \in \{0, 1\}^n)$: *It takes a string A as input. This procedure takes $O(n \log k + k \log^3 k)$*
 194 *time.*
- 196 • $\text{QUERY}(B \in \{0, 1\}^n)$: *for any B with $z := D_{\text{ham}}(A, B) \leq k$, $\text{QUERY}(B)$ outputs a value*
 197 *\tilde{z} such that $|\tilde{z} - z| = \tilde{O}(k/e^{\epsilon} \log k)$ with probability 0.99, and the result is ϵ -DP. This*
 198 *procedure takes $O(n \log k + k \log^3 k)$ time.*

200 To achieve the results above, we set parameters $M_1 = O(\log k)$, $M_2 = O(k)$, $M_3 = O(\log^2 k)$ in
 201 Algorithm 1. We divide the proof of Theorem 3.1 into the following subsections:

202 3.1 TIME COMPLEXITY

203 Note that both the initializing and query run ENCODE (Algorithm 1) exactly once, we show that the
 204 running time of ENCODE is $O(n \log k)$.

206 **Lemma 3.2.** *Given $M_1 = O(\log k)$, the running time of ENCODE (Algorithm 1) is $O(n \log k)$.*

207 *Proof.* In ENCODE, for each character in the input string, the algorithm iterates M_1 times. Therefore
 208 the total time complexity is $O(n \cdot M_1) = O(n \log k)$. \square

209 3.2 PRIVACY GUARANTEE

210 Next we prove our data structure is ϵ -DP.

212 **Lemma 3.3.** *Let A and A' be two strings that differ on only one position. Let $\mathcal{A}(A)$ and $\mathcal{A}(A')$ be*
 213 *the output of INIT (Algorithm 1) given A and A' . For any output S , we have:*

$$214 \Pr[\mathcal{A}(A) = S] \leq e^\epsilon \cdot \Pr[\mathcal{A}(A') = S].$$

215 We defer the proof to Appendix D.

216 **Algorithm 1** Differential Private Hamming Distance Query

217 1: **data structure** DPHAMMINGDISTANCE ▷ Theorem 3.1
218 2: **members**
219 3: $M_1, M_2, M_3 \in \mathbb{N}_+$
220 4: $h(x) : [2n] \rightarrow [M_2]$ ▷ h and g are public random hash function
221 5: $g(x, i) : [2n] \times [M_1] \rightarrow [M_3]$
222 6: $S_{i,j,c} \in \{0, 1\}^{M_1 \times M_2 \times M_3}$ for all $i \in [M_1], j \in [M_2], c \in [M_3]$ ▷ S represents the sketch
223 7: **end members**
224 8:
225 9: **procedure** ENCODE($A \in \{0, 1\}^n, n$) ▷ Lemma 3.2
226 10: $S_{i,j,c}^* \leftarrow 0$ for all i, j, c
227 11: **for** $p \in [n]$ **do**
228 12: **for** $i \in [M_1]$ **do**
229 13: $j \leftarrow h(2(p-1) + A_p)$
230 14: $c \leftarrow g(2(p-1) + A_p, i)$
231 15: $S_{i,j,c}^* \leftarrow S_{i,j,c}^* \oplus 1$
232 16: **end for**
233 17: **end for**
234 18: **return** S^*
235 19: **end procedure**
236 20:
237 21: **procedure** INIT($A \in \{0, 1\}^n, n \in \mathbb{N}_+, k \in \mathbb{N}_+, \epsilon \in \mathbb{R}_+$) ▷ Lemma 3.3
238 22: $M_1 \leftarrow 10 \log k$
239 23: $M_2 \leftarrow 2k$
240 24: $M_3 \leftarrow 400 \log^2 k$
241 25: $S \leftarrow \text{ENCODE}(A, n)$
242 26: Flip each $S_{i,j,c}$ with independent probability $1/(1 + e^{\epsilon/(2M_1)})$
243 27: **end procedure**
244 28:
245 29: **procedure** QUERY($B \in \{0, 1\}^n$) ▷ Lemma 3.7
246 30: $S^B \leftarrow \text{ENCODE}(B, n)$
247 31: **return** $0.5 \cdot \sum_{j=1}^{M_2} \max_{i \in [M_1]} (\sum_{c=1}^{M_3} |S_{i,j,c} - S_{i,j,c}^B|)$
248 32: **end procedure**
249 33: **end data structure**

250 3.3 UTILITY GUARANTEE

251 The utility analysis is much more involved than privacy and runtime analysis. We defer the proofs
252 to the appendix, while stating key lemmas.

253 We first consider the distance between sketches of A and B without the random flipping process.
254 Let $E(A), E(B)$ be $\text{ENCODE}(A)$ and $\text{ENCODE}(B)$. We prove with probability 0.99, $D_{\text{ham}}(A, B) =$
255 $0.5 \cdot \sum_{j=1}^{M_2} \max_{i \in [M_1]} (\sum_{c=1}^{M_3} |E(A)_{i,j,c} - E(B)_{i,j,c}|)$. Before we present the error guarantee, we
256 will first introduce two technical lemmas. We let $T = \{p \subseteq [n] \mid A_p \neq B_p\}$ denote the set of “bad”
257 coordinates. Then we prove that for each coordinate in the sketch, the number of bad coordinates is
258 bounded.

259 **Lemma 3.4.** Define set $T := \{p \in [n] \mid A_p \neq B_p\}$. Define set $T_j := \{p \subseteq T \mid h(p) = j\}$. When
260 $M_2 = 2k$, with probability 0.99, for all $j \in [M_2]$, we have $|T_j| \leq 10 \log k$, i.e.,

$$261 \Pr[\forall j \in [M_2], |T_j| \leq 10 \log k] \geq 0.99.$$

263 The next lemma shows that with high probability, the second level hashing g will hash bad coordinates to distinct buckets.

265 **Lemma 3.5.** When $M_1 = 10 \log k, M_2 = 2k, M_3 = 400 \log^2 k$, with probability 0.98, for all
266 $j \in [M_2]$, there is at least one $i \in [M_1]$, such that all values in $\{g(2(p-1) + A_p, i) \mid p \in T_j\} \cup \{g(2(p-1) + B_p, i) \mid p \in T_j\}$ are distinct.

268 With these two lemmas in hand, we are in the position to prove the error bound before the random
269 bit flipping process.

270 **Lemma 3.6.** Let $E(A), E(B)$ be the output of $\text{ENCODE}(A)$ and $\text{ENCODE}(B)$. With probability
 271 0.98, $D_{\text{ham}}(A, B) = 0.5 \cdot \sum_{j=1}^{M_2} \max_{i \in [M_1]} (\sum_{c=1}^{M_3} |E(A)_{i,j,c} - E(B)_{i,j,c}|)$.
 272

273 Our final result provides utility guarantees for Algorithm 1.
 274

275 **Lemma 3.7.** Let z be $D_{\text{ham}}(A, B)$, \tilde{z} be the output of $\text{QUERY}(B)$ (Algorithm 1). With probability
 276 0.98, $|z - \tilde{z}| = O(k \log^3 k / e^{\epsilon / \log k})$.
 277

278 *Proof.* From Lemma 3.6, we know with probability 0.98, when $\epsilon \rightarrow \infty$ (i.e. without the random
 279 flip process), the output of $\text{QUERY}(B)$ (Algorithm 1) equals the exact hamming distance.
 280

281 We view the random flip process as random variables. Let random variables $R_{i,j,c}$ be 1 with proba-
 282 bility $1/(1 + e^{\epsilon / M_1})$, or 0 with probability $1 - 1/(1 + e^{\epsilon / M_1})$. So we have
 283

$$284 |\tilde{z} - z| = \sum_{j=1}^{M_2} \max_{i \in [M_1]} (\sum_{c=1}^{M_3} R_{i,j,c}) \leq \sum_{j=1}^{M_2} \sum_{i=1}^{M_1} (\sum_{c=1}^{M_3} R_{i,j,c}),$$

285 where the second step follows from $\max_i \leq \sum_i$ when all the summands are non-negative.
 286

287 Therefore, the expectation of $\tilde{z} - z$ is:
 288

$$289 \mathbb{E}[|\tilde{z} - z|] = M_1 M_2 M_3 \cdot \mathbb{E}[R_{i,j,c}] = k \log^3 k \cdot \frac{1}{(1 + e^{\epsilon / \log k})} = O(\frac{k \log^3 k}{e^{\epsilon / \log k}}),$$

290 where the last step follows from simple algebra. The variance of $\tilde{z} - z$ is:
 291

$$293 \text{Var}[|\tilde{z} - z|] = M_1 M_2 M_3 \cdot \text{Var}[R_{i,j,c}] = k \log^3 k \cdot \frac{1}{(1 + e^{\epsilon / \log k})} \cdot (1 - \frac{1}{(1 + e^{\epsilon / \log k})}).$$

294 Using Chebyshev's inequality (Lemma 2.1), we have
 295

$$297 \Pr[|\tilde{z} - z| \notin O(\frac{k \log^3 k}{e^{\epsilon / \log k}})] \leq 0.01.$$

298 Thus we complete the proof. □
 299

300 **Remark 3.8.** We will show how to generalize Theorem 3.1 to m bit strings, and how to boost the
 301 success probability to $1 - \beta$. To boost the success probability, we note that individual data structure
 302 succeeds with probability 0.99, we could take $\log(1/\beta)$ independent copies of the data structure,
 303 and query all of them. By a standard Chernoff bound argument, with probability at least $1 - \beta$, at
 304 least $3/4$ fraction of these data structures would output the correct answer, hence what we could do
 305 is to take the median of these answers. These operations blow up both INIT and QUERY by a factor
 306 of $\log(1/\beta)$ in its runtime. Generalizing for a database of m strings is relatively straightforward:
 307 we will run the INIT procedure to A_1, \dots, A_m , this would take $O(mn \log k + mk \log^3 k)$ time.
 308 For each query, note we only need to ENCODE the query once, and we can subsequently compute
 309 the Hamming distance from the sketch for m sketched database strings, therefore the total time for
 310 query is $O(n \log k + mk \log^3 k)$. It is important to note that as long as $k \log^3 k < n$, the query
 311 time is sublinear in the combined size of the query and its output. Finally, we could use the success
 312 probability boosting technique described before, that uses $\log(m/\beta)$ data structures to account for
 313 a union bound over the success of all distance estimates.
 314

4 DIFFERENTIALLY PRIVATE EDIT DISTANCE DATA STRUCTURE

315 Our algorithm for edit distance follows from the dynamic programming method introduced by
 316 (Ukkonen, 1985; Landau et al., 1998; Landau & Vishkin, 1988; Myers, 1986). We note that a
 317 key procedure in these algorithms is a subroutine to estimate *longest common prefix* (LCP) between
 318 two strings A and B and their substrings. We design an ϵ -DP data structure for LCP based on our
 319 ϵ -DP Hamming distance data structure. Due to space limitation, we defer the details of the DP-LCP
 320 data structure to Appendix E. In the following discussion, we will assume access to a DP-LCP data
 321 structure with the following guarantees:
 322

323 **Theorem 4.1.** Given a string A of length n . There exists an ϵ -DP data structure DPLCP (Algo-
 324 rithm 3 and Algorithm 4) supporting the following operations

- 324 • $\text{INIT}(A \in \{0, 1\}^n)$: It preprocesses an input string A . This procedure takes $O(n(\log k +$
325 $\log \log n))$ time.
- 326
- 327 • $\text{INITQUERY}(B \in \{0, 1\}^n)$: It preprocesses an input query string B . This procedure take
328 $O(n(\log k + \log \log n))$ time.
- 329
- 330 • $\text{QUERY}(i, j)$: Let w be the longest common prefix of $A[i : n]$ and $B[j : n]$ and \tilde{w} be
331 the output of $\text{QUERY}(i, j)$. With probability $1 - 1/(300k^2)$, we have: 1). $\tilde{w} \geq w$; 2).
332 $\mathbb{E}[D_{\text{ham}}(A[i : i + \tilde{w}], B[j : j + \tilde{w}])] = O((\log k + \log \log n)/e^{\epsilon/(\log k \log n)})$. This procedure
333 takes $O(\log^2 n(\log k + \log \log n))$ time.
- 334

We will be basing our edit distance data structure on the following result, which achieves the optimal dependence on n and k assuming SETH:

Lemma 4.2 ((Landau et al., 1998)). *Given two strings A and B of length n . If the edit distance between A and B is no more than k , there is an algorithm which computes the edit distance between A and B in time $O(k^2 + n)$.*

We start from a naive dynamic programming approach. Define $D(i, j)$ to be the edit distance between string $A[1 : i]$ and $B[1 : j]$. We could try to match $A[i]$ and $B[j]$ by inserting, deleting and substituting, which yields the following recurrence:

$$D(i, j) = \min \begin{cases} D(i-1, j) + 1, & \text{if } i > 0; \\ D(i-1, j-1) + 1, & \text{if } j > 0; \\ D(i-1, j-1) + \mathbf{1}[A[i] \neq B[j]], & \text{if } i, j > 0. \end{cases}$$

The edit distance between A and B is then captured by $D(n, n)$. When $k < n$, for all $D(i, j)$ such that $|i - j| > k$, because the length difference between $A[1 : i]$ and $B[1 : j]$ is greater than k , $D(i, j) > k$. Since the final answer $D(n, n) \leq k$, those positions with $|i - j| > k$ won't affect $D(n, n)$. Therefore, we only need to consider the set $\{D(i, j) : |i - j| \leq k\}$.

For $d \in [-k, k]$, $r \in [0, k]$, we define $F(r, d) = \max_i \{i : D(i, i + d) \leq r\}$ and let $\text{LCP}(i, j)$ denote the length of the longest common prefix of $A[i : n]$ and $B[j : n]$. The algorithm of (Landau et al., 1998) defines $\text{EXTEND}(r, d) := F(r, d) + \text{LCP}(F(r, d) + 1, F(r, d) + d + 1)$. Then we have the recurrence relation:

$$F(r, d) = \max \begin{cases} \text{EXTEND}(r-1, d) + 1, & \text{if } r-1 \geq 0; \\ \text{EXTEND}(r-1, d-1), & \text{if } d-1 \geq -k, r-1 \geq 0; \\ \text{EXTEND}(r-1, d+1) + 1, & \text{if } d+1, r+1 \leq k. \end{cases}$$

The edit distance between A and B equals $\min_r \{r : F(r, 0) = n\}$.

To implement LCP , (Landau et al., 1998) uses a suffix tree data structure with initialization time $O(n)$ and query time $O(1)$, thus the total time complexity is $O(k^2 + n)$. In place of their suffix tree data structure, we use our DP-LCP data structure (Theorem 4.1). This leads to Algorithm 2.

Theorem 4.3. *Given a string A of length n . There exists an ϵ -DP data structure DPEDITDISTANCE (Algorithm 2) supporting the following operations:*

- 363 • $\text{INIT}(A \in \{0, 1\}^n)$: It preprocesses an input string A . This procedure takes $O(n(\log k +$
364 $\log \log n))$ time.
- 365
- 366 • $\text{QUERY}(B \in \{0, 1\}^n)$: For any query string B with $w := D_{\text{edit}}(A, B) \leq k$, QUERY
367 outputs a value \tilde{w} such that $|w - \tilde{w}| \leq \tilde{O}(k/e^{\epsilon/(\log k \log n)})$ with probability 0.99. This
368 procedure takes $O(n(\log k + \log \log n) + k^2 \log^2 n(\log k + \log \log n)) = \tilde{O}(k^2 + n)$ time.
- 369

Again, we divide the proof into runtime, privacy and utility.

4.1 TIME COMPLEXITY

We prove the time complexity of INIT and QUERY respectively.

Lemma 4.4. *The running time of INIT (Algorithm 2) is $O(n(\log k + \log \log n))$.*

Proof. The INIT runs DPLCP.INIT . From Theorem 4.1, the init time is $O(n(\log k + \log \log n))$. \square

Lemma 4.5. *QUERY (Algorithm 2) runs in time $O((n + k^2 \log n)(\log k + \log \log n))$.*

Proof. The QUERY runs DPLCP.QUERYINIT once and DPLCP.QUERY k^2 times. From Theorem 4.1, the query time is $O(n(\log k + \log \log n) + k^2 \log^2 n(\log k + \log \log n))$. \square

4.2 PRIVACY GUARANTEE

Lemma 4.6. *The data structure DPEDITDISTANCE (Algorithm 2) is ϵ -DP.*

Proof. The data structure only stores a DPLCP(Algorithm 3, 4). From Theorem 4.1 and the post-processing property (Lemma 2.5), it is ϵ -DP. \square

4.3 UTILITY GUARANTEE

Before analyzing the error of the output of **QUERY** (Algorithm 2), we first introduce a lemma:

Lemma 4.7. Let A, B be two strings. Let $\text{LCP}(i, d)$ be the length of the true longest common prefix of $A[i : n]$ and $B[i + d : n]$. For $i_1 \leq i_2, d \in [-k, k]$, we have $i_1 + \text{LCP}(i_1, d) \leq i_2 + \text{LCP}(i_2, d)$.

Proof. Let $w_1 = \text{LCP}(i_1, d)$, $w_2 = \text{LCP}(i_2, d)$. Then for $j \in [i_1, i_1 + w_1 - 1]$, $A[j] = B[j + d]$. On the other side, w_2 is the length of the longest common prefix for $A[i_2 : n]$ and $B[i_2 + d : n]$. So $A[i_2 + w_2] \neq B[i_2 + w_2 + d]$. Therefore, $(i_2 + w_2) \notin [i_1, i_1 + w_1 - 1]$. Since $i_2 + w_2 \geq i_2 \geq i_1$, we have $i_2 + w_2 > i_1 + w_1$. \square

Lemma 4.8. Let \tilde{r} be the output of `QUERY` (Algorithm 2), r be the true edit distance $D_{\text{edit}}(A, B)$. With probability 0.99, we have $|r - \tilde{r}| = O(k(\log k + \log \log n)/(1 + e^{\epsilon/(\log k \log n)}))$.

Proof. We divide the proof into two parts. In part one, we prove that with probability 0.99, $\tilde{r} \leq r$. In part two, we prove that with probability 0.99, $r - \tilde{r} = O(k(\log k + \log \log n)/(1 + e^{\epsilon/(\log k \log n)}))$. In Theorem 4.1, with probability $1 - 1/(300k^2)$, DPLCP.QUERY satisfies two conditions. Our following discussion supposes all DPLCP.QUERY satisfies the two conditions. There are $3k^2$ LCP queries, by union bound, the probability is at least 0.99.

432 **Part I.** If no noise is added to guarantee differential privacy (using original LCP function instead
 433 of our DPLCP data structure), the dynamic programming method outputs the true edit distance. We
 434 define $F'(i, j)$ as the dynamic programming array F without privacy guarantee, $\text{EXTEND}'(i, j)$ be
 435 the result of $\text{EXTEND}(i, j)$ without privacy guarantee. Then we prove that for all $i \in [0, k], j \in$
 436 $[-k, k]$, $F(i, j) \geq F'(i, j)$ holds true.

437 We prove the statement above by math induction on i . For $i = 0$, $F(0, 0) = F'(0, 0) = 0$. Suppose
 438 for $i - 1$, $F(i - 1, j) \geq F'(i - 1, j)$, then for i ,

$$440 \quad F(i, j) = \max \begin{cases} \text{EXTEND}(i - 1, j) + 1, & \text{if } i - 1 \geq 0; \\ \text{EXTEND}(i - 1, j - 1), & \text{if } j - 1 \geq -k, i - 1 \geq 0; \\ \text{EXTEND}(i - 1, j + 1) + 1, & \text{if } j + 1, i + 1 \leq k, i - 1 \geq 0. \end{cases}$$

443 For $\text{EXTEND}(i - 1, j)$, we have

$$445 \quad \begin{aligned} \text{EXTEND}(i - 1, j) &= F(i - 1, j) + \text{DPLCP.QUERY}(F(i - 1, j), F(i - 1, j) + j) \\ 446 &\geq F(i - 1, j) + \text{LCP}(F(i - 1, j), F(i - 1, j) + j) \\ 447 &\geq F'(i - 1, j) + \text{LCP}(F'(i - 1, j), F'(i - 1, j) + j) \\ 448 &= \text{EXTEND}'(i - 1, j) \end{aligned}$$

450 The second step is because in QUERY (Theorem 4.1), $\tilde{w} \geq w$. The third step follows from
 451 $F(i - 1, j) \geq F'(i - 1, j)$ and Lemma 4.7. Thus, $F(i, j) = \max_{j_2 \in [j, j-1, j+1]} \{\text{EXTEND}(i, j_2)\} \geq$
 452 $\max_{j_2 \in [j, j-1, j+1]} \{\text{EXTEND}'(i, j_2)\} = F'(i, j)$. Since $\tilde{r} = \min\{\tilde{r} : F(\tilde{r}, 0) = n\}$, $r = \min\{r : F'(r, 0) = n\}$, we have $F(r, 0) \geq F'(r, 0) = n$. Therefore $\tilde{r} \leq r$.

454 **Part II.** Let $G(L, R, j) := D_{\text{edit}}(A[L : R], B[L + j, R + j])$. In this part, we prove that the edit
 455 distance $G(1, F(i, j), j) \leq i \cdot (1 + O((\log k + \log \log n)/(1 + e^{\epsilon/(\log k \log n)})))$ by induction on i .

457 For $i = 0$, $F(0, 0) = \text{LCP}(A, B)$. The statement holds true. Suppose for $i - 1$, $G(1, F(i - 1, j), j) \leq (i - 1) \cdot (1 + O((\log k + \log \log n)/(1 + e^{\epsilon/(\log k \log n)})))$, then we prove this holds for i .

460 Because $F(i, j) = \max_{j_2 \in [j, j-1, j+1]} \{\text{EXTEND}(i, j_2)\}$, there is some $j_2 \in \{j, j - 1, j + 1\}$
 461 such that $F(i, j) = F(i - 1, j_2) + \text{DPLCP.QUERY}(F(i - 1, j_2), F(i - 1, j_2) + j_2)$. Let
 462 $Q := \text{DPLCP.QUERY}(F(i - 1, j_2), F(i - 1, j_2) + j_2)$. Therefore

$$463 \quad \begin{aligned} G(1, F(i, j), j) &\leq G(1, F(i - 1, j_2) + Q, j_2) + 1 \\ 464 &\leq G(1, F(i - 1, j_2), j_2) + G(F(i - 1, j_2), F(i - 1, j_2) + Q, j_2) + 1 \\ 465 &\leq G(1, F(i - 1, j_2), j_2) + 1 + O((\log k + \log \log n)/(1 + e^{\epsilon/(\log k \log n)})) \\ 466 &\leq i \cdot (1 + O((\log k + \log \log n)/(1 + e^{\epsilon/(\log k \log n)}))) \end{aligned}$$

468 The third step follows from Theorem 4.1, and the fourth step follows from the induction hypothesis.
 469 Therefore, $r = G(1, F(\tilde{r}, 0), 0) \leq \tilde{r} \cdot (1 + O((\log k + \log \log n)/(1 + e^{\epsilon/(\log k \log n)})))$ and the proof
 470 is completed. \square

472 **Remark 4.9.** To the best of our knowledge, this is the first edit distance algorithm, based on noisy
 473 LCP implementations. In particular, we prove a structural result: if the LCP has query (additive)
 474 error δ , then we could implement an edit distance data structure with (additive) error $O(k\delta)$. Com-
 475 pared to standard relative error approximation, additive error approximation for edit distance is
 476 relatively less explored (see e.g., (Bringmann et al., 2022) for using additive approximation to solve
 477 the gap edit distance problem). We hope this structural result sheds light on additive error edit
 478 distance algorithms.

479 5 CONCLUSION

481 We study the differentially private Hamming distance and edit distance data structure problem in the
 482 function release communication model. This type of data structures are ϵ -DP against any sequence
 483 of queries of arbitrary length. For Hamming distance, our data structure has query time $\tilde{O}(mk + n)$
 484 and error $\tilde{O}(k/e^{\epsilon/(\log k)})$. For edit distance, our data structure has query time $\tilde{O}(mk^2 + n)$ and error
 485 $\tilde{O}(k/e^{\epsilon/(\log k \log n)})$. While the runtime of our data structures (especially edit distance) is nearly-
 optimal, it is interesting to design data structures with better utility in this model.

486 ETHIC STATEMENT

487
 488 This paper does not involve human subjects, personally identifiable data, or sensitive applications.
 489 We do not foresee direct ethical risks. We follow the ICLR Code of Ethics and affirm that all aspects
 490 of this research comply with the principles of fairness, transparency, and integrity.

491 REPRODUCIBILITY STATEMENT

492
 493 We ensure reproducibility of our theoretical results by including all formal assumptions, definitions,
 494 and complete proofs in the appendix. The main text states each theorem clearly and refers to the
 495 detailed proofs. No external data or software is required.

496 REFERENCES

497
 498 Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and
 499 Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC*
 500 *conference on computer and communications security*, pp. 308–318, 2016.

501
 502 Francesco Aldà and Benjamin I.P. Rubinstein. The bernstein mechanism: function release under
 503 differential privacy. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence*,
 504 AAAI’17, pp. 1705–1711. AAAI Press, 2017.

505
 506 Martin Aumüller, Anders Bourgeat, and Jana Schmurr. Differentially private sketches for jaccard
 507 similarity estimation. In *International Conference on Similarity Search and Applications*, pp.
 508 18–32. Springer, 2020.

509
 510 Arturs Backurs, Zinan Lin, Sepideh Mahabadi, Sandeep Silwal, and Jakub Tarnawski. Efficiently
 511 computing similarities to private datasets. In *The Twelfth International Conference on Learning*
 512 *Representations*, 2024.

513
 514 Eugene Bagdasaryan, Omid Poursaeed, and Vitaly Shmatikov. Differential privacy has disparate
 515 impact on model accuracy. *Advances in Neural Information Processing Systems (NeurIPS)*, 32:
 516 15479–15488, 2019.

517
 518 Amos Beimel, Haim Kaplan, Yishay Mansour, Kobbi Nissim, Thatchaphol Saranurak, and Uri
 519 Stemmer. Dynamic algorithms against an adaptive adversary: Generic constructions and lower
 520 bounds. In *Proceedings of the 54th Annual ACM SIGACT Symposium on Theory of Computing*,
 521 pp. 1671–1684, 2022.

522
 523 Djamal Belazzougui and Qin Zhang. Edit distance: Sketching, streaming, and document exchange.
 524 In *2016 IEEE 57th Annual Symposium on Foundations of Computer Science (FOCS)*, pp. 51–60.
 525 IEEE, 2016.

526
 527 Bonnie Berger, Michael S. Waterman, and Yun William Yu. Levenshtein distance, sequence com-
 528 parison and biological database search. *IEEE Transactions on Information Theory*, 2021.

529
 530 Giulia Bernardini, Philip Bille, Inge Li Gørtz, and Teresa Anna Steiner. Differentially private sub-
 531 string and document counting. *Proceedings of the ACM on Management of Data*, 3(2):1–27,
 532 2025.

533
 534 Sudatta Bhattacharya and Michal Koucký. Locally consistent decomposition of strings with appli-
 535 cations to edit distance sketching. In *Proceedings of the 55th Annual ACM Symposium on Theory*
 536 *of Computing*, pp. 219–232, 2023.

537
 538 Mahdi Boroujeni, Soheil Ehsani, Mohammad Ghodsi, MohammadTaghi HajiAghayi, and Saeed
 539 Seddighin. Approximating edit distance in truly subquadratic time: Quantum and mapreduce.
 540 *Journal of the ACM (JACM)*, 68(3):1–41, 2021.

541
 542 Karl Bringmann, Alejandro Cassis, Nick Fischer, and Vasileios Nakos. Improved Sublinear-
 543 Time Edit Distance for Preprocessed Strings. In Mikołaj Bojańczyk, Emanuela Merelli, and
 544 David P. Woodruff (eds.), *49th International Colloquium on Automata, Languages, and Program-
 545 ming (ICALP 2022)*, Leibniz International Proceedings in Informatics (LIPIcs), pp. 32:1–32:20,
 546 Dagstuhl, Germany, 2022. Schloss Dagstuhl – Leibniz-Zentrum für Informatik.

540 Diptarka Chakraborty, Elazar Goldenberg, and Michal Koucký. Streaming algorithms for embed-
 541 ding and computing edit distance in the low distance regime. In *Proceedings of the forty-eighth*
 542 *annual ACM symposium on Theory of Computing*, pp. 712–725, 2016a.

543

544 Diptarka Chakraborty, Elazar Goldenberg, and Michal Koucký. Streaming algorithms for computing
 545 edit distance without exploiting suffix trees. *arXiv preprint arXiv:1607.03718*, 2016b.

546

547 Moses Charikar. Similarity estimation techniques from rounding algorithms. In *Proceedings of the*
 548 *thiry-fourth annual ACM symposium on Theory of computing*, pp. 380–388, 2002.

549

550 Kamalika Chaudhuri and Claire Monteleoni. Privacy-preserving logistic regression. In *NIPS*, vol-
 551 ume 8, pp. 289–296. Citeseer, 2008.

552

553 Yeshwanth Cherapanamjeri, Sandeep Silwal, David Woodruff, Fred Zhang, Qiuyi Zhang, and Sam-
 554 son Zhou. Robust algorithms on adaptive inputs from bounded adversaries. In *The Eleventh*
 555 *International Conference on Learning Representations*, 2023.

556

557 Benjamin Coleman and Anshumali Shrivastava. A one-pass distributed and private sketch for kernel
 558 sums with applications to machine learning at scale. In *Proceedings of the 2021 ACM SIGSAC*
 559 *Conference on Computer and Communications Security*, CCS ’21, pp. 3252–3265, New York,
 560 NY, USA, 2021. Association for Computing Machinery.

561

562 Wenxin Ding, Gautam Kamath, Weina Wang, and Nihar B. Shah. Calibration with privacy in peer
 563 review. In *2022 IEEE International Symposium on Information Theory (ISIT)*, 2022.

564

565 Jasha Droppo and Alex Acero. Context dependent phonetic string edit distance for automatic speech
 566 recognition. In *2010 IEEE International Conference on Acoustics, Speech and Signal Processing*,
 567 2010.

568

569 Cynthia Dwork. Differential privacy. In *International Colloquium on Automata, Languages, and*
 570 *Programming (ICALP)*, pp. 1–12, 2006.

571

572 Cynthia Dwork, Krishnaram Kenthapadi, Frank McSherry, Ilya Mironov, and Moni Naor. Our data,
 573 ourselves: Privacy via distributed noise generation. In *Annual International Conference on the*
 574 *Theory and Applications of Cryptographic Techniques*, pp. 486–503. Springer, 2006.

575

576 Jonathan G. Fiscus, Jerome Ajot, Nicolas Radde, and Christophe Laprun. Multiple dimension Lev-
 577 enshtein edit distance calculations for evaluating automatic speech recognition systems during
 578 simultaneous speech. In Nicoletta Calzolari, Khalid Choukri, Aldo Gangemi, Bente Maegaard,
 579 Joseph Mariani, Jan Odijk, and Daniel Tapia (eds.), *Proceedings of the Fifth International Con-
 580 ference on Language Resources and Evaluation (LREC’06)*, Genoa, Italy, 2006. European Lan-
 581 guage Resources Association (ELRA).

582

583 Yeqi Gao, Zhao Song, and Xin Yang. Differentially private attention computation. *arXiv preprint*
 584 *arXiv:2305.04701*, 2023.

585

586 Omer Gold and Micha Sharir. Dynamic time warping and geometric edit distance: Breaking the
 587 quadratic barrier. *ACM Trans. Algorithms*, 2018.

588

589 Elazar Goldenberg, Robert Krauthgamer, and Barna Saha. Sublinear algorithms for gap edit dis-
 590 tance. In *2019 IEEE 60th Annual Symposium on Foundations of Computer Science (FOCS)*,
 591 2019.

592

593 Elazar Goldenberg, Aviad Rubinstein, and Barna Saha. Does preprocessing help in fast sequence
 594 comparisons? In *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Com-
 595 puting (STOC)*, pp. 657–670, 2020.

596

597 Elazar Goldenberg, Tomasz Kociumaka, Robert Krauthgamer, and Barna Saha. An Algorithmic
 598 Bridge Between Hamming and Levenshtein Distances. In Yael Tauman Kalai (ed.), *14th In-
 599 novations in Theoretical Computer Science Conference (ITCS 2023)*, Leibniz International Pro-
 600 ceedings in Informatics (LIPIcs), pp. 58:1–58:23, Dagstuhl, Germany, 2023. Schloss Dagstuhl –
 601 Leibniz-Zentrum für Informatik.

594 Rob Hall, Alessandro Rinaldo, and Larry Wasserman. Differential privacy for functions and func-
 595 tional data. *J. Mach. Learn. Res.*, 2013.

596

597 Richard W Hamming. Error detecting and error correcting codes. *The Bell System Technical Journal*,
 598 29(2):147–160, 1950.

599

600 Avinatan Hassidim, Haim Kaplan, Yishay Mansour, Yossi Matias, and Uri Stemmer. Adversarially
 601 robust streaming algorithms via differential privacy. *J. ACM*, 69(6), 2022.

602

603 Zhexue Huang. Extensions to the k-means algorithm for clustering large data sets with categorical
 604 values. *Data Mining and Knowledge Discovery*, 2(3):283–304, 1997.

605

606 Zhexue Huang and Mingkui Ng. A fuzzy k-modes algorithm for clustering categorical data. *IEEE
 607 Transactions on Fuzzy Systems*, 7(4):446–452, 1999.

608

609 Zhiyi Huang and Aaron Roth. Exploiting metric structure for efficient private query release. In
 610 *Proceedings of the Twenty-Fifth Annual ACM-SIAM Symposium on Discrete Algorithms*, SODA
 611 '14, pp. 523–534, 2014.

612

613 Russell Impagliazzo and Ramamohan Paturi. On the complexity of k-sat. *Journal of Computer and
 614 System Sciences*, 62(2):367–375, 2001.

615

616 Russell Impagliazzo, Ramamohan Paturi, and Francis Zane. Which problems have strongly expo-
 617 nential complexity? *Journal of Computer and System Sciences*, 63(4):512–530, 2001.

618

619 Piotr Indyk and Rajeev Motwani. Approximate nearest neighbors: Towards removing the curse of
 620 dimensionality. In *Proceedings of the Thirtieth Annual ACM Symposium on Theory of Computing*,
 621 pp. 604–613, 1998.

622

623 Bargav Jayaraman and David Evans. Evaluating differentially private machine learning in practice.
 624 In *28th USENIX Security Symposium (USENIX Security 19)*, pp. 1895–1912, 2019.

625

626 Ce Jin, Jelani Nelson, and Kewen Wu. An Improved Sketching Algorithm for Edit Distance. In *38th
 627 International Symposium on Theoretical Aspects of Computer Science (STACS)*, pp. 45:1–45:16,
 628 2021.

629

630 Kunho Kim, Sivakanth Gopi, Janardhan Kulkarni, and Sergey Yekhanin. Differentially private n-
 631 gram extraction. *Advances in neural information processing systems (NeurIPS)*, 34:5102–5111,
 632 2021.

633

634 Tomasz Kociumaka and Barna Saha. Sublinear-time algorithms for computing & embedding gap
 635 edit distance. In *2020 IEEE 61st Annual Symposium on Foundations of Computer Science
 (FOCS)*, 2020.

636

637 Tomasz Kociumaka, Ely Porat, and Tatiana Starikovskaya. Small-space and streaming pattern
 638 matching with k edits. In *2021 IEEE 62nd Annual Symposium on Foundations of Computer
 639 Science (FOCS)*, pp. 885–896. IEEE, 2021.

640

641 Michal Koucký and Michael E Saks. Almost linear size edit distance sketch. In *Proceedings of the
 642 56th Annual ACM Symposium on Theory of Computing*, pp. 956–967, 2024.

643

644 Gad M. Landau and Uzi Vishkin. Fast string matching with k differences. *Journal of Computer and
 645 System Sciences*, 37, 1988.

646

647 Gad M. Landau, Eugene Wimberly Myers, and Jeanette P. Schmidt. Incremental string comparison.
 648 *SIAM J. Comput.*, 27:557–582, 1998.

649

650 Vladimir I Levenshtein. Binary codes capable of correcting deletions, insertions and reversals.
 651 *Soviet Physics Doklady*, 10:707–710, 1966.

652

653 Erzhi Liu, Jerry Yao-Chieh Hu, Alex Reneau, Zhao Song, and Han Liu. Differentially private kernel
 654 density estimation. *arXiv preprint arXiv:2409.01688*, 2024.

648 Zelun Luo, Daniel J Wu, Ehsan Adeli, and Fei-Fei Li. Scalable differential privacy with sparse
 649 network finetuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
 650 *Recognition (CVPR)*, pp. 5059–5068, 2021.

651

652 Pierre-François Marteau. Time warp edit distance with stiffness adjustment for time series matching.
 653 *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2009.

654

655 Eugene W. Myers. An $O(ND)$ difference algorithm and its variations. *Algorithmica*, 1986.

656

657 Timothy Naumovitz, Michael Saks, and C Seshadhri. Accurate and nearly optimal sublinear approx-
 658 imations to ulam distance. In *Proceedings of the Twenty-Eighth Annual ACM-SIAM Symposium*
 659 *on Discrete Algorithms*, pp. 2012–2031. SIAM, 2017.

660

661 Gonzalo Navarro. A guided tour to approximate string matching. *ACM Comput. Surv.*, 2001.

662

663 Mohammad Norouzi, Ali Punjani, and David J Fleet. Fast search in hamming space with multi-index
 664 hashing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*
 665 (*CVPR*), pp. 3108–3115. IEEE, 2012.

666

667 Natalia Ponomareva, Hussein Hazimeh, Alex Kurakin, Zheng Xu, Carson Denison, H Brendan
 668 McMahan, Sergei Vassilvitskii, Steve Chien, and Abhradeep Guha Thakurta. How to dp-fy ml:
 669 A practical guide to machine learning with differential privacy. *Journal of Artificial Intelligence*
 670 *Research*, 77:1113–1201, 2023.

671

672 Ely Porat and Ohad Lipsky. Improved sketching of hamming distance with error correcting. In *Com-
 673 binatorial Pattern Matching*, pp. 173–182, Berlin, Heidelberg, 2007. Springer Berlin Heidelberg.

674

675 Aviad Rubinstein and Zhao Song. Reducing approximate longest common subsequence to approxi-
 676 mate edit distance. In *Proceedings of the Fourteenth Annual ACM-SIAM Symposium on Discrete*
 677 *Algorithms*, pp. 1591–1600. SIAM, 2020.

678

679 Aviad Rubinstein, Saeed Seddighin, Zhao Song, and Xiaorui Sun. Approximation algorithms for lcs
 680 and lis with truly improved running times. *FOCS*, 2019.

681

682 Grigori Sidorov, Helena Gómez-Adorno, Ilia Markov, David Pinto, and Nahun Loya. Computing
 683 text similarity using tree edit distance. In *2015 Annual Conference of the North American Fuzzy*
 684 *Information Processing Society (NAFIPS) held jointly with 2015 5th World Conference on Soft*
 685 *Computing (WConSC)*, 2015.

686

687 Zhao Song, Yitan Wang, Zheng Yu, and Lichen Zhang. Sketching for first order method: efficient
 688 algorithm for low-bandwidth channel and vulnerability. In *International Conference on Machine*
 689 *Learning (ICML)*, pp. 32365–32417. PMLR, 2023a.

690

691 Zhao Song, Xin Yang, Yuanyuan Yang, and Lichen Zhang. Sketching meets differential privacy:
 692 fast algorithm for dynamic kronecker projection maintenance. In *International Conference on*
 693 *Machine Learning (ICML)*, pp. 32418–32462. PMLR, 2023b.

694

695 Teresa Anna Steiner. Differentially private approximate pattern matching. In *Proceedings of the 15th*
 696 *Innovations in Theoretical Computer Science (ITCS 2024)*, volume 287 of *Leibniz International*
 697 *Proceedings in Informatics (LIPIcs)*, pp. 94:1–94:18, 2024. doi: 10.4230/LIPIcs.ITCS.2024.94.

698

699 Jiankai Sun, Xin Yang, Yuanshun Yao, Junyuan Xie, Di Wu, and Chong Wang. Dpauc: differen-
 700 tially private auc computation in federated learning. In *Proceedings of the Thirty-Seventh AAAI*
 701 *Conference on Artificial Intelligence and Thirty-Fifth Conference on Innovative Applications of*
 702 *Artificial Intelligence and Thirteenth Symposium on Educational Advances in Artificial Intelli-
 703 gence, AAAI’23/IAAI’23/EAAI’23*. AAAI Press, 2023.

704

705 Reihaneh Torkzadehmahani, Peter Kairouz, and Benedict Paten. Dp-cgan: Differentially private
 706 synthetic data and label generation. In *Proceedings of the IEEE/CVF Conference on Computer*
 707 *Vision and Pattern Recognition Workshops (CVPR Workshop)*, 2019.

708

709 Aleksei Triastcyn and Boi Faltings. Bayesian differential privacy for machine learning. In *Inter-
 710 national Conference on Machine Learning*, pp. 9583–9592. PMLR, 2020.

702 Esko Ukkonen. Finding approximate patterns in strings. *Journal of Algorithms*, 6(1):132–137, 1985.
 703

704 Tal Wagner, Yonatan Naamad, and Nina Mishra. Fast private kernel density estimation via locality
 705 sensitive quantization. In *Proceedings of the 40th International Conference on Machine Learning*,
 706 ICML’23. JMLR.org, 2023.

707 Xiao Shaun Wang, Yan Huang, Yongan Zhao, Haixu Tang, XiaoFeng Wang, and Diyue Bu. Effi-
 708 cient genome-wide, privacy-preserving similar patient query based on private edit distance. In
 709 *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security*,
 710 CCS ’15, pp. 492–503, New York, NY, USA, 2015. Association for Computing Machinery.

711 Ziteng Wang, Chi Jin, Kai Fan, Jiaqi Zhang, Junliang Huang, Yiqiao Zhong, and Liwei Wang.
 712 Differentially private data releasing for smooth queries. *Journal of Machine Learning Research*,
 713 2016.

714 Benjamin Weggenmann and Florian Kerschbaum. Syntf: Synthetic and differentially private term
 715 frequency vectors for privacy-preserving text mining. In *The 41st International ACM SIGIR
 716 Conference on Research & Development in Information Retrieval*, pp. 305–314, 2018.

717 Oliver Williams and Frank McSherry. Probabilistic inference and differential privacy. *Advances in
 718 Neural Information Processing Systems (NeurIPS)*, 23:2451–2459, 2010.

719 Ruihan Wu, Xin Yang, Yuanshun Yao, Jiankai Sun, Tianyi Liu, Kilian Q Weinberger, and Chong
 720 Wang. Differentially private multi-party data release for linear regression. In *The 38th Conference
 721 on Uncertainty in Artificial Intelligence*, 2022.

722 Xin Yang, Jiankai Sun, Yuanshun Yao, Junyuan Xie, and Chong Wang. Differentially private label
 723 protection in split learning. *arXiv preprint arXiv:2203.02073*, 2022.

724 Brian Young, Tom Faris, and Luigi Armogida. Levenshtein distance as a measure of accuracy and
 725 precision in forensic pcr-mps methods. *Forensic Science International: Genetics*, 55:102594,
 726 2021.

727 Da Yu, Saurabh Naik, Arturs Backurs, Sivakanth Gopi, Huseyin A. Inan, Gautam Kamath, Janard-
 728 han Kulkarni, Yin Tat Lee, Andre Manoel, Lukas Wutschitz, Sergey Yekhanin, and Huishuai
 729 Zhang. Differentially private fine-tuning of language models. In *The Tenth International Confer-
 730 ence on Learning Representations, ICLR 2022*, 2022.

731 Xiang Yue, Minxin Du, Tianhao Wang, Yaliang Li, Huan Sun, and Sherman S. M. Chow. Differential
 732 privacy for text analytics via natural text sanitization. In *Findings, ACL-IJCNLP 2021*, 2021.

733 Fuheng Zhao, Dan Qiao, Rachel Redberg, Divyakant Agrawal, Amr El Abbadi, and Yu-Xiang Wang.
 734 Differentially private linear sketches: Efficient implementations and applications. *Advances in
 735 Neural Information Processing Systems*, 35:12691–12704, 2022.

736 Yuqing Zhu, Xiang Yu, Manmohan Chandraker, and Yu-Xiang Wang. Private-knn: Practical dif-
 737 ferential privacy for computer vision. In *Proceedings of the IEEE/CVF Conference on Computer
 738 Vision and Pattern Recognition (CVPR)*, pp. 11854–11862, 2020.

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756 757 758 759 760 761 762 763 764 765 766 767 768 769 770 771 772 773 774 775 776 777 778 779 780 781 782 783 784 785 786 787 788 789 790 791 792 793 794 795 796 797 798 799 800 801 802 803 804 805 806 807 808 809 Appendix

LLM USAGE DISCLOSURE

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

A RELATED WORK

Differential Privacy. Differential privacy is a ubiquitous notion for protecting the privacy of database. (Dwork et al., 2006) first introduced this concept, which characterizes a class of algorithms such that when inputs are two neighboring datasets, with high probability the output distributions are similar. Differential privacy has a wide range of applications in general machine learning (Chaudhuri & Monteleoni, 2008; Williams & McSherry, 2010; Jayaraman & Evans, 2019; Triastcyn & Faltings, 2020; Ponomareva et al., 2023), training deep neural networks (Abadi et al., 2016; Bagdasaryan et al., 2019), computer vision (Zhu et al., 2020; Luo et al., 2021; Torkzadehmahani et al., 2019), natural language processing (Yue et al., 2021; Weggenmann & Kerschbaum, 2018), large language models (Gao et al., 2023; Yu et al., 2022), label protect (Yang et al., 2022), multiple data release (Wu et al., 2022), n -gram extraction (Kim et al., 2021), federated learning (Sun et al., 2023; Song et al., 2023a) and peer review (Ding et al., 2022). In recent years, differential privacy has been playing an important role for data structure design, both in making these data structures robust against adaptive adversary (Beimel et al., 2022; Hassidim et al., 2022; Song et al., 2023b; Cherapanamjeri et al., 2023) and in the function release communication model (Hall et al., 2013; Huang & Roth, 2014; Wang et al., 2016; Aldà & Rubinstein, 2017; Coleman & Shrivastava, 2021; Wagner et al., 2023; Backurs et al., 2024; Liu et al., 2024). (Steiner, 2024) studied the pattern matching problem under differential privacy setting. (Bernardini et al., 2025) provided a differentially private algorithm for substring and document counting. (Zhao et al., 2022) studied the how to build differential private algorithm for linear sketch. (Aumüller et al., 2020) proposed an algorithm for jaccard similarity estimation.

Hamming Distance and Edit Distance. Given bit strings A and B , many distance measurements have been proposed that capture various characteristics of bit strings. Hamming distance was first studied by Hamming (Hamming, 1950) in the context of error correction for codes. From an algorithmic perspective, Hamming distance is mostly studied in the context of approximate nearest-neighbor search and locality-sensitive hashing (Indyk & Motwani, 1998; Charikar, 2002). When it is known that the query B has the property $D_{\text{ham}}(A, B) \leq k$, (Porat & Lipsky, 2007) shows how to construct a sketch of size $\tilde{O}(k)$ in $\tilde{O}(n)$ time, and with high probability, these sketches preserve Hamming distance. Edit distance, proposed by Levenshtein (Levenshtein, 1966), is a more robust notion of distance between bit strings. It has applications in computational biology (Wang et al., 2015; Young et al., 2021; Berger et al., 2021), text similarity (Navarro, 2001; Sidorov et al., 2015) and speech recognition (Fiscus et al., 2006; Droppo & Acero, 2010). From a computational perspective, it is known that under the Strong Exponential Time Hypothesis (SETH), no algorithm can solve edit distance in $O(n^{2-o(1)})$ time, even its approximate variants (Belazzougui & Zhang, 2016; Chakraborty et al., 2016a;b; Naumovitz et al., 2017; Rubinstein et al., 2019; Rubinstein & Song, 2020; Goldenberg et al., 2020; Jin et al., 2021; Boroujeni et al., 2021; Kociumaka et al., 2021; Bhattacharya & Koucký, 2023; Koucký & Saks, 2024). Hence, various assumptions have been imposed to enable more efficient algorithm design. The most related assumption to us is that $D_{\text{edit}}(A, B) \leq k$, and in this regime various algorithms have been proposed (Ukkonen, 1985; Myers, 1986; Landau & Vishkin, 1988; Goldenberg et al., 2019; Kociumaka & Saha, 2020; Goldenberg et al., 2023). Under SETH, it has been shown that the optimal dependence on n and k is $O(n + k^2)$, up to sub-polynomial factors (Goldenberg et al., 2023).

B LIMITATIONS

Our proposed differentially private (DP) data structures focus specifically on Hamming and edit distances, which, while fundamental, do not generalize to other distance metrics such as Euclidean or Jaccard similarity. Extending our techniques to these metrics remains an open challenge. Furthermore, the efficiency of our approach relies on the assumption that the maximum distance between the query and any string in the database is bounded by k , which may not hold for datasets with

highly diverse bit strings. Although our data structures support sublinear query time for moderate k , the query time grows with k , potentially limiting scalability for large k in high-dimensional datasets. Additionally, while the privacy guarantees are strong (ϵ -DP for arbitrary query sequences), the constants involved in our error bounds may be non-negligible in practical scenarios, especially for smaller privacy budgets (ϵ).

C FUTURE WORK

We plan to explore extending our framework to broader metrics, improving scalability for large k , and tightening the constants in our error bounds for enhanced practical utility.

We currently lack a formal lower bound proof that precisely matches our upper bound result. We suspect this gap arises from a technical issue in the existing analysis of the algorithm. Any improvement in the error would likely require entirely new ideas. The relatively large additive error stems from our use of sketching: each sketch requires at least $\tilde{\Omega}(k)$ bits, and our method adds independent noise to each bit, resulting in an additive error of $\tilde{O}(k/e^\epsilon)$. Independent noise appears unavoidable, since the number of possible query strings grows exponentially with k and each sketch dimension behaves independently. For this reason, we do not see a path to improving beyond the current additive error of $\tilde{O}(k/e^\epsilon)$.

D PROOFS FOR HAMMING DISTANCE DATA STRUCTURE

In this section, we include all proof details in Section 3.

D.1 PROOF OF LEMMA 3.3

Proof of Lemma 3.3. Let $E(A), E(A')$ be $\text{ENCODE}(A)$ and $\text{ENCODE}(A')$. Let $\#(E(A) = S)$ be the number of the same bits between $E(A)$ and S , $\#(E(A) \neq S)$ be the number of the different bits between $E(A)$ and S . Then the probability that the random flip process transforms $E(A)$ into S is:

$$\begin{aligned} \Pr[\mathcal{A}(A) = S] &= \left(\frac{1}{1 + e^{\epsilon/(2M_1)}} \right)^{\#(E(A) \neq S)} \left(\frac{e^{\epsilon/(2M_1)}}{1 + e^{\epsilon/(2M_1)}} \right)^{\#(E(A) = S)} \\ &= \frac{(e^{\epsilon/(2M_1)})^{\#(E(A) = S)}}{(1 + e^{\epsilon/(2M_1)})^n} \end{aligned}$$

Since for each position, ENCODE changes at most M_1 bits, and A and A' only have one different position. Therefore there are at most $2M_1$ different bits between $E(A)$ and $E(A')$. So we have

$$\begin{aligned} \frac{\Pr[\mathcal{A}(A) = S]}{\Pr[\mathcal{A}(A') = S]} &\leq (e^{\epsilon/(2M_1)})^{|\#(E(A) = S) - \#(E(A') = S)|} \\ &\leq (e^{\epsilon/(2M_1)})^{2M_1} \\ &= e^\epsilon \end{aligned}$$

Thus we complete the proof. \square

D.2 PROOF OF LEMMA 3.4

Proof of Lemma 3.4. h is a hash function randomly drawn from all functions $[2n] \rightarrow [M_2]$. For certain j , $h(p) = j$ for all p are independent random variables, each of them equals 1 with probability $1/M_2$, or 0 with probability $1 - 1/M_2$. So we have

$$\begin{aligned} \Pr[|T_j| \geq 10 \log k] &= \Pr\left[\sum_{p \in T} [h(p) = j] \geq 10 \log k\right] \\ &= \sum_{d=10 \log k}^{|T|} \binom{|T|}{d} \left(\frac{1}{M_2}\right)^d \left(1 - \frac{1}{M_2}\right)^{|T|-d} \\ &\leq \sum_{d=10 \log k}^{|T|} \frac{|T|!}{d!(|T|-d)!} \left(\frac{1}{M_2}\right)^d \end{aligned}$$

$$\begin{aligned}
&\leq \sum_{d=10 \log k}^{|T|} \frac{|T|^d}{d!} \left(\frac{1}{M_2}\right)^d \\
&\leq \sum_{d=10 \log k}^{|T|} \frac{1}{d!} \left(\frac{1}{2}\right)^d \\
&\leq \frac{1}{(10 \log k)!} \sum_{d=10 \log k}^{|T|} \left(\frac{1}{2}\right)^d \\
&\leq \frac{1}{200k}
\end{aligned}$$

The fifth step follows from that fact that $|T| \leq k$, $M_2 = 2k$.

Therefore, by union bound over all $j \in [M_2]$, we can show

$$\begin{aligned}
\Pr[\forall j \in [M_2], |T_j| < 10 \log k] &\geq 1 - 2k \cdot \left(\frac{1}{200k}\right) \\
&= 0.99.
\end{aligned}$$

Thus, we complete the proof. \square

D.3 PROOF OF LEMMA 3.5

Proof of Lemma 3.5. g is a hash function randomly drawn from all functions $[2n] \times [M_1] \rightarrow [M_3]$. For every single $i \in [M_1]$, define event E_i as the event that the $2|T_j|$ values in $\{g(2(p-1) + A_p, i) \mid p \in T_j\} \cup \{g(2(p-1) + B_p, i) \mid p \in T_j\}$ are mapped into distinct positions.

$$\begin{aligned}
\Pr[E_i] &= \prod_{c=1}^{2|T_j|} \left(1 - \frac{c}{M_3}\right) \\
&\geq 1 - \sum_{c=1}^{2|T_j|} \frac{c}{M_3} \\
&= 1 - \frac{2|T_j|(|T_j|+1)}{M_3} \\
&> 1 - \frac{2(10 \log^2 k)}{400 \log^k} \\
&= 0.5
\end{aligned}$$

The fourth step follows from Lemma 3.4. It holds true with probability 0.99.

For different $i \in [M_1]$, E_i are independent. Therefore, the probability that all E_i are false is $(0.5)^{M_1} < 1/(1000k)$. By union bound, the probability that for every $j \in [M_2]$ there exists at least one i such that E_i is true is at least

$$1 - M_2 \cdot 0.5^{M_1} \geq 1 - M_2/(1000k) \geq 0.98. \quad \square$$

D.4 PROOF OF LEMMA 3.6

Proof of Lemma 3.6. From Lemma 3.5, for all j , there is at least one i , such that the set $\{g(2(p-1) + A_p, i) \mid p \in T_j\} \cup \{g(2(p-1) + B_p, i) \mid p \in T_j\}$ contains $2|T_j|$ distinct values. Therefore, for that i , $E(A)_{i,j,1 \sim M_3}$ and $E(B)_{i,j,1 \sim M_3}$ have exactly $2|T_j|$ different bits. For the rest of i , the different bits of $E(A)_{i,j,1 \sim M_3}$ and $E(B)_{i,j,1 \sim M_3}$ is no more than $2|T_j|$. So we have $0.5 \cdot \sum_{j=1}^{M_2} \max_{i \in [M_1]} (\sum_{c=1}^{M_3} |E(A)_{i,j,c} - E(B)_{i,j,c}|) = 0.5 \cdot \sum_{j=1}^{M_2} 2|T_j| = |T| = D_{\text{ham}}(A, B)$. \square

E DIFFERENTIAL PRIVATE LONGEST COMMON PREFIX

We design an efficient, ϵ -DP longest common prefix (LCP) data structure in this section. Specifically, for two positions i and j in A and B respectively, we need to calculate the maximum l , so that $A[i : i+l] = B[j : j+l]$. For this problem, we build a differentially private data structure (Algorithm 3 and Algorithm 4). The main contribution is a novel utility analysis that accounts for the error incurred by differentially private bit flipping.

918 **Algorithm 3** Differential Private Longest Common Prefix, Part 1

919 1: **data structure** DPLCP ▷ Theorem 4.1

920 2: **members**

921 3: $T_{i,j}^A, T_{i,j}^B$ for all $i \in [\log n], j \in [2^i]$

922 4: $\triangleright T_{i,j}$ represents the hamming sketch (Algorithm 1) of the interval $[i \cdot n/2^j, (i+1) \cdot n/2^j]$

923 5: **end members**

924 6:

925 7: **procedure** BUILDTREE($A \in \{0,1\}^n, n \in \mathbb{N}_+, k \in \mathbb{N}_+, \epsilon \in \mathbb{R}$) ▷ Lemma E.3

926 8: $M_1 \leftarrow \log k + \log \log n + 10, M_2 \leftarrow 1, M_3 \leftarrow 10, \epsilon' \leftarrow \epsilon / \log n$

927 9: **for** i from 0 to $\log n$ **do**

928 10: **for** j from 0 to $2^i - 1$ **do**

929 11: $T_{i,j}^* \leftarrow \text{DPHAMMINGDISTANCE.INIT}(A[j \cdot n/2^i : (j+1) \cdot n/2^i], M_1, M_2, M_3, \epsilon')$ ▷ Algorithm 1

930 12:

931 13: **end for**

932 14: **end for**

933 15: **return** T^*

934 16: **end procedure**

935 17:

936 18: **procedure** INIT($A \in \{0,1\}^n, n \in \mathbb{N}_+, k \in \mathbb{N}_+, \epsilon \in \mathbb{R}$) ▷ Lemma E.1

937 19: $T^A \leftarrow \text{BUILDTREE}(A, n, k, \epsilon)$

938 20: **end procedure**

939 21:

940 22: **procedure** QUERYINIT($B \in \{0,1\}^n, n \in \mathbb{N}_+, k \in \mathbb{N}_+$) ▷ Lemma E.1

941 23: $T^B \leftarrow \text{BUILDTREE}(B, n, k, 0)$

942 24: **end procedure**

943 25:

944 26: **procedure** INTERVALSKETCH($T, p_l \in [n], p_r \in [n]$)

945 27: Divide the interval $[p_l, p_r]$ into $O(\log n)$ intervals. Each of them is stored on a node of the tree T .

946 28: Retrieve the Hamming distance sketches of these nodes as S_1, S_2, \dots, S_t .

947 29: Initialize a new sketch $S \leftarrow 0$ with the same size of the sketches above.

948 30: **for** every position w in the sketch S **do**

949 31: $S[w] \leftarrow S_1[w] \oplus S_2[w] \oplus S_3[w] \oplus \dots \oplus S_t[w]$

950 32: **end for**

951 33: **return** S

952 34: **end procedure**

953 35:

954 36: **procedure** SKETCHHAMMINGDISTANCE($S^A, S^B \in \mathbb{R}^{M_1 \times M_2 \times M_3}$) ▷ Lemma E.4 and E.5

955 37: Let M_1, M_2, M_3 be the size of dimensions of the sketches S^A and S^B .

956 38: **return** $0.5 \cdot \sum_{j=1}^{M_2} \max_{i \in [M_1]} (\sum_{c=1}^{M_3} |S_{i,j,c}^A - S_{i,j,c}^B|)$

957 39: **end procedure**

958 40: **end data structure**

E.1 TIME COMPLEXITY

We prove the running time of the three operations above.

Lemma E.1. *The running time of INIT and INITQUERY (Algorithm 3) are $O(n \log n(\log k + \log \log n))$*

Proof. From Lemma 3.2, the running time of building node $T_{i,j}$ is $O((n/2^i)M_1)$. Therefore the total building time of all nodes is

$$\sum_{i=0}^{\log n} \sum_{j=0}^{2^i-1} (n/2^i) M_1 = \sum_{i=0}^{\log n} 2^i \cdot (n/2^i) M_1 = O(n \log n (\log k + \log \log n)).$$

Thus, we complete the proof. \square

Lemma E.2. *The running time of QUERY (Algorithm 4) is $O(\log^2 n(\log k + \log \log n))$*

972 **Algorithm 4** Differential Private Longest Common Prefix, Part 2

```

973 1: data structure DPCLP
974 2: procedure QUERY( $i \in [n], j \in [n]$ ) ▷ Theorem 4.1
975 3:    $L \leftarrow 0, R \leftarrow n$  ▷ Lemma E.2 and E.6
976 4:   while  $L \neq R$  do
977 5:      $mid \leftarrow \lceil \frac{L+R}{2} \rceil$ 
978 6:      $S^A \leftarrow \text{INTERVALSKETCH}(T^A, i, i + mid)$  ▷ Algorithm 3
979 7:      $S^B \leftarrow \text{INTERVALSKETCH}(T^B, j, j + mid)$  ▷ Algorithm 3
980 8:     threshold  $\leftarrow 1.5M_1M_3/(1 + e^{\epsilon}/(\log k \log n))$ 
981 9:     if  $\text{SKETCHHAMMINGDISTANCE}(S^A, S^B) \leq \text{threshold}$  then ▷ Algorithm 3
982 10:       $L \leftarrow mid$ 
983 11:    else
984 12:       $R \leftarrow mid - 1$ 
985 13:    end if
986 14:  end while
987 15:  return  $L$ 
988 16: end procedure
989 17: end data structure

```

990
991 *Proof.* In QUERY (Algorithm 4), we use binary search. There are totally $\log n$ checks. In each
992 check, we need to divide the interval into $\log n$ intervals and merge their sketches of size $M_1M_2M_3$.
993 So the time complexity is $O(\log^2 n(\log k + \log \log n))$. □
994

995 E.2 PRIVACY GUARANTEE

996 **Lemma E.3.** *The data structure DPLCP (Algorithm 3 and Algorithm 4) is ϵ -DP.*

997
998 *Proof.* On each node, we build a hamming distance data structure DPHAMMINGDISTANCE that is
999 $(\epsilon/\log n)$ -DP. For two strings A and A' that differ on only one bit, since every position is in at most
1000 $\log n$ nodes on the tree, for any output S , the probability
1001

$$\frac{\Pr[\text{BUILDTREE}(A) = S]}{\Pr[\text{BUILDTREE}(A') = S]} \leq (e^{\epsilon/\log n})^{\log n} = e^\epsilon$$

1002 Thus we complete the proof. □
1003

1004 E.3 UTILITY GUARANTEE

1005 Before analyzing the error of the query, we first bound the error of SKETCHHAMMINGDISTANCE
1006 (Algorithm 3).
1007

1008 **Lemma E.4.** *We select $M_1 = \log k + \log \log n + 10, M_2 = 1, M_3 = 10$ for DPHAMMINGDIS-
1009 TANCE data structure in BUILDTREE(Algorithm 3). Let z be the true hamming distance of the
1010 two strings $A[i : i + mid]$ and $B[i : i + mid]$. Let \tilde{z} be the output of SKETCHHAMMINGDIS-
1011 TANCE(Algorithm 3). When $\epsilon = +\infty$ (without the random flip process), then we have*

- 1012 • if $z = 0$, then with probability 1, $\tilde{z} = 0$.
- 1013 • if $z \neq 0$, then with probability $1 - 1/(300k^2 \log n)$, $\tilde{z} \neq 0$.

1014
1015 *Proof.* Our proof follows from the proof of Lemma 3.6. We prove the case of $z = 0$ and $z \neq 0$
1016 respectively.
1017

1018 When $z = 0$, it means the string $A[i : i + mid]$ and $B[i : i + mid]$ are identical. Therefore, the
1019 output of the hash function is also the same. Therefore, the output S^A and S^B from INTERVALSKETCH
1020 (Algorithm 3) are identical. Then $\tilde{z} = 0$.
1021

1022 When $z \neq 0$, define set $Q := \{p \in [mid] \mid A[i + p - 1] \neq B[j + p - 1]\}$ as the positions where
1023 string A and B are different. $|Q| = z$. Note that $M_1 = \log k + \log \log n + 10, M_2 = 1, M_3 = 10$,
 $S^A, S^B \in \{0, 1\}^{M_1 \times M_2 \times M_3}$. For every $i' \in [M_1]$, the probability that $S_{i'}^A$ and $S_{i'}^B$ are identical is

the probability that all $c \in [M_3]$ is mapped exactly even times from the position set Q . Formally, define event E as $[\forall j', |\{p \in Q \mid g(p) = j'\}| \bmod 2 = 0]$. Define another event E' as there is only one position mapped odd times from set $Q_{1 \sim z-1}$. Then the probability equals

$$\begin{aligned} \Pr[E] &= \Pr[E'] \cdot \Pr[E|E'] \\ &\leq \Pr[E|E'] \\ &= 1/M_3 \end{aligned}$$

The last step is because $\Pr[E|E']$ is the probability that $g(Q_z)$ equals the only position that mapped odd times. There are totally M_3 positions and the hash function g is uniformly distributed, so the probability is $1/M_3$.

For different $i' \in [M_1]$, the event E are independent. So the total probability that $z' \neq 0$ is the probability that for at least one i' , event E holds true. So the probability is $1 - (1/M_3)^{M_1} = 1 - (1/10)^{\log k + \log \log n + 10} > 1 - 1/(300k^2 \log n)$. \square

Lemma E.5. *Let $M_1 = \log k + \log \log n + 10$, $M_2 = 1$, $M_3 = 10$. Let z be the true hamming distance of the two strings $A[i : i + \text{mid}]$ and $B[i : i + \text{mid}]$. Let \tilde{z} be the output of `SKETCHHAMMINGDISTANCE`(Algorithm 3). With the random flip process with DP parameter ϵ , we have:*

- When $z = 0$, with probability $1 - 1/(300k^2 \log n)$, $\tilde{z} < (1 + o(1))M_1M_3/(1 + e^{\epsilon/(\log k \log n)})$.
- When $z > 3M_1M_3/(1 + e^{\epsilon/(\log k \log n)})$, with probability $1 - 1/(300k^2 \log n)$, $\tilde{z} > (2 - o(1))M_1M_3/(1 + e^{\epsilon/(\log k \log n)})$.

Proof. In the random flip process in `DPHAMMINGDISTANCE` (Algorithm 1,3), the privacy parameter $\epsilon' = \epsilon/\log n$. We flip each bit of the sketch with independent probability $1/(1 + e^{\epsilon/\log n \log k})$. Then we prove the case of $z = 0$ and $z > 3M_1M_3/(1 + e^{\epsilon/(\log k \log n)})$ respectively.

When $z = 0$, similar to the proof of Lemma 3.7, we view the flipping operation as random variables. Let random variable $R_{i,j,c}$ be 1 if the sketch $S_{i,j,c}^A$ is flipped, otherwise 0. From Lemma E.4, S^A and S^B are identical. Then we have

$$\begin{aligned} |z - \tilde{z}| &= \max_{i \in [M_1]} \sum_{c=1}^{M_3} R_{i,j,c} \\ &\leq \sum_{i=1}^{M_1} \sum_{c=1}^{M_3} R_{i,j,c} \end{aligned}$$

Since $R_{i,j,c}$ are independent Bernoulli random variables, using Hoeffding's inequality (Lemma 2.2), we have

$$\Pr\left[\left|\sum_{i=1}^{M_1 \times M_3} R_{i,j,c} - M_1M_3 \mathbb{E}[R_{i,j,c}]\right| > L\right] \leq e^{-2L^2/(M_1 \times M_3)}$$

When $L = M_1\sqrt{M_3}$,

$$\begin{aligned} \Pr\left[\left|\sum_{i=1}^{M_1 \times M_3} R_{i,j,c} - M_1M_3 \mathbb{E}[R_{i,j,c}]\right| > L\right] &\leq e^{-2M_1^2 M_3 / (M_1 M_3)} \\ &\leq e^{-2(\log k + \log \log n)} \\ &\leq 1/(300k^2 \log n) \end{aligned}$$

Thus we complete the $z = 0$ case.

When $z > 3M_1M_3/(1 + e^{\epsilon/(\log k \log n)})$, the proof is similar to $z = 0$. With probability $1 - 1/(300k^2 \log n)$, we have $|z - \tilde{z}| < (1 + o(1))M_1M_3/(1 + e^{\epsilon/(\log k \log n)})$. Thus $\tilde{z} > (2 - o(1))M_1M_3/(1 + e^{\epsilon/(\log k \log n)})$. \square

1080
 1081 **Lemma E.6.** Let \tilde{w} be the output of $\text{QUERY}(i, j)$ (Algorithm 4), w be the longest common prefix
 1082 of $A[i : n]$ and $B[j : n]$. With probability $1 - 1/(300k^2)$, we have: 1. $\tilde{w} \geq w$. 2. $D_{\text{ham}}(A[i : i + \tilde{w}], B[j : j + \tilde{w}]) \leq 3M_1M_3/(1 + e^{\epsilon/(\log k \log n)})$.
 1083

1084 *Proof.* In $\text{QUERY}(i, j)$ (Algorithm 4), we use a binary search to find the optimal w . In binary
 1085 search, there are totally $\log n$ calculations of $\text{SKETCHHAMMINGDISTANCE}$. Define threshold :=
 1086 $1.5M_1M_3/(1 + e^{\epsilon/(\log k \log n)})$. Define a return value of $\text{SKETCHHAMMINGDISTANCE}$ is good if:
 1087 1). when $z = 0$, $\tilde{z} < \text{threshold}$. 2). when $z > 2 \cdot \text{threshold}$, $\tilde{z} < \text{threshold}$. z and \tilde{z} are defined in
 1088 Lemma E.5.

1089 Therefore, by Lemma E.5, each $\text{SKETCHHAMMINGDISTANCE}$ is good with probability at least
 1090 $1 - 1/(k^2 \log n)$. There are $\log n$ $\text{SKETCHHAMMINGDISTANCE}$ in the binary search, by union
 1091 bound, the probability that all of them are good is at least $1 - 1/(300k^2)$.
 1092

1093 When all answers for $\text{SKETCHHAMMINGDISTANCE}$ are good, from the definition of binary search,
 1094 for any two positions L, R such that $D_{\text{ham}}(A[i : i+L], B[j, j+L]) = 0$, $D_{\text{ham}}(A[i : i+R], B[j, j+R]) \geq 2 \cdot \text{threshold}$, we have $L \leq \tilde{w} \leq R$. Next, we prove $w \leq \tilde{w}$ and $D_{\text{ham}}(A[i : i + \tilde{w}], B[j : j + \tilde{w}]) \leq 3M_1M_3/(1 + e^{\epsilon/(\log k \log n)})$ respectively.
 1095

1096 w is the true longest common prefix of $A[i : n]$ and $B[j : n]$, so we have $D_{\text{ham}}(A[i : i+L], B[j, j+L]) = 0$. Let $L = w$, we have $w = L \leq \tilde{w}$.
 1097

1098 Let R be the minimum value that $D_{\text{ham}}(A[i : i+R], B[j, j+R]) \geq 2 \cdot \text{threshold}$. Because
 1099 $D_{\text{ham}}(A[i : i+R], B[j, j+R])$ is monotone for R , and $\tilde{w} \leq R$, we have $D_{\text{ham}}(A[i : i + \tilde{w}], B[j : j + \tilde{w}]) \leq D_{\text{ham}}(A[i : i+R], B[j, j+R]) = 2 \cdot \text{threshold} = 3M_1M_3/(1 + e^{\epsilon/(\log k \log n)})$.
 1100

1101 Thus we complete the proof. □
 1102

1103
 1104
 1105
 1106
 1107
 1108
 1109
 1110
 1111
 1112
 1113
 1114
 1115
 1116
 1117
 1118
 1119
 1120
 1121
 1122
 1123
 1124
 1125
 1126
 1127
 1128
 1129
 1130
 1131
 1132
 1133