DIFFERENTIALLY PRIVATE ONE PERMUTATION HASHING

Anonymous authors

003 004

010

011

012

013

014

015

016

017

018

019

021

023

Paper under double-blind review

Abstract

Minwise hashing (MinHash) is a standard hashing algorithm for large-scale search and learning with the binary Jaccard similarity. One permutation hashing (OPH) is an effective and efficient alternative of MinHash which splits the data into *K* bins and generates hash values within each bin. In this paper, to protect the privacy of the output sketches, we combine differential privacy (DP) with OPH, and propose DP-OPH framework with three variants: DP-OPH-fix, DP-OPH-re and DP-OPH-rand, depending on the densification strategy to deal with empty bins in OPH. Detailed algorithm design and privacy and utility analysis are provided. The proposed DP-OPH methods significantly improves the DP minwise hashing (DP-MH) alternative in the literature. Experiments on similarity search confirm the effectiveness of our proposed algorithms. We also provide an extension to real-value data, named DP-BCWS, in the appendix.

1 INTRODUCTION

025 Let $u, v \in \{0, 1\}^D$ be two D-dimensional binary vectors. In this paper, we focus on the 026 hashing algorithms for the Jaccard similarity (a.k.a. the "resemblance") defined as J(u, v) =027 $\frac{\sum_{i=1}^{D} \mathbb{1}\{u_i = v_i = 1\}}{\sum_{i=1}^{D} \mathbb{1}\{u_i + v_i \ge 1\}}$. This is a widely used similarity measure in machine learning applications. u and 029 \overline{v} can also be viewed as two sets of items represented by the locations of non-zero entries. In industrial applications with massive data size, directly calculating the pairwise Jaccard similarity among the data points becomes too expensive. To accelerate large-scale search and learning, the celebrated 031 "minwise hashing" (MinHash) algorithm (Broder, 1997; Broder et al., 1997) has been a standard hashing technique for approximating the Jaccard similarity in massive binary datasets. It has seen 033 numerous applications such as near neighbor search, duplicate detection, malware detection, cluster-034 ing, large-scale learning, social networks, and computer vision (Indyk & Motwani, 1998; Charikar, 035 2002; Fetterly et al., 2003; Das et al., 2007; Buehrer & Chellapilla, 2008; Bendersky & Croft, 2009; Chierichetti et al., 2009; Pandey et al., 2009; Lee et al., 2010; Deng et al., 2012; Chum & Matas, 037 2012; Tamersov et al., 2014; Shrivastava & Li, 2014; Zhu et al., 2017; Nargesian et al., 2018; 038 Wang et al., 2019; Lemiesz, 2021; Feng & Deng, 2021; Li & Li, 2022). The output of MinHash is an integer. For large-scale applications, to store and use the hash values (or called sketches) 040 more conveniently and efficiently, Li & König (2010) proposed b-bit MinHash that only stores the 041 last b bits of the hashed integers, which is memory-efficient and convenient for similarity search and machine learning. Thus, it has been a popular coding strategy for the MinHash values and its 042 alternatives (Li et al., 2011; 2015; Shah & Meinshausen, 2017; Yu & Weber, 2022). 043

044 045

1.1 ONE PERMUTATION HASHING (OPH) FOR JACCARD SIMILARITY APPROXIMATION

To use MinHash in practice, we need to generate K hash values to achieve good utility. This requires applying K random permutations (or hash functions as approximations) per data point, yielding an O(Kf) complexity where f is the number of non-zero elements. One permutation hashing (OPH) (Li et al., 2012) provides a strategy to significantly reduce the complexity to O(f). The idea of OPH is: to generate K hashes, we split the data vector into K non-overlapping bins, and conduct MinHash within each bin. Yet, empty bins may arise which breaks the alignment of the hashes. To deal with empty bins, densification schemes (Shrivastava, 2017; Li et al., 2019) are proposed that fill the empty bins with some non-empty bin. It is shown that OPH with densification provides unbiased Jaccard estimator, and the estimation variance can often be smaller than that of MinHash. OPH has

been widely used as an improved method over MinHash for the Jaccard similarity (Dahlgaard et al., 055 2017; Zhao et al., 2020; Jia et al., 2021; Tseng et al., 2021; Jiang et al., 2022). 056

057 1.2 HASHING/SKETCHING AND DIFFERENTIAL PRIVACY

058 MinHash and OPH both belong to the broad family of probabilistic hashing/sketching methods 059 designed for various purposes and tasks. Examples of more sketching methods include the ran-060 dom projection (RP) based methods for cosine estimation (Charikar, 2002; Vempala, 2005), the 061 count-sketch (CS) for frequency estimation (Charikar et al., 2004), and the Flajolet-Martin (FM) 062 sketch (Flajolet & Martin, 1985) and HyperLogLog sketch (Flajolet et al., 2007) for cardinality es-063 timation, etc. Since the data sketches produce "summaries" of the original data, sketching/hashing 064 may also cause data privacy leakage. Therefore, protecting the privacy of the data sketches has 065 become an important topic which has gained growing research interests in recent years.

066 Differential privacy (DP) (Dwork et al., 2006b) has become a popular privacy definition with rigor-067 ous mathematical formulation, which has been widely applied to clustering, regression and classifi-068 cation, principle component analysis, matrix completion, optimization, deep learning (Blum et al., 069 2005; Chaudhuri & Monteleoni, 2008; Feldman et al., 2009; Gupta et al., 2010; Chaudhuri et al., 2011: Kasiviswanathan et al., 2013: Zhang et al., 2012: Abadi et al., 2016: Agarwal et al., 2018; 071 Ge et al., 2018; Wei et al., 2020; Dong et al., 2022), etc. Prior efforts have also been con-072 ducted on combining differential privacy with the aforementioned hashing algorithms, e.g., for RP (Blocki et al., 2012; Kenthapadi et al., 2013; Stausholm, 2021), count-sketch (Zhao et al., 2022), 073 and FM sketch (Smith et al., 2020; Dickens et al., 2022). Some works (e.g., Blocki et al. (2012); 074 Smith et al. (2020); Dickens et al. (2022)) assumed "internal randomness", i.e., the randomness of 075 the hash functions are kept private, and showed that many hashing methods themselves already pos-076 sess strong DP property under some data conditions. However, this setting is more restrictive in 077 practice as it requires that the hash keys or projection matrices cannot be accessed by any adver-078 sary. In another setup (e.g., Kenthapadi et al. (2013); Stausholm (2021); Zhao et al. (2022)), both 079 the randomness of the hash functions and the algorithm outputs are treated as public information, 080 and perturbation mechanisms are developed to make the algorithms differentially private.

081 **Contributions.** While prior works have proposed DP algorithms for some sketching methods men-082 tioned earlier, the differential privacy of OPH and MinHash for the Jaccard similarity has not been 083 well studied. In this paper, we mainly focus on the differential privacy of one permutation hash-084 ing (OPH), the state-of-the-art framework for hashing the Jaccard similarity. We consider the more 085 practical and general setup where the randomness of the algorithm is external/public.

We develop three variants under the DP-OPH framework, DP-OPH-fix, DP-OPH-re, and DP-OPH-087 rand, corresponding to fixed densification, re-randomized densification, and no densification for 088 OPH, respectively. We provide detailed algorithm design and privacy analysis for each variant, and 089 compare them with a DP MinHash (DP-MH) method. In our retrieval experiments, we show that the proposed DP-OPH method substantially improves DP-MH, and re-randomized densification is 091 superior over fixed densification in terms of differential privacy. DP-OPH-rand performs the best 092 when ϵ is small, while DP-OPH-re is the most performant in when larger ϵ is allowed. We also 093 extend our algorithms to real-value datasets and develop DP-BCWS algorithm in Appendix A.

094 095 096

098

099

100

101

102

103

BACKGROUND: MINHASH, b-BIT CODING, AND DIFFERENTIAL PRIVACY 2

| Algorithm I williwise hashing (willimiash) | Algorith | m 1 N | Minwise | hashing | (MinHash) |
|--|----------|-------|---------|---------|-----------|
|--|----------|-------|---------|---------|-----------|

Input: Binary vector $\boldsymbol{u} \in \{0, 1\}^D$; number of hash values K

Output: K MinHash values $h_1(\boldsymbol{u}), ..., h_K(\boldsymbol{u})$

1: Generate K independent permutations $\pi_1, ..., \pi_K$: $[D] \to [D]$ with seeds 1, ..., K respectively 2: **for** k = 1 to *K* **do**

3: $h_k(\boldsymbol{u}) \leftarrow \min_{i:u_i \neq 0} \pi_k(i)$

4: **end for** 104

105

Minwise hashing (MinHash). The MinHash method is summarized in Algorithm 1. We first 106 generate K independent permutations $\pi_1, ..., \pi_K : [D] \mapsto [D]$. Here, [D] denotes $\{1, ..., D\}$. For 107 each permutation, the hash value is the first non-zero location in the permuted vector, i.e., $h_k(u) =$

J

108 min_{i: $v_i \neq 0$} $\pi_k(i)$, $\forall k = 1, ..., K$. Analogously, for another data vector $v \in \{0, 1\}^D$, we also obtain K hash values, $h_k(v)$. The MinHash estimator of J(u, v) is the average over the hash collisions:

113

$$\hat{H}_{MH}(\boldsymbol{u},\boldsymbol{v}) = \frac{1}{K} \sum_{k=1}^{K} \mathbb{1}\{h_k(\boldsymbol{u}) = h_k(\boldsymbol{v})\},\tag{1}$$

where $\mathbb{1}\{\cdot\}$ is the indicator function. By standard probability calculation, we can show that $\mathbb{E}[\hat{J}_{MH}] = J$ and $Var[\hat{J}_{MH}] = \frac{J(1-J)}{K}$. In practice, K does not need to be very large to achieve good utility. For instance, usually 128 ~ 1024 hash values would be sufficient for search and learning problems (Indyk & Motwani, 1998; Li et al., 2011; Shrivastava & Li, 2014).

118 b-bit coding of the hash value. Li & König (2010) proposed "b-bit minwise hashing" as a conve-119 nient coding strategy for the integer hash value h(u) generated by MinHash (or by OPH which will 120 be introduced later). Basically, we only keep the last b-bits of each hash value. In our analysis, for 121 convenience, we assume that "taking the last b-bits" can be achieved by some "rehashing" trick to map the integer values onto $\{0, ..., 2^{b} - 1\}$ uniformly. There are at least three benefits of this coding 122 strategy: (i) storing only b bits saves the storage cost compared with storing the full 32 or 64 bit 123 integers; (ii) the last few bits are more convenient for the purpose of indexing, e.g., in approximate 124 nearest neighbor search (Indyk & Motwani, 1998); (iii) we can transform the last few bits into a 125 positional representation, allowing us to approximate the Jaccard similarity by inner product, which 126 is required by training large-scale linear models (Li et al., 2011). Given these advantages, in this 127 work, we will adopt this *b*-bit coding strategy in our private algorithm design. 128

Differential privacy (DP). We formally define differential privacy (DP) as follows.

Definition 2.1 (Differential privacy (Dwork et al., 2006b)). For a randomized algorithm $\mathcal{M} : \mathcal{U} \mapsto Range(\mathcal{M})$ and $\epsilon, \delta \ge 0$, if for any two neighboring datasets U and U', the following holds

$$Pr[\mathcal{M}(U) \in Z] \le e^{\epsilon} Pr[\mathcal{M}(U') \in Z] + \delta$$

133 134

132

for $\forall Z \subset Range(\mathcal{M})$, then algorithm \mathcal{M} is said to satisfy (ϵ, δ) -DP. If $\delta = 0$, \mathcal{M} is called ϵ -DP.

Intuitively, DP requires that the distributions of the outputs before and after a small change in the data are close, so that an adversary cannot detect the change based on the outputs. Smaller ϵ and δ implies stronger privacy. The parameter δ is usually interpreted as the "failure probability" allowed for the ϵ -DP guarantee to be violated.

Privacy statement and applications. We follow the standard attribute-level DP setup in aforemen-140 tioned related works on DP hashing/sketching: $u, u' \in \{0, 1\}^D$ are called neighboring if they differ 141 in one dimension. Treating the binary vectors as sets, with our proposed DP-OPH algorithms, an ad-142 versary cannot detect from the output sketches whether any item exists in the set or not, which holds 143 *independently for all the data vectors in the database.* DP-OPH can naturally be applied as a private 144 variant of OPH in cases where MinHash-type methods are found to be useful. As a concrete exam-145 ple application, the bioinformatics community releases sets of MinHashes for all known genomes 146 on a regular basis (Ondov et al., 2016; Brown & Irber, 2016), which are used for downstream ML 147 tasks like similarity search, classification, clustering, etc. (Berlin et al., 2015) In this type of data, 148 each data point corresponds to (a large set of) genes of a human, which contains the biological information of an individual which is highly sensitive and confidential. Our methods protect the 149 identification of any gene from the released sketches in the DP sense. 150

151 152

153

154

3 HASHING FOR JACCARD SIMILARITY WITH DIFFERENTIAL PRIVACY

In this section, we present DP-OPH algorithms based on privatizing the *b*-bit hash values from OPH and utility analysis, and demonstrate its advantage over a DP-MinHash alternative.

155 156 157

3.1 ONE PERMUTATION HASHING (OPH)

Algorithm 2 outlines the steps of OPH: we first use a permutation π (same for all data vectors) to randomly split the feature dimensions [D] into K bins $\mathcal{B}_1, ..., \mathcal{B}_K$ with equal length d = D/K(assuming integer division holds). Then, for each bin \mathcal{B}_k , we set the smallest permuted index of "1" as the k-th OPH hash value. If \mathcal{B}_k is empty (i.e., it does not contain any "1"), we record an "E" representing empty bin. Li et al. (2012) showed that we can construct statistically unbiased Jaccard

162 Algorithm 2 One Permutation Hashing (OPH) 163 **Input:** Binary vector $\boldsymbol{u} \in \{0, 1\}^D$; number of hash values K 164 **Output:** K OPH hash values $h_1(\boldsymbol{u}), ..., h_K(\boldsymbol{u})$ 165 1: Let d = D/K. Use a permutation $\pi : [D] \mapsto [D]$ with fixed seed to randomly split [D] into K 166 equal-size bins $\mathcal{B}_1, ..., \mathcal{B}_K$, with $\mathcal{B}_k = \{j \in [D] : (k-1)d + 1 \le \pi(j) \le kd\}$ 167 2: **for** k = 1 to *K* **do** 168 3: if Bin \mathcal{B}_k is non-empty then 169 4: $h_k(\boldsymbol{u}) \leftarrow \min_{i \in \mathcal{B}_k, u_i \neq 0} \pi(j)$ 170 5: else 171 $h_k(\boldsymbol{u}) \leftarrow E$ 6: end if 7: 172 8: end for 173 174 175 Algorithm 3 OPH-fix and OPH-re: OPH with fixed and re-randomized densification 176 **Input:** OPH hash values $h_1(u), ..., h_K(u)$ each in $[D] \cup \{E\}$; bins $\mathcal{B}_1, ..., \mathcal{B}_K$; d = D/K177 **Output:** K densified OPH hash values $h_1(\boldsymbol{u}), ..., h_K(\boldsymbol{u})$ 178 1: Let $NonEmptyBin = \{k \in [K] : h_k(u) \neq E\}$ 179 2: for k = 1 to K do 3: if $h_k(\boldsymbol{u}) = E$ then 181 4: Uniformly randomly select $k' \in NonEmptyBin$ 5: ▷ fixed densification $h_k(\boldsymbol{u}) \leftarrow h_{k'}(\boldsymbol{u})$ 6: Or 183 $MapToIndex = SortedIndex (\pi(\mathcal{B}_k)) + (k'-1)d$ 7: $\pi^{(k)}: \pi(\mathcal{B}_{k'}) \mapsto MapToIndex$ ▷ within-bin partial permutation 8: 185 $h_k(\boldsymbol{u}) \leftarrow \min_{j \in \mathcal{B}_{k'}, u_j \neq 0} \pi^{(k)}(\pi(j))$ 9: ▷ re-randomized densification 186 187 end if $10 \cdot$ 188 11: end for 189 190 191 estimators by ignoring the empty bins. However, this estimator is unstable when the data is relatively 192 sparse; moreover, since empty bins are different for every distinct data vector, the vanilla OPH hash 193 values do not form a metric space (i.e., the hash values of different data points are not aligned). 194 Densification for OPH. To tackle the issue caused by empty bins, a series of works has been con-195 ducted to densify the OPH. The general idea is to "borrow" the data/hash from non-empty bins, with 196 some careful design. In Algorithm 3, we present two recent representatives of OPH densification 197 methods: fixed densification (Shrivastava, 2017) and re-randomized densification (Li et al., 2019), 198 noted as OPH-fix and OPH-re, respectively. Given an OPH hash vector from Algorithm 2 (possibly 199 containing "E"s), we denote the set of non-empty bins $NonEmptyBin = \{k \in [K] : h_k(u) \neq E\}$. The densification procedure scans over k = 1, ..., K. For each k with $h_k(u) = E$, we do: 200 201 1. Uniformly randomly pick a bin $k' \in NonEmptyBin$ that is non-empty. 202 2. (a) OPH-fix: we directly copy the k'-th hash value: $h_k(u) \leftarrow h_{k'}(u)$. 203 (b) OPH-re: we apply an additional minwise hashing to bin $\mathcal{B}_{k'}$ using the "partial permu-204 tation" of \mathcal{B}_k to get the hash for $h_k(\boldsymbol{u})$. 205 206 Specifically, In Algorithm 2, for re-randomized densification, SortedIndex and MapToIndex are 207 used to define the within bin "partial permutation" $\pi^{(k)}$ of bin \mathcal{B}_k for re-randomizing the empty bins. 208 It is shown that for both variants, the Jaccard estimator of the same form as (1) is unbiased. Li et al. 209 (2019) showed that re-randomized densification always achieves smaller Jaccard estimation variance 210 than that of fixed densification, and the improvement is especially significant when the data is sparse. 211 Similar to b-bit MinHash, we can also keep the last b bits of the OPH hash values for convenient use. 212 213 3.2 DIFFERENTIAL PRIVATE ONE PERMUTATION HASHING (DP-OPH) 214 **DP-OPH with densification.** To privatize densified OPH, in Algorithm 4, we first take the last b 215 bits of the hash values. Since the output space is finite with cardinality 2^b , we apply the randomized

216 Algorithm 4 Differentially Private Densified One Permutation Hashing (DP-OPH-fix, DP-OPH-re) 217 **Input:** Densified OPH hash values $h_1(u), \dots, h_K(u)$; number of bits $b; \epsilon > 0, 0 < \delta < 1$ 218 f: lower bound on the number of non-zeros in each data vector 219 **Output:** b-bit DP-OPH values $h(\boldsymbol{u}) = [h_1(\boldsymbol{u}), ..., h_K(\boldsymbol{u})]$ 220 ▷ After which $h_k(u) \in \{0, ..., 2^b - 1\}$ 1: Take the last *b* bits of all hash values 221 2: Set $N = F_{fix}^{-1}(1-\delta; D, K, f)$ (for DP-OPH-fix) or $N = F_{re}^{-1}(1-\delta; D, K, f)$ (for DP-OPH-re), 222 and let $\epsilon' = \epsilon/N$ 223 3: for k = 1 to K do 224 $\tilde{h}_k(\boldsymbol{u}) = \begin{cases} h_k(\boldsymbol{u}), & \text{with prob.} \ \frac{e^{\epsilon'}}{e^{\epsilon'} + 2^b - 1}\\ i, & \text{with prob.} \ \frac{e^{\epsilon'}}{e^{\epsilon'} + 2^b - 1}, \text{for } i \in \{0, ..., 2^b - 1\}, \ i \neq h_k(\boldsymbol{u}) \end{cases}$ 225 4: 226 5: end for 227

228 229 230

231

232

233 234

235

236 237

259 260

261 262

264 265 266

267

response technique (Dwork & Roth, 2014; Wang et al., 2017) to flip the bits to achieve DP. After running Algorithm 3, suppose a densified OPH hash value $h_k(\boldsymbol{u}) = j, j \in 0, ..., 2^b - 1$. With some $\epsilon' > 0$ that will be specified later, we output $\tilde{h}_k(\boldsymbol{u}) = j$ with probability $\frac{e^{\epsilon'}}{e^{\epsilon'}+2^b-1}$, and $\tilde{h}_k(\boldsymbol{u}) = i$ for $i \neq j$ with probability $\frac{1}{e^{\epsilon'}+2^b-1}$. It is easy to verify that, for a neighboring data \boldsymbol{u}' , when $h_k(\boldsymbol{u}') = j$, for $\forall i \in 0, ..., 2^b - 1$, we have $\frac{P(\tilde{h}_k(\boldsymbol{u})=i)}{P(\tilde{h}_k(\boldsymbol{u}')=i)} = 1$; when $h_k(\boldsymbol{u}') \neq j$, we have $e^{-\epsilon'} \leq \frac{P(\tilde{h}_k(\boldsymbol{u})=i)}{P(\tilde{h}_k(\boldsymbol{u}')=i)} \leq e^{\epsilon'}$. Therefore, for a single hash value, this bit flipping satisfies ϵ' -DP.

It remains to determine ϵ' . Naively, since the perturbations (flipping) of the hash values are indepen-238 dent, by the composition property of DP (Dwork et al., 2006a), simply setting $\epsilon' = \epsilon/K$ for all K 239 MinHash values would achieve overall ϵ -DP (for the hashed vector). However, since K is usually 240 around hundreds, a very large ϵ value is required for this strategy to be useful. To this end, we can 241 trade a small δ in the DP definition for a significantly reduced ϵ . Note that, not all the K hashed bits 242 will change after we switch from u to its neighbor u'. Assume each data vector contains at least f 243 non-zeros, which is realistic since many data in practice have both high dimensionality D as well as 244 many non-zero elements. Intuitively, when the data is not too sparse, u and u' tends to be similar. 245 Therefore, the number of different hash values from Algorithm 3, $X = \sum_{k=1}^{K} \mathbb{1}\{h_k(\boldsymbol{u}) \neq h_k(\boldsymbol{u}')\}$ 246 can be upper bounded by some N with probability $1 - \delta$. In the proof, this allows us to set $\epsilon' = \epsilon/N$ 247 in the flipping probability and count δ as the failure probability in (ϵ, δ) -DP. 248

249 Next, we derive the distribution of X. Accordingly, in Algorithm 4, we set $N = F_{fix}^{-1}(1-\delta; D, f, K)$ 250 for DP-OPH-fix, $N = F_{re}^{-1}(1-\delta; D, f, K)$ for DP-OPH-re, where $F_{fix}(x) = P(X \le x)$ is the 251 cumulative mass function (CMF) of X with OPH-fix ((2) + (3)), and F_{re} is the CMF of X with 252 OPH-re ((2) + (4)), and F^{-1} is the inverse CMF. The proof can be found in Appendix **B**.

Lemma 3.1. Let $u, u' \in \{0, 1\}^D$ be neighbors. Denote $X = \sum_{k=1}^K \mathbb{1}\{h_k(u) \neq h_k(u')\}$ where the hashes are generated by Algorithm 3. Denote f = |u|, d = D/K. We have

$$P(X = x) = \sum_{j=\max(0,K-f)}^{K - \lceil f/d \rceil} \sum_{z=1}^{\min(f,d)} \Theta(x,j,z|K), \quad \text{for } x = 0, ..., K - \lceil f/d \rceil,$$
(2)

with $\Theta(x, j, z|K) = \tilde{P}(x|z, j)P\left(\tilde{f} = z|K - j\right)P(N_{emp} = j)$, where $P\left(\tilde{f} = z|K - j\right)$ is given in Lemma B.2, and $P(N_{emp} = j)$ is from Lemma B.1. Moreover,

For OPH-fix:
$$\tilde{P}(x|z,j) = \mathbb{1}\{x=0\} (1-P_{\neq}) + \mathbb{1}\{x>0\}P_{\neq} \cdot g_{bino}\left(x-1;\frac{1}{K-j},j\right),$$
 (3)

For OPH-re:
$$\tilde{P}(x|z,j) = (1 - P_{\neq}) \cdot g_{bino}\left(x; \frac{P_{\neq}}{K-j}, j\right) + P_{\neq} \cdot g_{bino}\left(x - 1; \frac{P_{\neq}}{K-j}, j\right),$$
 (4)

where $g_{bino}(x; p, n)$ is the CMF of Binomial(p, n), and $P_{\neq}(z, b) = \left(1 - \frac{1}{2^b}\right) \frac{1}{z}$.

The privacy guarantee of DP-OPH with densification is shown as below.

Theorem 3.2. Both DP-OPH-fix and DP-OPH-re in Algorithm 4 achieve (ϵ, δ) -DP.

270 Algorithm 5 Differentially Private One Permutation Hashing with Random Bits (DP-OPH-rand) 271 **Input:** OPH hash values $h_1(u), ..., h_K(u)$ from Algorithm 2; number of bits $b; \epsilon > 0$ 272 **Output:** DP-OPH-rand hash values $\tilde{h}(\boldsymbol{u}) = [\tilde{h}_1(\boldsymbol{u}), ..., \tilde{h}_K(\boldsymbol{u})]$ 273 1: Take the last *b* bits of all hash values ▷ After which $h_k(u) \in \{0, ..., 2^b - 1\}$ 274 2: for k = 1 to K do 275 3: if $h_k(\boldsymbol{u}) \neq E$ then 276 $\tilde{h}_k(\boldsymbol{u}) = \begin{cases} h_k(\boldsymbol{u}), & \text{with prob.} \ \frac{e^{\epsilon}}{e^{\epsilon}+2^b-1}\\ i, & \text{with prob.} \ \frac{1}{e^{\epsilon'}+2^b-1}, \ \text{for } i \in \{0, ..., 2^b-1\}, \ i \neq h_k(\boldsymbol{u}) \end{cases}$ 277 4: 278 5: 279 $\tilde{h}_k(u) = i$ with probability $\frac{1}{2^b}$, for $i = 0, ..., 2^b - 1 \, \triangleright$ Assign random bits to empty bin 6: 7: end if 281 8: end for 282

283 284 285

286

287

288

289

290

291 292 293

295 296

319

Algorithm 6 Differentially Private MinHash (DP-MH)Input: MinHash values $h_1(\boldsymbol{u}), ..., h_K(\boldsymbol{u})$; number of bits $b; \epsilon > 0, 0 < \delta < 1$
f: lower bound on the number of non-zeros in each data vectorOutput: DP-MH values $\tilde{h}(\boldsymbol{u}) = [\tilde{h}_1(\boldsymbol{u}), ..., \tilde{h}_K(\boldsymbol{u})]$ 1: Take the last b bits of all hash values
2: Set $N = F_{bino}^{-1}(1 - \delta; \frac{1}{f}, K)$, and $\epsilon' = \epsilon/N$ 3: for k = 1 to K do
4: $\tilde{h}_k(\boldsymbol{u}) = \begin{cases} h_k(\boldsymbol{u}), & \text{with prob. } \frac{e^{\epsilon'}}{e^{\epsilon'} + 2^b - 1}, \text{ for } i \in \{0, ..., 2^b - 1\}, i \neq h_k(\boldsymbol{u})$ 5: end for

DP-OPH without densification. From the practical perspective, we may also choose to privatize the OPH without densification (i.e., add DP to the output of Algorithm 2). The first step is to take the last *b* bits of every non-empty hash and get *K* hash values from $\{0, ..., 2^b - 1\} \cup \{E\}$. Then, for non-empty bins, we keep the hash value with probability $\frac{e^{\epsilon}}{e^{\epsilon}+2^{b}-1}$, and randomly flip it otherwise. For empty bins (i.e., $h_k(u) = E$), we simply assign a random value in $\{0, ..., 2^b - 1\}$ to $\tilde{h}_k(u)$. The formal procedure of this so-called DP-OPH-rand method is summarized in Algorithm 5.

Theorem 3.3. Algorithm 5 achieves ϵ -DP.

Compared with Algorithm 4, DP-OPH-rand achieves strict DP with smaller flipping probability (effectively, $N \equiv 1$ in Algorithm 4). This demonstrates the essential benefit of the binning operation in OPH, since the change in one data coordinate will only affect one hash value (if densification is not applied). As a result, the non-empty hash values are less perturbed in DP-OPH-rand than in DP-OPH-fix or DP-OPH-re. But this comes with an extra cost as we have to assign random bits to empty bins which do not provide any useful information, and this extra cost does not diminish as ϵ increases because the number of empty bins only depends on the data itself and K.

DP-MinHash. While we have presented our main DP-OPH algorithms, we also present a DP Min-Hash (DP-MH) method (Algorithm 6) as a baseline comparison. The mechanism of DP-MH is similar to that of densified DP-OPH. The difference between Algorithm 6 and Algorithm 4 is in the calculation of N. In Algorithm 6, we set $N = F_{bino}^{-1}(1-\delta; \frac{1}{f}, K)$ where $F_{bino}^{-1}(x; p, n)$ is the inverse cumulative mass function of Binomial(p, n) with n trials and success probability p.

Theorem 3.4. Algorithm 6 is (ϵ, δ) -DP.

The proof strategy is similar to Theorem 3.2 by noting that $X = \sum_{k=1}^{K} \mathbb{1}\{h_k(u) \neq h_k(u')\}$ for neighboring u and u' follows $Binomial(\frac{1}{f}, K)$. In a related work, Aumüller et al. (2020) also proposed to apply randomized response to MinHash. However, the authors incorrectly used a tail bound for the binomial distribution (see their Lemma 1) which is only valid for small deviation. In DP, δ is often very small (e.g., 10^{-6}), so the large deviation tail bound should be used which is looser than the one used therein¹. That said, in their paper, the perturbation is underestimated and their method does not satisfy DP. In our Algorithm 6, we fix it by using the exact probability mass function to compute the tail probability, avoiding any loss due to the concentration bounds.

3.3 COMPARISON OF DP-OPH AND DP-MH

We first analyze the mean of the Jaccard estimators and derive unbiased estimators of J. To simplify the formula, we assume that u and v have the same "privacy discount factor" N (which implies that u and v have similar sparsity). The results can be easily extended to the general case.

Theorem 3.5. For $u, v \in \{0, 1\}^D$, denote $f_u = |u|$, $f_v = |v|$, $a = |u \cap v|$. Suppose u and v have the same privacy discount factor N in Algorithm 4 or Algorithm 6. Then, $J = \frac{a}{f_u + f_v - a}$. Denote $p = \frac{\exp(\epsilon/N)}{\exp(\epsilon/N) + 2^b - 1}.$ For DP-OPH-fix, DP-OPH-re, and DP-MH, define $\hat{J} = \frac{1}{K} \sum_{k=1}^{K} \mathbb{1}\{\tilde{h}_k(\boldsymbol{u}) = \tilde{h}_k(\boldsymbol{v})\}.$ We have $\mathbb{E}[\hat{J}] = \frac{(2^b p + 1)^2}{2^b (2^b - 1)} J + \frac{1}{2^b}.$ Thus, an unbiased estimator is $\hat{J}_{unbias} = \frac{(2^b - 1)(2^b \hat{J} - 1)}{(2^b p - 1)^2}.$

The variances of the unbiased estimators defined in Theorem 3.5 are given as below.

Theorem 3.6. Define $J_B = J + (1 - J)\frac{1}{2^b}$, $\tilde{J} = \frac{a-1}{f_u + f_v - a - 1}$. Denote $c_1 = p^2 + \frac{(1-p)^2}{2^b - 1}$, and $c_2 = \frac{2p(1-p)}{2^b-1} + \frac{2^b-2}{(2^b-1)^2}(1-p)^2$. Define $\zeta(m) = \mathbb{E}[\frac{1}{\tilde{f}}|m]$ where the conditional distribution of \tilde{f} is given in Lemma B.2, and:

$$\begin{aligned} \tau_{11} &= JJ, \quad \tau_{10} = J - JJ, \quad \tau_{00} = 1 - 2J + JJ, \\ \tau_{11,f}(m) &= \frac{1}{m}J + \frac{m-1}{m}J\tilde{J}, \quad \tau_{10,f}(m) = \frac{m-1}{m}(J - J\tilde{J}), \quad \tau_{00,f}(m) = 1 - (2 - \frac{1}{m})J + \frac{m-1}{m}J\tilde{J}, \\ \tau_{11,r}(m) &= \frac{\zeta(m)}{m}J + \frac{m-\zeta(m)}{m}J\tilde{J}, \quad \tau_{10,r}(m) = \frac{m-\zeta(m)}{m}(J - J\tilde{J}), \\ \tau_{00,r}(m) &= 1 - (2 - \frac{\zeta(m)}{m})J + \frac{m-\zeta(m)}{m}J\tilde{J}. \end{aligned}$$

Further denote $P_{11} = \tau_{11} + \frac{1}{2^{b-1}}\tau_{10} + \frac{1}{2^{2b}}\tau_{00}, \quad P_{10} = J_B - P_{11}, \quad P_{00} = 1 - 2J_B + P_{11},$ and $(P_{11,f}, P_{10,f}, P_{00,f})$ and $(P_{11,r}, P_{10,r}, P_{00,r})$ analogously by replacing $(\tau_{11}, \tau_{10}, \tau_{00})$ with $(\tau_{11,f}, \tau_{10,f}, \tau_{00,f})$ and $(\tau_{11,r}, \tau_{10,r}, \tau_{00,r})$, respectively. We have for DP-MH:

$$Var[\hat{J}_{unbias,MH}] = \frac{1}{K} \left(\frac{(2^b p + 1)^2}{(2^b p - 1)^2} J + \frac{2^b - 1}{(2^b p - 1)^2} \right) \left(\frac{(2^b - 1)^2}{(2^b p - 1)^2} - \frac{(2^b p + 1)^2}{(2^b p - 1)^2} J \right).$$

For DP-OPH: Let $m = K - N_{emp}$ where N_{emp} is distributed as Lemma B.1.

$$Var[\hat{J}_{unbias,OPH}] = \frac{2^{2b}(2^b - 1)^2}{(2^b p - 1)^4} \left[\frac{1}{K} \left(\frac{(2^b p + 1)^2}{2^b (2^b - 1)} J + \frac{1}{2^b} \right) + \frac{1}{K^2} A - \left(\frac{(2^b p + 1)^2}{2^b (2^b - 1)} J + \frac{1}{2^b} \right)^2 \right],$$

$$A = \mathbb{E}_{m} [m(m-1)H_{N} + (K-m)(K+m-1)H_{E}]$$

with $H_N = c_1^2 P_{11} + 2c_1 c_2 P_{10} + c_2^2 P_{00}$. For DP-OPH-fix, $H_E = c_1^2 P_{11,f} + 2c_1 c_2 P_{10,f} + c_2^2 P_{00,f}$; for DP-OPH-re, $H_E = c_1^2 P_{11,r} + 2c_1 c_2 P_{10,r} + c_2^2 P_{00,r}$.

Comparison: Densified DP-OPH vs. DP-MH. We show that OPH is a better method than MinHash from the privacy perspective. Firstly, we compare N, the privacy discount factor, in DP-OPH-fix, DP-OPH-re, and DP-MH. Smaller N leads to smaller bit flipping probability which benefits the utility. In Figure 1, we plot N vs. f, for D = 1024, K = 64, and $\delta = 10^{-6}$. Similar comparison also holds for other D, K combinations. We observe that N in DP-OPH is typically smaller than that in DP-MH. Moreover, N for DP-OPH-re is consistently smaller than that of DP-OPH-fix. This illustrates that re-randomization in densification is an important step to achieve stronger privacy.

In Figure 2, we plot the empirical MSE of the unbiased estimators. The data is simulated with $f_u = f_v = f$, and a = f/2 (see notations in Theorem 3.5). The empirical MSE matches the variances in Theorem 3.6. DP-OPH-re has smaller variance than DP-OPH-fix and DP-MH.

¹For X following a Binomial distribution with mean μ , Aumüller et al. (2020) used the concentration in-equality $P(X \ge (1+\xi)\mu) \le \exp(-\frac{\xi^2\mu}{3})$, which only holds when $0 \le \xi \le 1$. For large deviations (large ξ), the valid Binomial tail bound should be $P(X \ge (1+\xi)\mu) \le \exp(-\frac{\xi^2\mu}{\xi+2})$.



Figure 1: Comparison of the privacy discount factor N for densified DP-OPH and DP-MH, against the number of non-zero elements in the data vector f. $D = 1024, K = 64, \delta = 10^{-6}$.



Figure 2: MSE comparison of the unbiased Jaccard estimators (Theorem 3.5). The dash curves are theoretical variances in Theorem 3.6. $D = 1024, K = 64, \delta = 10^{-6}$. b = 4.

4 EXPERIMENTS

405 We conduct similarity search on two datasets from genome and web where MinHash-type algorithms 406 are widely used: (1) the Leukemia gene expression dataset (https://sbcb.inf.ufrgs.br/cumida); (2) the 407 Webspam (Chang & Lin, 2011) dataset for spam detection. Both datasets are binarilized to 0/1. For Leukemia, we first standardize the features columns (to mean 0 and std 1), and then keep entries 408 larger than 1 to be 1 and zero out the others. For Webspam, the entries are non-negative and we 409 simply set the positive elements to 1. For Leukemia, we treat each data point as the query and other 410 points as the database for search. For Webspam, we use the training set as the database, and the test 411 set as queries. For each query point, we set the ground truth ("gold-standard") neighbors as the top 412 50 data points in the database with the highest Jaccard similarities to the query. 413

414 Setup. To search with DP-OPH and DP-MH, we generate the private hash values and compute the 415 collision estimator between the query and each data point. Then, we retrieve the data points with the 416 highest estimated Jaccard similarities to the query. For densified DP-OPH (Algorithm 4) and DP-417 MH (Algorithm 6), we ensure the lower bound f on the number of non-zero elements by filtering 418 the data points with at least f non-zeros. We use f = 1000, 500 for Leukemia and Webspam, 419 respectively, which cover 100% and 90% of the total data points.

Results. In Figure 3, we report the precision for Leukemia with b = 1, 2, 4 and $\epsilon \in [1, 30]$. The ϵ range is common in the literature of DP hashing, e.g., the [2.45, 33.5] reported in Zhao et al. (2022) which studied private count-sketch. The recall comparisons are similar. The results are averaged over all query points and over 5 runs. We observe that:

424 425 426

388

389

400

401 402 403

404

• DP-OPH-re outperforms DP-MH and DP-OPH-fix, at all ϵ levels. That is, DP-OPH-re is a uniformly more superior method than the existing DP-MH method for private hashing.

- DP-OPH-rand achieves good accuracy with small ϵ (e.g., $\epsilon < 5$), but stops improving with ϵ afterwards (due to the random bits for the empty bins), justifying the trade-off discussed in Section 3.2. When ϵ gets larger (e.g., $\epsilon = 5 \sim 15$), DP-OPH-re performs the best.
- 428 429 430





Figure 3: Precision@1 results on Leukemia gene expression dataset with b = 1, 2, 4. $\delta = 10^{-6}$. Dotted curves are for non-private OPH-re.



Figure 4: Precision@10 results on Webspam dataset with b = 2. $\delta = 10^{-6}$.

CONCLUSION 5

461

473 474

475

476 In this paper, we study differentially privatized one permutation hashing (DP-OPH) methods. We 477 develop three variants depending on the densification procedure of OPH, and provide privacy and 478 utility analyses of our algorithms. We show the significant advantages of our DP-OPH over the DP 479 MinHash alternative proposed in prior literature for hashing the Jaccard similarity at various privacy 480 levels. Experiments are conducted on retrieval tasks to justify the effectiveness of the proposed DP-OPH, and provide guidance on the appropriate choice of the DP-OPH variant in different sce-481 narios. In Appendix A, we also provide DP-BCWS which is based on bin-wise consistent weighted 482 samples (BCWS) (Li et al., 2019) for weighted Jaccard similarity (for non-negative data). Given 483 the efficiency and good performance, we expect DP-OPH to serve as a useful privatized alternative 484 in practical applications where MinHash-type methods are heavily used. In the appendix, we also 485 provide an extension of DP-OPH to real-value data called DP-BCWS.

486 REFERENCES 487

504

521

523

529

- 488 Martín Abadi, Andy Chu, Ian J. Goodfellow, H. Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep learning with differential privacy. In Proceedings of the 2016 ACM SIGSAC 489 Conference on Computer and Communications Security (CCS), pp. 308–318, Vienna, Austria, 490 2016. 491
- 492 Naman Agarwal, Ananda Theertha Suresh, Felix X. Yu, Sanjiv Kumar, and Brendan McMahan. 493 cpSGD: Communication-efficient and differentially-private distributed SGD. In Advances in Neu-494 ral Information Processing Systems (NeurIPS), pp. 7575–7586, Montréal, Canada, 2018. 495
- A. Asuncion and D.J. Newman. UCI machine learning repository, 2007. URL 496 http://www.ics.uci.edu/\$\sim\$mlearn/{MLR}epository.html. 497
- 498 Martin Aumüller, Anders Bourgeat, and Jana Schmurr. Differentially private sketches for jaccard 499 similarity estimation. CoRR, abs/2008.08134, 2020. 500
- 501 Michael Bendersky and W. Bruce Croft. Finding text reuse on the web. In Proceedings of the 502 Second International Conference on Web Search and Web Data Mining (WSDM), pp. 262–271, 503 Barcelona, Spain, 2009.
- Konstantin Berlin, Sergey Koren, Chen-Shan Chin, James P Drake, Jane M Landolin, and Adam M 505 Phillippy. Assembling large genomes with single-molecule sequencing and locality-sensitive 506 hashing. *Nature biotechnology*, 33(6):623–630, 2015. 507
- 508 Jeremiah Blocki, Avrim Blum, Anupam Datta, and Or Sheffet. The johnson-lindenstrauss transform 509 itself preserves differential privacy. In Proceedings of the 53rd Annual IEEE Symposium on 510 Foundations of Computer Science (FOCS), pp. 410-419, New Brunswick, NJ, 2012. 511
- 512 Avrim Blum, Cynthia Dwork, Frank McSherry, and Kobbi Nissim. Practical privacy: the SuLQ framework. In Proceedings of the Twenty-fourth ACM SIGACT-SIGMOD-SIGART Symposium 513 on Principles of Database Systems (PODS), pp. 128-138, Baltimore, MD, 2005. 514
- 515 Andrei Z Broder. On the resemblance and containment of documents. In Proceedings of the Com-516 pression and Complexity of Sequences (SEQUENCES), pp. 21-29, Salerno, Italy, 1997. 517
- 518 Andrei Z. Broder, Steven C. Glassman, Mark S. Manasse, and Geoffrey Zweig. Syntactic clustering 519 of the web. Comput. Networks, 29(8-13):1157-1166, 1997. 520
- C. Titus Brown and Luiz Irber. sourmash: a library for minhash sketching of DNA. J. Open Source Softw., 1(5):27, 2016. 522
- Gregory Buehrer and Kumar Chellapilla. A scalable pattern mining approach to web graph com-524 pression with communities. In Proceedings of the International Conference on Web Search and 525 Web Data Mining (WSDM), pp. 95–106, Stanford, CA, 2008. 526
- 527 Chih-Chung Chang and Chih-Jen Lin. Libsvm: a library for support vector machines. ACM Trans-528 actions on Intelligent Systems and Technology (TIST), 2(3):27, 2011.
- Moses Charikar, Kevin Chen, and Martin Farach-Colton. Finding frequent items in data streams. 530 Theor. Comput. Sci., 312(1):3-15, 2004. 531
- 532 Moses S Charikar. Similarity estimation techniques from rounding algorithms. In *Proceedings of the* 533 Thiry-Fourth Annual ACM Symposium on Theory of Computing (STOC), pp. 380–388, Montreal, 534 Canada, 2002. 535
- Kamalika Chaudhuri and Claire Monteleoni. Privacy-preserving logistic regression. In Advances in 537 Neural Information Processing Systems (NIPS), pp. 289–296, Vancouver, Canada, 2008.
- Kamalika Chaudhuri, Claire Monteleoni, and Anand D. Sarwate. Differentially private empirical 539 risk minimization. J. Mach. Learn. Res., 12:1069-1109, 2011.

540 Flavio Chierichetti, Ravi Kumar, Silvio Lattanzi, Michael Mitzenmacher, Alessandro Panconesi, and 541 Prabhakar Raghavan. On compressing social networks. In Proceedings of the 15th ACM SIGKDD 542 International Conference on Knowledge Discovery and Data Mining (KDD), pp. 219–228, Paris, 543 France, 2009. 544 Ondrej Chum and Jiri Matas. Fast computation of min-hash signatures for image collections. In 545 Proceedings of the 2012 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 546 pp. 3077-3084, Providence, RI, 2012. 547 Søren Dahlgaard, Mathias Bæk Tejs Knudsen, and Mikkel Thorup. Practical hash functions for 548 similarity estimation and dimensionality reduction. In Advances in Neural Information Processing 549 Systems (NIPS), pp. 6615–6625, Long Beach, CA, USA, 2017. 550 551 Abhinandan Das, Mayur Datar, Ashutosh Garg, and Shyamsundar Rajaram. Google news personal-552 ization: scalable online collaborative filtering. In Proceedings of the 16th International Confer-553 ence on World Wide Web (WWW), pp. 271-280, Banff, Alberta, Canada, 2007. 554 Fan Deng, Stefan Siersdorfer, and Sergej Zerr. Efficient jaccard-based diversity analysis of large 555 document collections. In Proceedings of the 21st ACM International Conference on Information 556 and Knowledge Management (CIKM), pp. 1402-1411, Maui, HI, 2012. 557 Charlie Dickens, Justin Thaler, and Daniel Ting. Order-invariant cardinality estimators are differen-559 tially private. In Advances in Neural Information Processing Systems (NeurIPS), New Orleans, LA, 2022. 560 561 Jinshuo Dong, Aaron Roth, and Weijie J Su. Gaussian differential privacy. Journal of the Royal 562 Statistical Society Series B: Statistical Methodology, 84(1):3–37, 2022. 563 Cynthia Dwork and Aaron Roth. The algorithmic foundations of differential privacy. Found. Trends 564 Theor. Comput. Sci., 9(3-4):211-407, 2014. 565 566 Cynthia Dwork, Krishnaram Kenthapadi, Frank McSherry, Ilya Mironov, and Moni Naor. Our data, 567 ourselves: Privacy via distributed noise generation. In Advances in Cryptology - EUROCRYPT 568 2006, 25th Annual International Conference on the Theory and Applications of Cryptographic 569 Techniques, St. Petersburg, Russia, May 28 - June 1, 2006, Proceedings, volume 4004 of Lecture 570 Notes in Computer Science, pp. 486–503. Springer, 2006a. 571 Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam D. Smith. Calibrating noise to sensitivity 572 in private data analysis. In Proceedings of the Third Theory of Cryptography Conference (TCC), 573 pp. 265-284, New York, NY, 2006b. 574 Dan Feldman, Amos Fiat, Haim Kaplan, and Kobbi Nissim. Private coresets. In Proceedings of 575 the 41st Annual ACM Symposium on Theory of Computing (STOC), pp. 361–370, Bethesda, MD, 576 2009. 577 578 Weiqi Feng and Dong Deng. Allign: Aligning all-pair near-duplicate passages in long texts. In 579 Proceedings of the International Conference on Management of Data (SIGMOD), pp. 541–553, 580 Virtual Event, China, 2021. 581 Dennis Fetterly, Mark Manasse, Marc Najork, and Janet L. Wiener. A large-scale study of the 582 evolution of web pages. In Proceedings of the Twelfth International World Wide Web Conference 583 (WWW), pp. 669–678, Budapest, Hungary, 2003. 584 585 Philippe Flajolet and G. Nigel Martin. Probabilistic counting algorithms for data base applications. J. Comput. Syst. Sci., 31(2):182-209, 1985. 586 Philippe Flajolet, Éric Fusy, Olivier Gandouet, and Frédéric Meunier. Hyperloglog: the analysis 588 of a near-optimal cardinality estimation algorithm. In Discrete Mathematics and Theoretical 589 Computer Science, pp. 137–156, 2007. 590 Jason Ge, Zhaoran Wang, Mengdi Wang, and Han Liu. Minimax-optimal privacy-preserving sparse PCA in distributed systems. In Proceedings of the International Conference on Artificial Intelli-592 gence and Statistics (AISTATS), pp. 1589–1598, Playa Blanca, Lanzarote, Canary Islands, Spain,

2018.

594 Anupam Gupta, Katrina Ligett, Frank McSherry, Aaron Roth, and Kunal Talwar. Differentially 595 private combinatorial optimization. In Proceedings of the Twenty-First Annual ACM-SIAM Sym-596 posium on Discrete Algorithms (SODA), pp. 1106–1125, Austin, TX, 2010. 597 Piotr Indyk and Rajeev Motwani. Approximate nearest neighbors: Towards removing the curse 598 of dimensionality. In Proceedings of the Thirtieth Annual ACM Symposium on the Theory of Computing (STOC), pp. 604-613, Dallas, TX, 1998. 600 601 Sergey Ioffe. Improved consistent sampling, weighted minhash and L1 sketching. In Proceedings of 602 the 10th IEEE International Conference on Data Mining (ICDM), pp. 246–255, Sydney, Australia, 603 2010. 604 Peng Jia, Pinghui Wang, Junzhou Zhao, Shuo Zhang, Yiyan Qi, Min Hu, Chao Deng, and Xiaohong 605 Guan. Bidirectionally densifying LSH sketches with empty bins. In Proceedings of the Interna-606 tional Conference on Management of Data (SIGMOD), pp. 830–842, Virtual Event, China, 2021. 607 608 Nan Jiang, Chen Luo, Vihan Lakshman, Yesh Dattatreya, and Yexiang Xue. Massive text normalization via an efficient randomized algorithm. In WWW '22: The ACM Web Conference 2022, 609 Virtual Event, Lyon, France, April 25 - 29, 2022, pp. 2946–2956. ACM, 2022. 610 611 Shiva Prasad Kasiviswanathan, Kobbi Nissim, Sofya Raskhodnikova, and Adam D. Smith. Analyz-612 ing graphs with node differential privacy. In Proceedings of the 10th Theory of Cryptography 613 Conference, TCC, volume 7785, pp. 457–476, Tokyo, Japan, 2013. 614 Krishnaram Kenthapadi, Aleksandra Korolova, Ilya Mironov, and Nina Mishra. Privacy via the 615 johnson-lindenstrauss transform. J. Priv. Confidentiality, 5(1), 2013. 616 617 David C. Lee, Qifa Ke, and Michael Isard. Partition min-hash for partial duplicate image discovery. 618 In Proceedings of the 11th European Conference on Computer Vision (ECCV), Part I, pp. 648– 619 662, Heraklion, Crete, Greece, 2010. 620 Jakub Lemiesz. On the algebra of data sketches. Proc. VLDB Endow., 14(9):1655–1667, 2021. 621 622 Jin Li, Sudipta Sengupta, Ran Kalach, Ronakkumar N Desai, Paul Adrian Oltean, and James Robert 623 Benton. Using index partitioning and reconciliation for data deduplication, August 18 2015. US 624 Patent 9,110,936. 625 Ping Li and Arnd Christian König. b-bit minwise hashing. In Proceedings of the 19th International 626 Conference on World Wide Web (WWW), pp. 671–680, Raleigh, NC, 2010. 627 628 Ping Li, Anshumali Shrivastava, Joshua L. Moore, and Arnd Christian König. Hashing algorithms 629 for large-scale learning. In Advances in Neural Information Processing Systems (NIPS), pp. 2672– 630 2680, Granada, Spain, 2011. 631 Ping Li, Art B Owen, and Cun-Hui Zhang. One permutation hashing. In Advances in Neural 632 Information Processing Systems (NIPS), pp. 3122–3130, Lake Tahoe, NV, 2012. 633 634 Ping Li, Xiaoyun Li, and Cun-Hui Zhang. Re-randomized densification for one permutation hash-635 ing and bin-wise consistent weighted sampling. In Advances in Neural Information Processing Systems (NeurIPS), pp. 15900–15910, Vancouver, Canada, 2019. 636 637 Ping Li, Xiaoyun Li, Gennady Samorodnitsky, and Weijie Zhao. Consistent sampling through ex-638 tremal process. In Proceedings of the Web Conference (WWW), pp. 1317–1327, Virtual Event / 639 Ljubljana, Slovenia, April 19-23, 2021, 2021. 640 Xiaoyun Li and Ping Li. C-MinHash: Improving minwise hashing with circulant permutation. In 641 Proceedings of the International Conference on Machine Learning (ICML), pp. 12857–12887, 642 Baltimore, MD, 2022. 643 644 Mark Manasse, Frank McSherry, and Kunal Talwar. Consistent weighted sampling. Technical 645 Report MSR-TR-2010-73, Microsoft Research, 2010. 646 Fatemeh Nargesian, Erkang Zhu, Ken Q. Pu, and Renée J. Miller. Table union search on open data. 647

Proc. VLDB Endow., 11(7):813-825, 2018.

657

| 648 | Brian D Ondoy, Todd J Treangen, Páll Melsted, Adam B Mallonee, Nicholas H Bergman, Sergey |
|-----|---|
| 649 | Koren, and Adam M Phillippy. Mash: fast genome and metagenome distance estimation using |
| 650 | minhash. Genome biology, 17(1):1–14, 2016. |
| 651 | |

- Sandeep Pandey, Andrei Broder, Flavio Chierichetti, Vanja Josifovski, Ravi Kumar, and Sergei Vassilvitskii. Nearest-neighbor caching for content-match applications. In *Proceedings of the 18th International Conference on World Wide Web (WWW)*, pp. 441–450, Madrid, Spain, 2009.
- Rajen Dinesh Shah and Nicolai Meinshausen. On b-bit min-wise hashing for large-scale regression
 and classification with sparse data. J. Mach. Learn. Res., 18:178:1–178:42, 2017.
- Anshumali Shrivastava. Optimal densification for fast and accurate minwise hashing. In *Proceed-ings of the 34th International Conference on Machine Learning (ICML)*, pp. 3154–3163, Sydney, Australia, 2017.
- Anshumali Shrivastava and Ping Li. In defense of minhash over simhash. In *Proceedings of the Seventeenth International Conference on Artificial Intelligence and Statistics (AISTATS)*, pp. 886–
 894, Reykjavik, Iceland, 2014.
- Adam D. Smith, Shuang Song, and Abhradeep Thakurta. The flajolet-martin sketch itself preserves
 differential privacy: Private counting with minimal space. In *Advances in Neural Information Processing Systems*, virtual, 2020.
- Nina Mesing Stausholm. Improved differentially private euclidean distance approximation. In *Proceedings of the 40th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems (PODS)*, pp. 42–56, Virtual Event, China, 2021.
- Acar Tamersoy, Kevin A. Roundy, and Duen Horng Chau. Guilt by association: large scale malware
 detection by mining file-relation graphs. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*, pp. 1524–1533, New York, NY, 2014.
- Tom Tseng, Laxman Dhulipala, and Julian Shun. Parallel index-based structural graph clustering and its approximation. In *Proceedings of the International Conference on Management of Data (SIGMOD)*, pp. 1851–1864, Virtual Event, China, 2021.
- Santosh S Vempala. *The random projection method*, volume 65. American Mathematical Soc., 2005.
- Pinghui Wang, Yiyan Qi, Yuanming Zhang, Qiaozhu Zhai, Chenxu Wang, John C. S. Lui, and Xiaohong Guan. A memory-efficient sketch method for estimating high similarities in streaming sets. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD)*, pp. 25–33, Anchorage, AK, 2019.
- Tianhao Wang, Jeremiah Blocki, Ninghui Li, and Somesh Jha. Locally differentially private protocols for frequency estimation. In *Proceedings of the 26th USENIX Security Symposium, USENIX Security (USENIX)*, pp. 729–745, Vancouver, Canada, 2017.
- Kang Wei, Jun Li, Ming Ding, Chuan Ma, Howard H. Yang, Farhad Farokhi, Shi Jin, Tony Q. S.
 Quek, and H. Vincent Poor. Federated learning with differential privacy: Algorithms and performance analysis. *IEEE Trans. Inf. Forensics Secur.*, 15:3454–3469, 2020.
- Jia Xu, Zhenjie Zhang, Xiaokui Xiao, Yin Yang, Ge Yu, and Marianne Winslett. Differentially private histogram publication. *VLDB J.*, 22(6):797–822, 2013.
- Yun William Yu and Griffin M. Weber. Hyperminhash: Minhash in loglog space. *IEEE Trans. Knowl. Data Eng.*, 34(1):328–339, 2022.
- Jun Zhang, Zhenjie Zhang, Xiaokui Xiao, Yin Yang, and Marianne Winslett. Functional mechanism: Regression analysis under differential privacy. *Proc. VLDB Endow.*, 5(11):1364–1375, 2012.
- Fuheng Zhao, Dan Qiao, Rachel Redberg, Divyakant Agrawal, Amr El Abbadi, and Yu-Xiang Wang.
 Differentially private linear sketches: Efficient implementations and applications. In Advances in Neural Information Processing Systems (NeurIPS), 2022.

| 702 703 704 | Weijie Zhao, Deping Xie, Ronglai Jia, Yulei Qian, Ruiquan Ding, Mingming Sun, and Ping Li. Distributed hierarchical GPU parameter server for massive scale deep learning ads systems. In Proceedings of Machine Learning and Systems 2020 (MLSys), Austin, TX, 2020. |
|-------------------|--|
| 705 | Erkang Zhu, Ken O. Pu, Fatemeh Nargesian, and Renée J. Miller. Interactive navigation of open |
| 707 | data linkages. Proc. VLDB Endow., 10(12):1837–1840, 2017. |
| 708 | |
| 709 | |
| 710 | |
| 711 | |
| 712 | |
| 713 | |
| 714 | |
| 715 | |
| 716 | |
| 717 | |
| 718 | |
| 719 | |
| 720 | |
| 721 | |
| 722 | |
| 723 | |
| 724 | |
| 725 | |
| 726 | |
| 727 | |
| 728 | |
| 729 | |
| 730 | |
| 731 | |
| 732 | |
| 733 | |
| 734 | |
| 730 | |
| 730 | |
| 738 | |
| 739 | |
| 740 | |
| 741 | |
| 742 | |
| 743 | |
| 744 | |
| 745 | |
| 746 | |
| 747 | |
| 748 | |
| 749 | |
| 750 | |
| 751 | |
| 752 | |
| 753 | |
| 754 | |
| 755 | |

A EXTENSION: DIFFERENTIALLY PRIVATE BIN-WISE CONSISTENT WEIGHTED SAMPLING (DP-BCWS) FOR WEIGHTED JACCARD SIMILARITY

Algorithm 7 Consistent Weighted Sampling (CWS)Input: Non-negative data vector $\boldsymbol{u} \in \mathbb{R}^D_+$ Output: Consistent weighted sampling hash $h^* = (i^*, t^*)$ 1: for every non-zero v_i do2: $r_i \sim Gamma(2, 1), c_i \sim Gamma(2, 1), \beta_i \sim Uniform(0, 1)$ 3: $t_i \leftarrow \lfloor \frac{\log u_i}{r_i} + \beta_i \rfloor, y_i \leftarrow \exp(r_i(t_i - \beta_i))$ 4: $a_i \leftarrow c_i/(y_i \exp(r_i))$ 5: end for6: $i^* \leftarrow arg\min_i a_i, \quad t^* \leftarrow t_{i^*}$

In our main paper, we focused on DP hashing algorithms for the binary Jaccard similarity. Indeed, our algorithm can also be extended to hashing the weighted Jaccard similarity: (recall the definition)

$$J_w(\boldsymbol{u}, \boldsymbol{v}) = \frac{\sum_{i=1}^{D} \min\{u_i, v_i\}}{\sum_{i=1}^{D} \max\{u_i, v_i\}},$$
(5)

for two non-negative data vectors $u, v \in \mathbb{R}_+$. The standard hashing algorithm for (5) is called 778 Consistent Weighted Sampling (CWS) as summarized in Algorithm 7 (loffe, 2010; Manasse et al., 779 2010; Li et al., 2021). To generate one hash value, we need three length-D random vectors $r \sim$ 780 $Gamma(2,1), c \sim Gamma(2,1)$ and $\beta \sim Uniform(0,1)$. We denote Algorithm 7 as a function 781 $CWS(u; r, c, \beta)$. Li et al. (2019) proposed bin-wise CWS (BCWS) which exploits the same idea 782 of binning as in OPH. The binning and densification procedure of BCWS is exactly the same as 783 OPH (Algorithm 2 and Algorithm 3), except that every time we apply CWS, instead of MinHash, 784 to the data in the bins to generate hash values. Note that in CWS, the output contains two values: 785 i^* is a location index similar to the output of OPH, and t^* is a real-value scalar. Prior studies (e.g., 786 Li et al. (2021)) showed that the second element has minimal impact on the estimation accuracy in most practical cases (i.e., only counting the collision of the first element suffices). Therefore, in our 787 study, we only keep the first integer element as the hash output for subsequent learning tasks. 788

789

760 761

762

763

764

765

766

767

768

769

770 771 772

773

For weighted data vectors, we follow the prior DP literature on weighted sets (e.g., Xu et al. (2013); 790 Smith et al. (2020); Dickens et al. (2022); Zhao et al. (2022)) and define the neighboring data vec-791 tors as those who differ in one element. To privatize BCWS, there are also three possible ways 792 depending on the densification option. Since the DP algorithm design for densified BCWS requires 793 rigorous and non-trivial computations which might be an independent study, here we empirically 794 test the (b-bit) DP-BCWS method with random bits for empty bins. The details are provided in 795 Algorithm 8. In general, we first randomly split the data entries into K equal length bins, and ap-796 ply CWS to the data $u_{\mathcal{B}_k}$ in each non-empty bin \mathcal{B}_k using the random numbers $(r_{\mathcal{B}_k}, c_{\mathcal{B}_k}, \beta_{\mathcal{B}_k})$ to 797 generated K hash values (possibly including empty bins). After each hash is truncated to b bits, we uniformly randomly assign a hash value in $\{0, ..., 2^b - 1\}$ to every empty bin. 798

⁷⁹⁹ Using the same proof arguments as Theorem 3.3, we have the following guarantee.

Theorem A.1. Algorithm 8 satisfies ϵ -DP.

802

Empirical evaluation. In Figure 5, we train an l_2 -regularized logistic regression on the DailySports dataset². and report the test accuracy with various b and K values. The l_2 regularization parameter λ is tuned over a fine grid from 10^{-4} to 10. Similar to the results in the previous section, the performance of DP-BCWS becomes stable as long as $\epsilon > 5$. Note that, linear logistic regression only gives $\approx 75\%$ accuracy on original DailySports dataset (without DP). With DP-BCWS, the accuracy can reach $\approx 98\%$ with K = 1024 and $\epsilon = 5$.

809

²https://archive.ics.uci.edu/ml/datasets/daily+and+sports+activities



Figure 5: Test classification accuracy of DP-BCWS on DailySports dataset (Asuncion & Newman, 2007) with l_2 -regularized logistic regression.

In Figure 6, we train a neural network with two hidden layers of size 256 and 128 respectively on MNIST. We use the ReLU activation function and the standard cross-entropy loss. We see that, in a reasonable privacy regime (e.g., $\epsilon < 10$), DP-BCWS is able to achieve $\approx 95\%$ test accuracy with proper K and b combinations (one can choose the values depending on practical scenarios and needs). For example, with b = 4 and K = 128, DP-BCWS achieves $\approx 97\%$ accuracy at $\epsilon = 8$.



Figure 6: Test classification accuracy of DP-BCWS on MNIST with 2-hidden layer neural network.

В PROOFS

Lemma B.1 (Li et al. (2012)). Let $f = |\{i : u_i = 1\}|$, and $I_{emp,k}$ be the indicator function that the k-th bin is empty, and $N_{emp} = \sum_{k=1}^{K} I_{emp,k}$. Suppose mod(D, K) = 0. We have

$$P(N_{emp}=j) = \sum_{\ell=0}^{K-j} (-1)^{\ell} {K \choose j} {K-j \choose \ell} {D(1-(j+\ell)/K) \choose f} / {D \choose f}.$$

Lemma B.2 (Li et al. (2019)). Conditional on the event that m bins are non-empty, let \tilde{f} be the number of non-zero elements in a non-empty bin. Denote d = D/K. The conditional probability distribution of f is given by

$$P\left(\tilde{f}=j\big|m\right) = \frac{\binom{d}{j}H(m-1,f-j|d)}{H(m,f|d)}, \quad j=\max\{1,f-(m-1)d\},...,\min\{d,f-m+1\},$$

where $H(\cdot)$ follows the recursion: for any $0 < k \le K$ and $0 \le n \le f$,

$$H(k,n|d) = \sum_{i=\max\{1,n-(k-1)d\}}^{\min\{d,n-k+1\}} {\binom{d}{i}} H(k-1,n-i|d), \quad H(1,n|d) = {\binom{d}{n}}.$$

B.1 PROOF OF LEMMA 3.1

Proof. Without loss of generality, suppose u and u' differ in the *i*-th dimension, and by the symmetry try of DP, we can assume that $u_i = 1$ and $u'_i = 0$. We know that i is assigned to the $\lfloor mod(\pi(i), d) \rfloor$ th bin. Among the K hash values, this change will affect all the bins that uses the data/hash of the $k^* = \lfloor mod(\pi(i), d) \rfloor$ -th bin (after permutation), both in the first scan (if it is non-empty) and in the densification process. Let N_{emp} be the number of empty bins in h(u), and f be the number of non-zero elements in the k^* -th bin. We have, for $x = 0, ..., K - \lfloor f/d \rfloor$,

$$P(X = x) = \sum_{j=\max(0,K-f)}^{K-\lceil f/d\rceil} \sum_{z=1}^{\min(f,d)} P\left(X = x \middle| \tilde{f} = z, N_{emp} = j\right) P\left(\tilde{f} = z, N_{emp} = j\right)$$

$$= \sum_{j=\max(0,K-f)}^{K-\lceil f/d\rceil} \sum_{j=1}^{\min(f,d)} P\left(X = x \middle| \tilde{f} = z, N_{emp} = j\right) P\left(\tilde{f} = z \middle| K = j\right) P(N_{emp} = j)$$

$$= \sum_{j=\max(0,K-f)}^{K-\lceil f/d\rceil} \sum_{z=1}^{\min(f,d)} P\left(X=x \middle| \tilde{f}=z, N_{emp}=j \right) P\left(\tilde{f}=z | K-j \right) P\left(N_{emp}=j \right),$$

where $P\left(\tilde{f}=z|K-j\right)$ is given in Lemma B.2 and $P\left(N_{emp}=j\right)$ can be calculated by Lemma B.1. To compute the first conditional probability, we need to compute the number of times the k^* -th bin is picked to generated hash values, and the hash values are different for u and u'. Conditional on $\{f = z, N_{emp} = j\}$, denote $\Omega = \{k : \mathcal{B}_k \text{ is empty}\}$, and let R_k be the non-empty bin used for the k-th hash value $h_k(u)$, which takes value in $[K] \setminus \Omega$. We know that $|\Omega| = j$. We can write

$$X = \mathbb{1}\{h_{k^*}(\boldsymbol{u}) \neq h_{k^*}(\boldsymbol{u}')\} + \sum_{k \in \Omega} \mathbb{1}\{R_k = k^*, h_k(\boldsymbol{u}) \neq h_k(\boldsymbol{u}')\}.$$

Here we separate out the first term because the k^* -th hash always uses the k^* -bin. Note that the densification bin selection is uniform, and the bin selection is independent of the permutation for hashing. For the fixed densification, since the hash value $h_{k^*}(u)$ is generated and used for all hash values that use \mathcal{B}_{k^*} , we have

$$P\left(X=x\middle|\tilde{f}=z, N_{emp}=j\right) = \mathbb{1}\{x=0\}\left(1-P_{\neq}\right) + \mathbb{1}\{x>0\}P_{\neq} \cdot g_{bino}\left(x-1; \frac{1}{K-j}, j\right),$$

where $g_{bino}(x; p, n)$ is the probability mass function of the binomial distribution with n trials and success rate p, and $P_{\neq} = P(h_{k^*}(u) \neq h_{k^*}(u')) = \left(1 - \frac{1}{2^b}\right) \frac{1}{z}$. Based on the same reasoning, for re-randomized densification, we have

915
916
$$P\left(X = x \middle| \tilde{f} = z, N_{emp} = j\right) = (1 - P_{\neq}) \cdot g_{bino}\left(x; \frac{P_{\neq}}{K - j}, j\right) + P_{\neq} \cdot g_{bino}\left(x - 1; \frac{P_{\neq}}{K - j}, j\right)$$
917 Combining all the parts together completes the proof.

gether completes the proo

B.2 PROOF OF THEOREM 3.2

Proof. Let u and u' be neighbors only differing in one element. Denote $S = \{k \in [K] : h_k(u) \neq k\}$ $h_k(u')$ and $S^c = [K] \setminus S$. As discussed before, we can verify that for $k \in S_c$, we have $\frac{P(\tilde{h}_k(\boldsymbol{u})=i)}{P(\tilde{h}_k(\boldsymbol{u}')=i)} = 1 \text{ for any } i = 0, ..., 2^b - 1. \text{ For } k \in S, e^{-\epsilon'} \leq \frac{P(\tilde{h}_k(\boldsymbol{u})=i)}{P(\tilde{h}_k(\boldsymbol{u}')=i)} \leq e^{\epsilon'} \text{ holds for any } i = 0, ..., 2^b - 1.$ $i = 0, ..., 2^b - 1$. Thus, for any $Z \in \{0, ..., 2^b - 1\}^K$, the absolute privacy loss can be bounded by

$$\left|\log \frac{P(\tilde{h}(\boldsymbol{u}) = Z)}{P(\tilde{h}(\boldsymbol{u}') = Z)}\right| = \left|\log \prod_{k \in S} \frac{P(\tilde{h}_k(\boldsymbol{u}) = i)}{P(\tilde{h}_k(\boldsymbol{u}') = i)}\right| \le |S|\epsilon' = |S|\frac{\epsilon}{N}.$$
(6)

By Lemma 3.1, with probability $1-\delta$, $|S| \le F_{fix}^{-1}(1-\delta) = N$ for DP-OPH-fix; $|S| \le F_{re}^{-1}(1-\delta) = N$ N for DP-OPH-re. Hence, (6) is bounded by ϵ with probability $1 - \delta$. This proves the (ϵ, δ) -DP.

B.3 PROOF OF THEOREM 3.3

Proof. The proof is similar to the proof of Theorem 3.2. Since the original hash vector h(u) is not densified, there only exists exactly one hash value such that $h_k(u) \neq h_k(u)$ may happen for u' that differs in one element from u. W.l.o.g., assume $u_i = 1$ and $u'_i = 0$, and $i \in \mathcal{B}_k$. If bin k is non-empty for both u and u' (after permutation), then for any $Z \in \{0, ..., 2^b - 1\}^K$, $\log \frac{P(\tilde{h}(\boldsymbol{u})=Z)}{P(\tilde{h}(\boldsymbol{u}')=Z)}$ $| \leq \epsilon$ according to our analysis in Theorem 3.2 (the probability of hash in $[K] \setminus \{k\}$ cancels out). If bin k is empty for u', since $1 \le \frac{e^{\epsilon}}{e^{\epsilon}+2^b-1}/\frac{1}{2^b} \le e^{\epsilon}$ and $e^{-\epsilon} \le \frac{1}{2^b}/\frac{1}{e^{\epsilon}+2^b-1} \le 1$, we also have $\left|\log \frac{P(\tilde{h}(\boldsymbol{u})=Z)}{P(\tilde{h}(\boldsymbol{u}')=Z)}\right| \leq \epsilon$. Therefore, the algorithm is ϵ -DP as claimed.

B.4 PROOF OF THEOREM 3.5

Proof. For the two densified DP-OPH variants, DP-OPH-fix and DP-OPH-re, and the DP MinHash (DP-MH) methods, each full-precision (and unprivatized) hash value of h(u) and h(v) has collision probability equal to $P(h(\boldsymbol{u}) = h(\boldsymbol{v})) = J(\boldsymbol{u}, \boldsymbol{v})$. Let $h^{(b)}(\boldsymbol{u})$ denote the b-bit hash values. Since we assume the last b bits are uniformly assigned, we have $P(h^{(b)}(\boldsymbol{u}) = h^{(b)}(\boldsymbol{v})) = J + (1 - J)\frac{1}{2^b}$. Denote $p = \frac{\exp(\epsilon/N)}{\exp(\epsilon/N) + 2^b - 1}$. By simple probability calculation, the privatized *b*-bit hash values has collision probability

$$\begin{aligned} P(\tilde{h}(\boldsymbol{u}) = \tilde{h}(\boldsymbol{v})) \\ &= P(\tilde{h}(\boldsymbol{u}) = \tilde{h}(\boldsymbol{v})|h^{(b)}(\boldsymbol{u}) = h^{(b)}(\boldsymbol{v}))P(h^{(b)}(\boldsymbol{u}) = h^{(b)}(\boldsymbol{v})) \\ &+ P(\tilde{h}(\boldsymbol{u}) = \tilde{h}(\boldsymbol{v})|h^{(b)}(\boldsymbol{u}) \neq h^{(b)}(\boldsymbol{v}))P(h^{(b)}(\boldsymbol{u}) \neq h^{(b)}(\boldsymbol{v})) \\ &= \left[p^2 + \frac{(1-p)^2}{2^b-1}\right] \left(\frac{1}{2^b} + \frac{2^b-1}{2^b}J\right) + \left[\frac{2p(1-p)}{2^b-1} + \frac{2^b-2}{(2^b-1)^2}(1-p)^2\right] \left(\frac{2^b-1}{2^b} - \frac{2^b-1}{2^b}J\right) \\ &= \left[p^2 + \frac{(1-p)^2}{2^b-1} - \frac{2p(1-p)}{2^b-1} - \frac{2^b-2}{(2^b-1)^2}(1-p)^2\right] \frac{2^b-1}{2^b}J \\ &+ \frac{1}{2^b} \left[p^2 + \frac{(1-p)^2}{2^b-1} + 2p(1-p) + \frac{2^b-2}{2^b-1}(1-p)^2\right] \\ &= \left[p^2 + \frac{(1-p)^2-2(2^b-1)p(1-p)}{(2^b-1)^2}\right] \frac{2^b-1}{2^b}J + \frac{1}{2^b} \left[p^2 + 2p(1-p) + (1-p)^2\right] \\ &= \frac{(2^bp+1)^2}{2^b(2^b-1)}J + \frac{1}{2^b}, \end{aligned}$$
which implies $J = \frac{(2^b-1)(2^bP(\tilde{h}(\boldsymbol{u})=\tilde{h}(\boldsymbol{v}))-1)}{(2^b-1)^2}$. Therefore, let $\hat{J} = \frac{1}{K} \sum_{k=1}^{K} \mathbbm{1}\{\tilde{h}_k(\boldsymbol{u})=\tilde{h}_k(\boldsymbol{v})\}$, then

which implies $J = \frac{(2-1)(2-P(h(u)=h(v))-1)}{(2^b v - 1)^2}$. Therefore, let $J = \frac{1}{K} \sum_{k=1}^{K} \mathbb{1}\{h_k(u) = h_k(v)\}$, then an unbiased estimator of J can be formulated as

$$\hat{J}_{unbias} = \frac{(2^b - 1)(2^b \hat{J} - 1)}{(2^b p - 1)^2}$$

972 B.5 PROOF OF THEOREM 3.6

Proof. As before, define $\hat{J} = \frac{1}{K} \sum_{k=1}^{K} \mathbb{1}\{\tilde{h}_{k}(u) = \tilde{h}_{k}(v)\}$. For all three methods, we know that 975 $\mathbb{E}[\hat{J}] = \frac{(2^{b}p+1)^{2}}{2^{b}(2^{b}-1)}J + \frac{1}{2^{b}}$. Denote $J_{B} = P(h^{(b)}(u) = h^{(b)}(v)) = J + (1-J)\frac{1}{2^{b}}$. $\hat{J}_{unbias} =$ 977 $\frac{(2^{b}-1)(2^{b}\hat{J}-1)}{(2^{b}p-1)^{2}}$.

MinHash. We have

 $Var[\hat{J}] = \mathbb{E}[\hat{J}^2] - \mathbb{E}[\hat{J}]^2$

 $= \frac{1}{K^2} \mathbb{E} \left[\sum_{i=1}^K \mathbb{1}\{\tilde{h}_i(u) = \tilde{h}_i(v)\} \right] + \sum_{i \neq j} \mathbb{1}\{\tilde{h}_i(u) = \tilde{h}_i(v)\} \mathbb{1}\{\tilde{h}_j(u) = \tilde{h}_j(v)\} \right] - \left(\frac{(2^b p + 1)^2}{2^b (2^b - 1)}J + \frac{1}{2^b}\right)^2$ $= \frac{1}{K} \left(\frac{(2^b p + 1)^2}{2^b (2^b - 1)}J + \frac{1}{2^b}\right) + \frac{K - 1}{K}A - \left(\frac{(2^b p + 1)^2}{2^b (2^b - 1)}J + \frac{1}{2^b}\right)^2,$

where $A = \mathbb{E}[\mathbb{1}\{\tilde{h}_i(u) = \tilde{h}_i(v), \tilde{h}_j(u) = \tilde{h}_j(v)\}]$ for $i \neq j$. The key is to calculate A. By symmetry,

A

$$\begin{split} &= P(\tilde{h}_{i}(u) = \tilde{h}_{i}(v), \tilde{h}_{j}(u) = \tilde{h}_{j}(v) | h_{i}^{(b)}(u) = h_{i}^{(b)}(v), h_{j}^{(b)}(u) = h_{j}^{(b)}(v)) P(h_{i}^{(b)}(u) = h_{i}^{(b)}(v), h_{j}^{(b)}(u) = h_{j}^{(b)}(v)) \\ &+ 2P(\tilde{h}_{i}(u) = \tilde{h}_{i}(v), \tilde{h}_{j}(u) = \tilde{h}_{j}(v) | h_{i}^{(b)}(u) = h_{i}^{(b)}(v), h_{j}^{(b)}(u) \neq h_{j}^{(b)}(v)) P(h_{i}^{(b)}(u) = h_{i}^{(b)}(v), h_{j}^{(b)}(u) \neq h_{j}^{(b)}(v)) \\ &+ P(\tilde{h}_{i}(u) = \tilde{h}_{i}(v), \tilde{h}_{j}(u) = \tilde{h}_{j}(v) | h_{i}^{(b)}(u) \neq h_{i}^{(b)}(v), h_{j}^{(b)}(u) \neq h_{j}^{(b)}(v)) P(h_{i}^{(b)}(u) \neq h_{i}^{(b)}(v), h_{j}^{(b)}(u) \neq h_{j}^{(b)}(v)) \\ &:= A_{11} + 2A_{01} + A_{00}. \end{split}$$

By independence, we have

$$A_{11} = \left(p^2 + \frac{(1-p)^2}{2^b - 1}\right)^2 J_B^2$$

$$A_{10} = \left(p^2 + \frac{(1-p)^2}{2^b - 1}\right) \left(\frac{2p(1-p)}{2^b - 1} + \frac{2^b - 2}{(2^b - 1)^2}(1-p)^2\right) J_B(1-J_B)$$

$$A_{00} = \left(\frac{2p(1-p)}{2^b - 1} + \frac{2^b - 2}{(2^b - 1)^2}(1-p)^2\right)^2 (1-J_B)^2,$$

which leads to

$$A = \left(\frac{(2^b p + 1)^2}{2^b (2^b - 1)}J + \frac{1}{2^b}\right)^2.$$

¹⁰¹³ Thus, we have

$$Var[\hat{J}] = \frac{1}{K} \left(\frac{(2^{b}p+1)^{2}}{2^{b}(2^{b}-1)} J + \frac{1}{2^{b}} \right) \left(\frac{2^{b}-1}{2^{b}} - \frac{(2^{b}p+1)^{2}}{2^{b}(2^{b}-1)} J \right)$$

1017 and

$$\begin{aligned} Var[\hat{J}_{unbias,MH}] &= \frac{2^{2b}(2^b - 1)^2}{(2^b p - 1)^4} Var[\hat{J}] \\ &= \frac{1}{K} \left(\frac{(2^b p + 1)^2}{(2^b p - 1)^2} J + \frac{2^b - 1}{(2^b p - 1)^2} \right) \left(\frac{(2^b - 1)^2}{(2^b p - 1)^2} - \frac{(2^b p + 1)^2}{(2^b p - 1)^2} J \right) \end{aligned}$$

DP-OPH-fix. We write $\hat{J} = \frac{1}{K} \sum_{k=1}^{K} (\tilde{I}_k^N + \tilde{I}_k^E)$, where \tilde{I}_k^N is the indicator function of hash collision at the *k*-th bin and when the bin is non-empty, and \tilde{I}_k^N is the indicator function of hash

collision at the k-th bin and when the bin is empty. Similar to previous analysis,

$$Var[\hat{J}] = \frac{1}{K^2} \mathbb{E}\left[(\sum_{k=1}^{K} (\tilde{I}_k^N + \tilde{I}_k^E))^2 \right] - \left(\frac{(2^b p + 1)^2}{2^b (2^b - 1)} J + \frac{1}{2^b} \right)^2$$

$$= \frac{1}{K} \left(\frac{(2^b p + 1)^2}{2^b (2^b - 1)} J + \frac{1}{2^b} \right) + \frac{1}{K^2} A - \left(\frac{(2^b p + 1)^2}{2^b (2^b - 1)} J + \frac{1}{2^b} \right)^2,$$

where

$$A = \mathbb{E}\left[\sum_{i \neq j} (\tilde{I}_{i}^{N} + \tilde{I}_{i}^{E}) (\tilde{I}_{j}^{N} + \tilde{I}_{j}^{E})\right]$$

= $\mathbb{E}_{m}\left[\mathbb{E}[m(m-1)\tilde{I}_{i}^{N}\tilde{I}_{j}^{N} + 2m(K-m)\tilde{I}_{i}^{N}\tilde{I}_{j}^{E} + (K-m)(K-m-1)\tilde{I}_{i}^{E}\tilde{I}_{j}^{E}]|m\right].$ (7)

Here the condition on "|m" means the event that there are m simultaneously non-empty bins. Denote $I_k = 1{h_k(u) = h_k(v)}$ be the collision indicator of the original hash values, and $I_k^{(b)} = \mathbb{1}\{h_k^{(b)}(u) = h_k^{(b)}(v)\}$ be the collision indicator of the *b*-bit hash values. For two non-empty bins *i* and *j*, we have

$$\begin{split} \tau_{11} &:= P(h_i(u) = h_i(v), h_j(u) = h_j(v) | m) = \mathbb{E}[I_i I_j | m] = J\tilde{J}, \\ \tau_{10} &:= P(h_i(u) = h_i(v), h_j(u) \neq h_j(v) | m) = \mathbb{E}[I_i(1 - I_j) | m] = J - J\tilde{J}, \\ \tau_{00} &:= P(h_i(u) \neq h_i(v), h_j(u) \neq h_j(v) | m) = \mathbb{E}[(1 - I_i)(1 - I_j) | m] = 1 - 2J + J\tilde{J}, \end{split}$$

and using total probability formula (conditional on h_i and h_j), $(b) \langle \rangle = (b) \langle \rangle = (b)$

$$\begin{array}{ll} \text{1048} & P(h_i^{(b)}(u) = h_i^{(b)}(v), h_j^{(b)}(u) = h_j^{(b)}(v) | m) = \mathbb{E}[I_i^{(b)} I_j^{(b)} | m] \\ \text{1049} \\ \text{1050} & = \tau_{11} + 2\frac{1}{2^b}\tau_{10} + \frac{1}{2^{2b}}\tau_{00} \\ \text{1051} & = J\tilde{J} + \frac{1}{2^{b-1}}(J - J\tilde{J}) + \frac{1}{2^{2b}}(1 - 2J + J\tilde{J}) := P_{11} \\ \text{1052} & = 0 \\ \text{1053} & = 0 \\ \text{1054} & = 0 \\ \text{1055} & = 0 \\ \text{1055} & = 0 \\ \text{1056} & = 0 \\ \text{1056} & = 0 \\ \text{1057} & = 0 \\ \text{1057} & = 0 \\ \text{1058} & = 0 \\ \text{1058} & = 0 \\ \text{1059} &$$

 $P(h_i^{(b)}(u) = h_i^{(b)}(v), h_j^{(b)}(u) \neq h_j^{(b)}(v) | m) = \mathbb{E}[I_i^{(b)}(1 - I_j^{(b)}) | m] = J_B - P_{11} := P_{10}$ $P(h_i^{(b)}(u) \neq h_i^{(b)}(v), h_i^{(b)}(u) \neq h_i^{(b)}(v)|m) = \mathbb{E}[(1 - I_i^{(b)})(1 - I_i^{(b)})|m] = 1 - 2J_B + P_{11} := P_{00}.$

Thus,

$$\begin{split} & \mathbb{E}[\tilde{I}_{i}^{N}\tilde{I}_{j}^{N}|m] = \left(p^{2} + \frac{(1-p)^{2}}{2^{b}-1}\right)^{2}P_{11} + 2\left(p^{2} + \frac{(1-p)^{2}}{2^{b}-1}\right)\left(\frac{2p(1-p)}{2^{b}-1} + \frac{2^{b}-2}{(2^{b}-1)^{2}}(1-p)^{2}\right)P_{10} \\ & + \left(\frac{2p(1-p)}{2^{b}-1} + \frac{2^{b}-2}{(2^{b}-1)^{2}}(1-p)^{2}\right)^{2}P_{00}. \end{split}$$

For two empty bins *i* and *j*, we have, for fixed densification,

$$\begin{aligned} & \tau_{11,f} = P(h_i(u) = h_i(v), h_j(u) = h_j(v) | m) = \frac{1}{m}J + \frac{m-1}{m}J\tilde{J}, \\ & \tau_{10,f} = P(h_i(u) = h_i(v), h_j(u) \neq h_j(v) | m) = \mathbb{E}[I_i(1 - I_j)] = \frac{m-1}{m}(J - J\tilde{J}), \\ & \tau_{00,f} = P(h_i(u) \neq h_i(v), h_j(u) \neq h_j(v) | m) = \mathbb{E}[(1 - I_i)(1 - I_j)] = 1 - (2 - \frac{1}{m})J + \frac{m-1}{m}J\tilde{J}. \\ & \text{Similarly,} \end{aligned}$$

where

$$\mathbb{E}[\tilde{I}_{i}^{E}\tilde{I}_{j}^{E}|m] = \left(p^{2} + \frac{(1-p)^{2}}{2^{b}-1}\right)^{2} P_{11,f} + 2\left(p^{2} + \frac{(1-p)^{2}}{2^{b}-1}\right) \left(\frac{2p(1-p)}{2^{b}-1} + \frac{2^{b}-2}{(2^{b}-1)^{2}}(1-p)^{2}\right) P_{10,f} + \left(\frac{2p(1-p)}{2^{b}-1} + \frac{2^{b}-2}{(2^{b}-1)^{2}}(1-p)^{2}\right)^{2} P_{00,f},$$

$$P_{11,f} = \tau_{11,f} + \frac{1}{2^{b-1}}\tau_{10,f} + \frac{1}{2^{2b}}\tau_{00,f}, \quad P_{10,f} = J_B - P_{11,f}, \quad P_{00,f} = 1 - 2J_B + P_{11,f}.$$

It is not hard to note that $\mathbb{E}[\tilde{I}_i^N \tilde{I}_j^E | m] = \mathbb{E}[\tilde{I}_i^E \tilde{I}_j^E | m]$. Putting pieces together into (7), we get the variance for DP-OPH-fix.

1080
 DP-OPH-re. For DP-OPH-re, most calculations are the same as DP-OPH-fix. According to Li et al. (2019), we have

$$\begin{aligned} \tau_{11,r} &= P(h_i(u) = h_i(v), h_j(u) = h_j(v) | m) = \mathbb{E}[I_i I_j] = \frac{\zeta(m)}{m} J + \frac{m - \zeta(m)}{m} J \tilde{J}, \\ \tau_{10,r} &= P(h_i(u) = h_i(v), h_j(u) \neq h_j(v) | m) = \mathbb{E}[I_i(1 - I_j)] = \frac{m - \zeta(m)}{m} (J - J \tilde{J}), \\ \tau_{00,r} &= P(h_i(u) \neq h_i(v), h_j(u) \neq h_j(v) | m) = \mathbb{E}[(1 - I_i)(1 - I_j)] \\ &= 1 - (2 - \frac{\zeta(m)}{m}) J + \frac{m - \zeta(m)}{m} J \tilde{J}, \end{aligned}$$

with $\zeta(m) = \mathbb{E}[\frac{1}{\tilde{f}}|m]$ where the conditional distribution of \tilde{f} is given in Lemma B.2. We then get $P_{11,r}, P_{10,r}, P_{00,r}$ correspondingly. Plugging them into the formula above completes the proof. \Box