



**HAL**  
open science

## Smaller, Faster & Lighter KNN Graph Constructions

Rachid Guerraoui, Anne-Marie Kermarrec, Olivier Ruas, François Taïani

► **To cite this version:**

Rachid Guerraoui, Anne-Marie Kermarrec, Olivier Ruas, François Taïani. Smaller, Faster & Lighter KNN Graph Constructions. WWW '20 - The Web Conference 2020, Apr 2020, Taipei Taiwan, France. pp.1060-1070, 10.1145/3366423.3380184 . hal-02888286

**HAL Id: hal-02888286**

**<https://inria.hal.science/hal-02888286>**

Submitted on 2 Jul 2020

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Smaller, Faster & Lighter KNN Graph Constructions

Rachid Guerraoui  
EPFL  
rachid.guerraoui@epfl.ch

Olivier Ruas  
Peking University  
olivier.ruas@pku.edu.cn

Anne-Marie Kermarrec  
Mediego / EPFL  
anne-marie.kermarrec@epfl.ch

François Taïani  
Univ Rennes, Inria, CNRS, IRISA  
francois.taiani@irisa.fr

## ABSTRACT

We propose *GoldFinger*, a new *compact* and *fast-to-compute* binary representation of datasets to approximate Jaccard’s index. We illustrate the effectiveness of GoldFinger on the emblematic big data problem of K-Nearest-Neighbor (KNN) graph construction and show that GoldFinger can drastically accelerate a large range of existing KNN algorithms with little to no overhead. As a side effect, we also show that the compact representation of the data protects users’ privacy *for free* by providing  $k$ -anonymity and  $l$ -diversity. Our extensive evaluation of the resulting approach on several realistic datasets shows that our approach delivers speedups of up to 78.9% compared to the use of raw data while only incurring a negligible to moderate loss in terms of KNN quality. To convey the practical value of such a scheme, we apply it to item recommendation and show that the loss in recommendation quality is negligible.

## ACM Reference Format:

Rachid Guerraoui, Anne-Marie Kermarrec, Olivier Ruas, and François Taïani. 2020. Smaller, Faster & Lighter KNN Graph Constructions. In *Proceedings of The Web Conference 2020 (WWW ’20)*, April 20–24, 2020, Taipei, Taiwan. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3366423.3380184>

## 1 INTRODUCTION

K-Nearest-Neighbor (KNN) graphs<sup>1</sup> play a fundamental role in many big data applications, including search [5, 6], recommendation [8, 29, 32] and classification [39]. A KNN graph is a directed graph of entities (e.g., users, documents etc.), in which each entity (or *node*) is connected to its  $k$  most similar counterparts or *neighbors*, according to a given *similarity metric*. In many applications, this similarity metric is computed from a second set of entities (termed *items*) associated with each node in a bipartite graph (often extended with weights, such as ratings or frequencies). For instance, in a movie rating database, nodes are users, and each user is associated with the movies (items) she has already rated [22].

<sup>1</sup>Note that the problem of computing a complete KNN graph (which we address in this paper) is related but different from that of answering a sequence of KNN queries.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

WWW ’20, April 20–24, 2020, Taipei, Taiwan

© 2020 Association for Computing Machinery.

ACM ISBN 978-1-4503-7023-3/20/04.

<https://doi.org/10.1145/3366423.3380184>

Being able to compute a KNN graph efficiently is crucial in situations that are constrained, either in terms of time or resources. This is the case of *real time*<sup>2</sup> web applications, such as news recommenders, that must regularly recompute their suggestions in short intervals on fresh data to remain relevant. This is also the case of privacy-preserving personal assistants executing learning tasks on personal devices with limited resources [1].

Computing an *exact* KNN graph rapidly becomes intractable on large datasets. Fortunately, many applications only require a good *approximation* of the KNN graph [25, 27]. Recent KNN construction algorithms [8, 19] have therefore sought to reduce the number of similarity computations by exploiting a *greedy strategy*. These techniques, among the most efficient to date, seem, however, to have reached their limits.

In this paper, rather than reducing an algorithm’s complexity (e.g. by decreasing the number of similarity computations), we propose to pursue an orthogonal strategy that is motivated by the system bottlenecks induced by large data volumes: large amounts of data not only stress complex algorithms, they also choke the underlying computation pipelines these algorithms execute on [12].

More precisely, we propose to fingerprint the set of items associated with each node into what we have termed a *Single Hash Fingerprint* (SHF), a 64- to 8096-bit vector summarizing a node’s profile (e.g. the movies the user have seen, the web pages she visited). SHFs are very quick to construct, provide a sufficient approximation of the similarity between two nodes using extremely cheap bit-wise operations. While the main purpose of SHFs is efficient computation, it turns out that SHFs also protect the privacy of users by hiding the original clear-text information and providing interesting properties such as  $k$ -anonymity and  $l$ -diversity. We use these SHFs to rapidly construct KNN graphs, in an overall approach we have dubbed *GoldFinger*.

In this paper, we make the following contributions: (1) we introduce GoldFinger, a *generic* and *efficient* approach to accelerate *any* KNN graph algorithm relying on Jaccard’s index, one of the most common metrics used to compute KNN graphs, which can be tuned to trade space and time for accuracy; (2) we propose a formal analysis of GoldFinger’s estimation mechanism; (3) we formally analyze the privacy protection granted as a side effect by GoldFinger in terms of  $k$ -anonymity and  $l$ -diversity; (4) we extensively evaluate our approach on a range of state-of-the-art KNN graph algorithms such as Locality Sensitive Hashing (LSH), on six representative datasets, and we show that GoldFinger is able to

<sup>2</sup>*Real time* is meant in the sense of *web real-time*, i.e. the proactive push of information to on-line users.

deliver speedups of up to 78.9% against existing approaches, while only incurring a small loss in terms of quality; (5) as a case-study, we use the constructed graphs to produce recommendations, and show that despite the small loss in KNN quality, there is close to no loss in the quality of the derived recommendations.

## 2 PROBLEM, INTUITION, AND APPROACH

For ease of exposition, we consider in the following that nodes are *users* associated with *items* (e.g. web pages, movies, locations), without loss of generality.

### 2.1 Notations and problem definition

We note  $U = \{u_1, \dots, u_n\}$  the set of all users, and  $I = \{i_1, \dots, i_m\}$  the set of all items. The subset of items associated with user  $u$  (a.k.a. its *profile*) is noted  $P_u \subseteq I$ .  $P_u$  is generally much smaller than  $I$  (the universe of all items).

Our aim is to approximate a KNN graph  $G_{\text{KNN}}$  over  $U$  relying on some function *sim* computed over user profiles:

$$\begin{aligned} \text{sim} : U \times U &\rightarrow \mathbb{R} \\ (u, v) &\quad \text{sim}(u, v) = f_{\text{sim}}(P_u, P_v). \end{aligned}$$

$f_{\text{sim}}$  is any similarity function over sets that is typically positively correlated with the number of common items between the two sets, and negatively correlated with the total number of items present in both sets. We focus on Jaccard's index in the rest of the paper [44].

Formally, a KNN graph  $G_{\text{KNN}}$  connects each user  $u \in U$  with a set  $\text{knn}(u)$  of  $k$  other users that maximize the similarity function  $\text{sim}(u, -)$ :

$$\text{knn}(u) \in \underset{v \in U \setminus \{u\}}{\text{argtop}^k} f_{\text{sim}}(P_u, P_v) \quad (1)$$

where  $\text{argtop}^k$  returns the set of  $k$ -tuples of  $U \setminus \{u\}$  that maximize the similarity function  $\text{sim}(u, -)$ <sup>3</sup>.

Computing an exact KNN graph is particularly expensive: an exhaustive search requires  $O(|U|^2)$  similarity computations. Many scalable approaches therefore seek to construct *an approximate KNN graph*  $\widehat{G}_{\text{KNN}}$ , i.e., to find for each user  $u$  a neighborhood  $\widehat{\text{knn}}(u)$  that is as close as possible to an exact KNN neighborhood [8, 19]. The meaning of 'close' depends on the context, but in most applications, a good approximate neighborhood  $\widehat{\text{knn}}(u)$  is one whose aggregate similarity (its *quality*) comes close to that of an exact KNN set  $\text{knn}(u)$ .

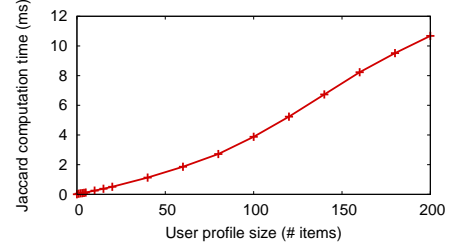
We capture how well the average similarity of an approximated graph  $\widehat{G}_{\text{KNN}}$  compares against that of an exact KNN graph  $G_{\text{KNN}}$  with the *average similarity* of  $\widehat{G}_{\text{KNN}}$ :

$$\text{avg\_sim}(\widehat{G}_{\text{KNN}}) = \mathbb{E}_{\substack{(u,v) \in U^2 \\ v \in \widehat{\text{knn}}(u)}} f_{\text{sim}}(P_u, P_v), \quad (2)$$

i.e. the average similarity of the edges of  $\widehat{G}_{\text{KNN}}$ . We then define the *quality* of  $\widehat{G}_{\text{KNN}}$  as

$$\text{quality}(\widehat{G}_{\text{KNN}}) = \frac{\text{avg\_sim}(\widehat{G}_{\text{KNN}})}{\text{avg\_sim}(G_{\text{KNN}})}. \quad (3)$$

<sup>3</sup>In other words,  $\text{argtop}^k$  generalizes the concept of argument of the maximum (usually noted  $\text{argmax}$ ) to the  $k$  top values of a function over a finite discrete set.



**Figure 1: Cost, averaged over 4.9 millions computations between randomly generated profiles on a Intel Xeon.**

A quality close to 1 indicates that the approximate neighborhoods have a quality close to that of ideal neighborhoods, and can replace them with little loss in most applications.

With the above notations, we can summarize our problem as follows: for a given dataset  $(U, I, (P_u)_{u \in U})$  and item-based similarity  $f_{\text{sim}}$ , we wish to compute an approximate  $\widehat{G}_{\text{KNN}}$  in the shortest time with the highest overall quality.

### 2.2 Intuition

A large portion of a KNN graph's construction time often comes from computing individual similarity values (up to 90% of the total construction time [9]). This is because computing explicit similarity values on even medium-sized profiles can be relatively expensive. Figure 1 shows on its y-axis the time required to compute Jaccard's index  $J(P_1, P_2) = \frac{|P_1 \cap P_2|}{|P_1 \cup P_2|}$  between two user profiles of the same size, when this size varies (x-axis). The cost of computing a single index is relatively high even for medium-size profiles: 2.7 ms for two random profiles of 80 items, a typical profile size of the datasets we have considered.

Earlier KNN graph construction approaches have therefore sought to limit the number of similarity computations [8, 19]. They typically adopt a greedy strategy, starting from a random graph, and progressively converging to a better KNN approximation. This strategy dramatically reduces the similarity computations they perform, and are easily parallelizable, but it is now difficult to see how their greedy component could be further improved.

In order to overcome the inherent cost of similarity computations, we propose in this paper to target the data on which computations run, rather than the algorithms that drive these computations. This strategy stems from the observation that *explicit* datastructures (hash tables, arrays) incur substantial costs. To avoid these costs, we advocate the use of *fingerprints*, a *compact*, *binary*, and *fast-to-compute* representation of data. Our intuition is that, with almost no overhead, fingerprints can capture enough of the characteristics of the data to provide a good approximation of similarity values, while drastically reducing the cost of computing these similarities.

### 2.3 GoldFinger and Single Hash Fingerprints

Our approach, dubbed *GoldFinger*, extracts from each user's profile a *Single Hash Fingerprint* (SHF for short). An SHF is a pair  $(B, c) \in \{0, 1\}^b \times \mathbb{N}$  comprising a bit array  $B = (\beta_x)_{x \in [0..b-1]}$  of  $b$  bits, and an integer  $c$ , which records the number of bits set to 1 in  $B$  (its L1 norm, which we call the *cardinality* of  $B$  in the following). The SHF of a user's profile  $P$  is computed by hashing each item of the profile

**Table 1: Effect of SHFs on computation time of Jaccard's index, compared to Fig. 1 (80 items).**

SHF length (bits)	Comp. Time (ms)	Speedup $ P  = 80$
64	0.011	253
256	0.032	84
1024	0.120	23
4096	0.469	6

into the array and setting to 1 the associated bit

$$\beta_x = \begin{cases} 1 & \text{if } \exists e \in P : h(e) = x, \\ 0 & \text{otherwise,} \end{cases}$$

$$c = \|\beta_x\|_1$$

where  $h()$  is a uniform hash function from  $I$  to  $\llbracket 0..b-1 \rrbracket$ , and  $\|\cdot\|_1$  counts the number of bits set to 1.

*Benefits in terms of space and speed.* The length  $b$  of the bit array  $B$  is usually much smaller than the total number of items, causing collisions, and a loss of information. This loss is counterbalanced by the highly efficient approximation SHFs can provide of any set-based similarity. The Jaccard's index of two user profiles  $P_1$  and  $P_2$  can be estimated from their respective SHFs ( $B_1, c_1$ ) and ( $B_2, c_2$ ) with

$$\widehat{J}(P_1, P_2) = \frac{\|B_1 \text{ AND } B_2\|_1}{c_1 + c_2 - \|B_1 \text{ AND } B_2\|_1}, \quad (4)$$

where  $B_1 \text{ AND } B_2$  represents the bitwise AND of the bit-arrays of the two profiles. This formula exploits two observations that hold generally with few collisions in the bit arrays (a point we revisit below). First, the size of a set of items  $P$  can be estimated from the cardinality of its SHF ( $B_P, c_P$ ):

$$|P| \approx \|B_P\|_1 = c_P. \quad (5)$$

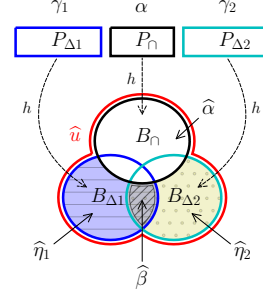
Second, the bit array  $B_{(P_1 \cap P_2)}$  of the intersection of two profiles  $P_1 \cap P_2$  can be approximated with the bitwise AND of their respective bit-arrays,  $B_1$  and  $B_2$ :

$$B_{(P_1 \cap P_2)} \approx (B_1 \text{ AND } B_2). \quad (6)$$

Equation (4) combines these two observations along with some simple set algebra ( $|P_1 \cup P_2| = |P_1| + |P_2| - |P_1 \cap P_2|$ ).

The computation incurred by (4) is much faster than on explicit profiles and is independent of the actual size of the explicit profiles. This is illustrated in Table 1 which shows the computation time of Eq. (4) on the same profiles as Figure 1 for SHFs of different lengths (as in Fig. 1, the values are averaged over 4.9 million computations). For instance, estimating Jaccard's index between two SHFs of 1024 bits (the default in our experiments) takes 0.120 ms, a 23-fold speedup compared to two explicit profiles of 80 items.

**The link with Bloom Filters and collisions:** SHFs can be interpreted as a highly degraded form of Bloom filters, and suffer from errors arising from collisions, as Bloom filters do. However, while Bloom filters are designed to test whether individual elements belong to a set, SHFs are designed to approximate set similarities. Bloom filters often employ multiple hash functions to minimize false positives. Those increase single-bit collisions and degrade the approximation provided by SHFs.



**Figure 2: Illustration of the collisions caused by  $h$  on two profiles  $P_1 = P_{\Delta 1} \cup P_{\cap}$  and  $P_2 = P_{\Delta 2} \cup P_{\cap}$ . Capital letters ( $P$  and  $B$ ) denote sets, Greek letters denote sizes, and the hat symbol ( $\hat{\cdot}$ ) represent random variables.**

## 2.4 Formal analysis of the Jaccard estimator

For readability in this section, we will generally treat bit arrays (i.e. belonging to  $\{0, 1\}^n$ ) as sets of bit positions (belonging to  $\mathcal{P}(\llbracket 0, b-1 \rrbracket)$ ). We will also note  $B_X$  the set of bits set to 1 by the (sub)profile  $P_X$ :  $B_X = h(P_X)$ .

For two given profiles  $P_1$  and  $P_2$ , the distribution of  $\widehat{J}(P_1, P_2)$  is governed by how the random hash function  $h$  maps the items of  $P_1$  and  $P_2$  onto the bit positions  $\llbracket 0, b-1 \rrbracket$ . The analysis of this random mapping is made simpler if we distinguish between the items that are present in both  $P_1$  and  $P_2$  ( $P_{\cap} = P_1 \cap P_2$ ) from those are unique to  $P_1$  (resp.  $P_2$ ), noted  $P_{\Delta 1} = P_1 \setminus P_{\cap}$  (resp.  $P_{\Delta 2} = P_2 \setminus P_{\cap}$ ). These three sets are represented by the three top rectangles of Figure 2.

$P_{\cap}$ ,  $P_{\Delta 1}$ , and  $P_{\Delta 2}$  are disjoint by definition, but their bit images  $B_{\cap}$ ,  $B_{\Delta 1}$ , and  $B_{\Delta 2}$  by  $h$  (shown as circles in Figure 2) are typically not, because of collisions. To analyze these collisions we introduce the following three helping sets:

- $B_{\hat{\eta}_1}$  are the bits of  $B_{\Delta 1}$  (corresponding to items only present in  $P_1$ ) that do not collide with those of  $B_{\cap}$ :  $B_{\hat{\eta}_1} = B_{\Delta 1} \setminus B_{\cap}$ ;
- $B_{\hat{\eta}_2}$  is the equivalent for  $P_2$ :  $B_{\hat{\eta}_2} = B_{\Delta 2} \setminus B_{\cap}$ ;
- and  $B_{\hat{\beta}}$  contains the collisions between  $B_{\hat{\eta}_1}$  and  $B_{\hat{\eta}_2}$ :  

$$B_{\hat{\beta}} = B_{\hat{\eta}_1} \cap B_{\hat{\eta}_2}.$$

If we note  $\hat{\alpha}$ ,  $\hat{\eta}_1$ ,  $\hat{\eta}_2$ ,  $\hat{\beta}$  the sizes of the sets  $B_{\cap}$ ,  $B_{\hat{\eta}_1}$ ,  $B_{\hat{\eta}_2}$ , and  $B_{\hat{\beta}}$ , respectively, we can express  $\widehat{J}(P_1, P_2)$  as

$$\widehat{J}(P_1, P_2) = \frac{\hat{\alpha} + \hat{\beta}}{\hat{\alpha} + \hat{\eta}_1 + \hat{\eta}_2 - \hat{\beta}} = \frac{2\hat{\alpha} + \hat{\eta}_A + \hat{\eta}_B}{\hat{u}} - 1, \quad (7)$$

where  $\hat{u} = \hat{\alpha} + \hat{\eta}_1 + \hat{\eta}_2 - \hat{\beta}$  is the number of bits set to 1 by either  $P_1$  or  $P_2$ , i.e.  $\hat{u} = |B_1 \cup B_2| = \|B_1 \text{ OR } B_2\|_1$ .

The distribution of  $\widehat{J}(P_1, P_2)$  is determined by the joint distribution of the random values  $(\hat{u}, \hat{\alpha}, \hat{\eta}_1, \hat{\eta}_2)$  when  $h$  (shown with dashed arrows in Figure 2) is chosen uniformly randomly among all functions  $h$  that map  $P_{\Delta 1} \cup P_{\cap} \cup P_{\Delta 2}$  onto  $\llbracket 0, b-1 \rrbracket$ . (Note that  $\hat{\beta} = \hat{\alpha} + \hat{\eta}_1 + \hat{\eta}_2 - \hat{u}$  is determined by the four other values, and therefore does not appear in the quadruplet.)

**THEOREM 2.1.** *The probability distribution of  $(\hat{u}, \hat{\alpha}, \hat{\eta}_1, \hat{\eta}_2)$  is*

$$\mathbb{P}(\hat{u}, \hat{\alpha}, \hat{\eta}_1, \hat{\eta}_2 | \alpha, \gamma_1, \gamma_2) = \frac{\text{Card}_h(\hat{u}, \hat{\alpha}, \hat{\eta}_1, \hat{\eta}_2, \alpha, \gamma_1, \gamma_2)}{b^{(\alpha + \gamma_1 + \gamma_2)}}, \quad (8)$$

where

- $\alpha$ ,  $\gamma_1$ , and  $\gamma_2$  are resp. the sizes of the sets  $P_\cap$ ,  $P_{\Delta_1}$ , and  $P_{\Delta_2}$  introduced earlier;
- $\text{Card}_h()$  is the number of hashing functions from  $P_{\Delta_1} \cup P_\cap \cup P_{\Delta_2}$  onto  $\llbracket 0, b-1 \rrbracket$  that produce the quadruplet  $(\widehat{u}, \widehat{\alpha}, \widehat{\eta}_1, \widehat{\eta}_2)$ :

$$\text{Card}_h(\widehat{u}, \widehat{\alpha}, \widehat{\eta}_1, \widehat{\eta}_2, \alpha, \gamma_1, \gamma_2) =$$

$$\binom{b}{\widehat{u}} \binom{b}{\widehat{\alpha}} \binom{\widehat{u}-\widehat{\alpha}}{\widehat{\beta}} \binom{\widehat{u}-\widehat{\alpha}-\widehat{\beta}}{\widehat{\eta}_1-\widehat{\beta}} \widehat{\alpha}! \binom{\alpha}{\widehat{\alpha}} \\ \times \xi(\gamma_1, \widehat{\eta}_1 + \widehat{\alpha}, \widehat{\eta}_1) \times \xi(\gamma_2, \widehat{\eta}_2 + \widehat{\alpha}, \widehat{\eta}_2);$$

- $\binom{\alpha}{\widehat{\alpha}}$  denotes Stirling's number of the second type;
- and  $\xi(x, y, z)$  is the number of functions  $f : X \mapsto Y$  from a finite set  $X$  onto a finite set  $Y$ , that are surjective on a subset  $Z \subseteq Y$  of  $Y$ , with  $x = |X|$ ,  $y = |Y|$ ,  $z = |Z|$ , defined by

$$\xi(x, y, z) = \sum_{k=0}^z (-1)^k \binom{z}{k} (y-k)^x.$$

PROOF. The formula for  $\mathbb{P}(\widehat{u}, \widehat{\alpha}, \widehat{\eta}_1, \widehat{\eta}_2 | \alpha, \gamma_1, \gamma_2)$  derives from a counting strategy on the functions from  $P_\cup = P_{\Delta_1} \cup P_\cap \cup P_{\Delta_2}$  onto  $\llbracket 0, b-1 \rrbracket$ : the denominator  $b^{(\alpha+\gamma_1+\gamma_2)}$  is the total number of such functions.

The numerator,  $\text{Card}_h()$ , is the number of functions  $h$  that yield precisely the quadruplet  $(\widehat{u}, \widehat{\alpha}, \widehat{\eta}_1, \widehat{\eta}_2)$ . We obtain  $\text{Card}_h()$  with a constructive argument. Because  $P_\cap$ ,  $P_{\Delta_1}$  and  $P_{\Delta_2}$  are disjoint,  $h$  can be seen as the piece-wise combination of three independent random functions  $h|_{P_\cap}$ ,  $h|_{P_{\Delta_1}}$ , and  $h|_{P_{\Delta_2}}$ , where  $h|_X$  denotes the restriction of  $h$  to a set  $X$ .

To construct  $h$ , we start by choosing  $B_\cup = h(P_\cup)$  within  $\llbracket 0, b-1 \rrbracket$ . As  $\widehat{u} = |h(P_\cup)|$ , there are  $\binom{b}{\widehat{u}}$  such choices.

Similar arguments for  $B_\cap$ ,  $B_{\widehat{\beta}}$ , and  $B_{\widehat{\eta}_1} \setminus B_\cap$  yield that overall the total number of choices for  $B_\cup$ ,  $B_\cap$ ,  $B_{\widehat{\beta}}$ ,  $B_{\widehat{\eta}_1} \setminus B_\cap$  (and  $B_{\widehat{\eta}_2} \setminus B_\cap$ ) is

$$\binom{b}{\widehat{u}} \binom{\widehat{u}-\widehat{\alpha}}{\widehat{\beta}} \binom{\widehat{u}-\widehat{\alpha}-\widehat{\beta}}{\widehat{\eta}_1-\widehat{\beta}}.$$

Once these supporting sets have been chosen, we pick  $h|_{P_\cap}$ ,  $h|_{P_{\Delta_1}}$ , and  $h|_{P_{\Delta_2}}$ .  $h|_{P_\cap}$  is a surjection from  $P_\cap$  to  $B_\cap$ . There are  $\widehat{\alpha}! \binom{\alpha}{\widehat{\alpha}}$  such surjections, where  $\binom{\alpha}{\widehat{\alpha}}$  is Stirling's number of the second type.

$h|_{P_{\Delta_1}}$  maps  $P_{\Delta_1}$  onto  $B_\cup$ , but only needs to be surjective on  $B_{\Delta_1} \setminus B_\cap = B_{\widehat{\eta}_1}$ . In addition,  $h|_{P_{\Delta_1}}$  only maps elements onto  $B_{\Delta_1} \subseteq B_{\widehat{\eta}_1} \cup B_\cap$ , whose cardinal is  $\widehat{\eta}_1 + \widehat{\alpha}$ .

Let us note  $\xi(x, y, z)$  the number of functions  $f : X \mapsto Y$  from a finite set  $X$  onto a finite set  $Y$ , that are surjective on a subset  $Z \subseteq Y$  of  $Y$ , with  $x = |X|$ ,  $y = |Y|$ ,  $z = |Z|$ . Using an inclusion-exclusion argument we have

$$\xi(x, y, z) = \sum_{k=0}^z (-1)^k \binom{z}{k} (y-k)^x$$

There are therefore  $\xi(\gamma_1, \widehat{\eta}_1 + \widehat{\alpha}, \widehat{\eta}_1)$  functions  $h|_{P_{\Delta_1}}$ . Similarly there are  $\xi(\gamma_2, \widehat{\eta}_2 + \widehat{\alpha}, \widehat{\eta}_2)$  functions  $h|_{P_{\Delta_2}}$ . The product of the above quantities yields  $\text{Card}_h()$ , and as a result (8).  $\square$

Combining the formula for  $\mathbb{P}(\widehat{u}, \widehat{\alpha}, \widehat{\eta}_1, \widehat{\eta}_2 | \alpha, \gamma_1, \gamma_2)$  provided by Theorem 2.1 and Eq. (7), we can compute the distribution and moments of the estimator  $\widehat{J}$  for any  $\alpha, \gamma_1, \gamma_2$ . For instance,  $\widehat{J}$ 's mean is  $\mathbb{E}(\widehat{J}(P_1, P_2)) = \sum(\widehat{u}, \widehat{\alpha}, \widehat{\eta}_1, \widehat{\eta}_2) \mathbb{P}(\widehat{u}, \widehat{\alpha}, \widehat{\eta}_1, \widehat{\eta}_2 | \alpha, \gamma_1, \gamma_2) \times \left( \frac{2\widehat{\alpha} + \widehat{\eta}_A + \widehat{\eta}_B}{\widehat{u}} - 1 \right)$ .

Figure 3 uses this strategy to plot the behavior of the estimator  $\widehat{J}$  against the real Jaccard index when comparing a profile,  $P_1$  of 100 items with other profiles  $P_2, P_3, P_4$  of varying sizes.

$\widehat{J}$  is biased. For instance when  $J(P_1, P_2) = 0.25$ —the right vertical dashed line in Fig. 3— $\widehat{J}$  returns an average value of 0.286. However, this absolute bias has little impact on KNN algorithms, which only need to *order* nodes correctly, rather than to predict exact similarities. Instead, the spread and overlap of the values returned by  $\widehat{J}$  drive its impact on a KNN approximation. This effect is thankfully limited: for instance, if we assume that  $P_2$  belongs to  $P_1$ 's exact KNN neighborhood with a similarity of  $J(P_1, P_2) = 0.25$ , an algorithm using  $\widehat{J}$  might exclude  $P_2$  if it finds another profile  $P_{2'}$  that  $\widehat{J}$  wrongly considers as more similar:

$$J(P_1, P_{2'}) < J(P_1, P_2) \text{ and } \widehat{J}(P_1, P_{2'}) > \widehat{J}(P_1, P_2).$$

Figure 3 shows this misordering has a very low probability of occurring (less than 2%) when  $J(P_1, P_{2'})$  is lower than 0.17 (left vertical dashed line). This is because  $\widehat{J}(P_1, P_2)$  has a 99% probability of being *higher* than 0.254 (1%-percentile value, shown as a solid horizontal line with a rectangle mark on the y-axis), while  $P_{2'}$  profiles with a real Jaccard's index lower than 0.17 have a 99% probability of being *lower* than this cut-off value (and this is roughly independent of  $P_{2'}$ 's size). This phenomenon is shown in more detail in Figure 4 when  $P_{2'}$  contains 100 items: the error in similarity caused by the use of SHFs to compute  $\widehat{J}$  is bounded with high probability by a quantity proportional to the spread of estimated values. When SHFs are large enough compared to the size of the profiles being compared (here with  $b = 1024$  bits), this spread is limited, but it increases as  $b$  gets smaller (Figure 5), highlighting a natural trade-off between compactness and accuracy that we will revisit in Section 4.

## 2.5 Privacy guarantees of GoldFinger

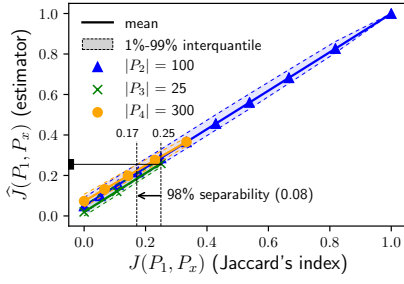
The noise introduced by collisions brings extra privacy benefits that are obtained at *no additional cost*, as an inherent side effect of our approach. Collisions obfuscate a user's profile and thus make it harder to guess this profile from its compacted SHF. This obfuscation can allow users to compute locally their SHF before sending it to some untrusted KNN-construction service.

We characterize the level of protection granted by GoldFinger along two standard measures of privacy, *k-anonymity* [42], and *ℓ-diversity* [34]. For this analysis, we assume an honest but curious attacker who wants to discover the profile  $P_u$  of a particular user  $u$ , knowing its corresponding SHF  $(B_u, c_u)$ . We assume the attacker knows the item set  $I$ , the user set  $U$  and the hash function  $h$ , but does not know the set of user profiles  $\{P_u\}_{u \in U}$ . More importantly, for a given bit position  $x \in \llbracket 0, b-1 \rrbracket$ , we assume the attacker can compute  $H_x = h^{-1}(x)$ , the preimage of  $x$  by  $h$ . How much information does  $(B_u, c_u)$  leak about the initial profile  $P_u$ ?

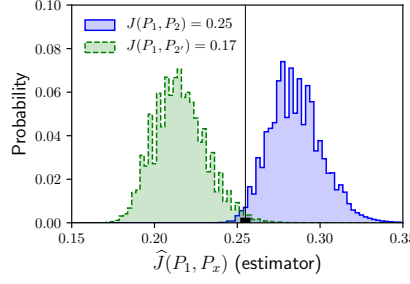
### 2.5.1 k-anonymity.

*Definition 2.2.* Consider an obfuscation mechanism  $obf : \mathcal{X} \mapsto \mathcal{Y}$  that maps a clear-text input  $x \in \mathcal{X}$  to an obfuscated value in  $\mathcal{Y}$ .  $obf()$  is *k-anonymous* for  $x \in \mathcal{X}$ , if the observed obfuscated value  $obf(x)$  is indistinguishable from that of at least  $k-1$  other explicit input values. Formally,  $obf()$  is *k-anonymous* for  $x \in \mathcal{X}$  iff

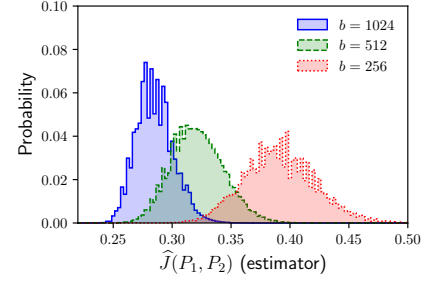
$$|obf^{-1}(obf(x))| \geq k. \quad (9)$$



**Figure 3: Mean and 1%-99% interquartile of  $\hat{J}$  between a profile  $P_1$  of 100 items,  $J(P_1, P_2) = 0.25$  and  $J(P_1, P_{2'}) = 0.17$  (resp- and 3 profiles  $P_2, P_3, P_4$  of varying sizes, resented in bins of 0.0025).  $|P_1| = |P_2| =$  using SHFs of  $b = 1024$  bits.**



**Figure 4: Distribution of  $\hat{J}(P_1, P_x)$  when  $J(P_1, P_2) = 0.25$  and  $J(P_1, P_{2'}) = 0.17$  (resp- and 3 profiles  $P_2, P_3, P_4$  of varying sizes, resented in bins of 0.0025).  $|P_1| = |P_2| =$  using SHFs of  $b = 1024$  bits.**



**Figure 5: Distribution of  $\hat{J}(P_1, P_2)$ , when  $J(P_1, P_2) = 0.25$ ,  $|P_1| = |P_2| = 100$  items, for different values of  $b$ .**

**THEOREM 2.3.** *GoldFinger ensures  $(2^{\frac{m}{b}-1})^{c_u}$ -anonymity in expectation for a given SHF  $(B_u, c_u)$  of length  $b$ , and cardinality  $c_u$ , where  $m = |I|$  is the size of the item set.*

**PROOF.** Let  $x$  be the index of a bit set to 1. Let  $H_x = h^{-1}(x)$  the set of all the items which are hashed by  $h$  to  $x$ , it is on average of size  $\frac{m}{b}$  with a uniformly random hashing function. Thus  $\mathcal{P}(H_x)$  (the powerset of  $H_x$ , sometimes noted  $\{0, 1\}^{H_x}$ ), whose cardinality is lower bounded in expectation by  $2^{\frac{m}{b}}$  (due to the convexity of  $x \rightarrow 2^x$ ), is the set of all possible sub-profiles that will set the bit  $x$  to 1. Since we do not consider the empty set, the total number of sets is bounded in expectation by  $2^{\frac{m}{b}-1}$ . All of these sub-profiles are indistinguishable once hashed, ensuring at least  $(2^{\frac{m}{b}-1})$ -anonymity in expectation for this bit. For every bit set to one, there are  $2^{\frac{m}{b}-1}$  possible set of items, leading to a  $(2^{\frac{m}{b}-1})^{c_u}$ -anonymity in expectation, since all pre-images  $(H_x)_x$  are pair-wise disjoint.  $\square$

The anonymity granted by GoldFinger increases with the size of the item set  $m = |I|$ . For instance, one of the datasets we consider, AmazonMovies, has 171,356 items. With 1024-bit SHFs (the typical size we use), GoldFinger provides  $2^{166}$ -anonymity, i.e. each compacted profile is indistinguishable from at least  $2^{166} \approx 9.35 \times 10^{49}$  possible profiles.

**2.5.2  $\ell$ -diversity.** Although  $k$ -anonymity provides a measure of the difficulty to recover the complete profile  $P_u$  of a user  $u$ , it does not cover cases in which an attacker would seek to guess some partial information about  $u$ . This type of question is better captured by a second metric,  $\ell$ -diversity [34]. The  $\ell$ -diversity model ensures that, for a given SHF  $(B_u, c_u)$ , the actual profile  $P_u$  it was created from is indistinguishable from  $\ell - 1$  other profiles  $\{P_i\}_{i \in [1.. \ell-1]}$ , and that these profiles form a *well-represented set*.

The difference with  $k$ -anonymity lies in this notion of *well-representedness*. In our case it means that we cannot infer any taste from possible profiles: for example in a movie dataset, if all the possible profiles of a given SHF include science-fiction movies, you can infer that the user enjoys science fiction.  $\ell$ -diversity measures the difficulty of such inferences.

**Definition 2.4.** Consider an obfuscation mechanism  $obf : \mathcal{P}(I) \mapsto \mathcal{Y}$  that maps a set of items  $P \subseteq I$  to an obfuscated

value in  $\mathcal{Y}$ .  $obf()$  is  $\ell$ -anonymity for  $P \subseteq I$ , if the observed obfuscated value  $obf(P)$  is indistinguishable from that of at least  $\ell - 1$  other explicit profiles  $Q = \{P_i\}_{i \in [1.. \ell-1]}$  that are pair-wise disjoint:  $\forall P_1, P_2 \in Q : P_1 \cap P_2 = \emptyset$ . Expressed formally<sup>4</sup>,  $obf()$  is  $\ell$ -anonymity for  $P \subseteq I$  iff

$$\min_{\substack{Q \subseteq obf^{-1}(obf(P)) \setminus \{P\} \\ \forall P_1, P_2 \in Q : P_1 \cap P_2 = \emptyset}} |Q| \geq \ell - 1. \quad (10)$$

**THEOREM 2.5.** *For a given SHF  $(B_u, c_u)$  of length  $b$  and cardinality  $c_u$ , SHF ensures  $(\frac{m}{b})$ -diversity for  $(B_u, c_u)$ , provided  $m \gg b \gg c_u$ .*

**PROOF.** The reasoning is similar to that of  $k$ -anonymity. Let  $x$  be the index of a bit set to 1. Because  $m \gg b$ ,  $h()$  closely approximates a perfectly uniform hash function, and  $|H_x| = \frac{m}{b}$  items are hashed into this bit w.h.p. Because  $b \gg c_u$ , we can assume these  $|H_x|$  are pairwise independent for different  $x \in B_u$ . Choosing any arbitrary order on items, let us note  $i_j^x$  the  $j^{\text{th}}$  element of the pre-image  $H_x$  for each bit  $x$  set to 1 in  $B_u$ . W.l.o.g., we can choose our order so that  $i_0^x \in P_u$  for all  $x$ . Consider now the profiles  $Q_j = \cup_{x: B_u[x]=1} \{i_j^x\}$  for  $j \in [1.. \frac{m}{b} - 1]$ . By construction (i)  $P_u \neq Q_j$  for  $j \geq 1$ , (ii) the  $\{Q_j\}_{j \in [1.. \frac{m}{b} - 1]}$  are pair-wise disjoint, and (iii) they are all indistinguishable from  $P_u$  once mapped onto their SHF.  $\square$

Applying Theorem 2.5 to the dataset AmazonMovies, using 1024 bit long SHFs, we ensure 167-diversity. It is important to note that SHFs were not designed to provide privacy but faster similarity computation. The above two privacy properties are therefore obtained “for free”, as a side effect of our approach. They do not preclude, however, the use of additional strengthening privacy-protection mechanisms—such as the insertion of random noise to the SHF [2] to obtain differential privacy guaranties [20]—albeit typically at the cost of reduced performances.

## 3 EXPERIMENTAL SETUP

### 3.1 Datasets

We evaluate GoldFinger on six publicly available datasets (Table 2). To apply Jaccard’s index, we binarize each dataset by only keeping in a user profile  $P_u$  those items that user  $u$  has rated higher than 3.

<sup>4</sup>This definition, adapted to our context, differs slightly from that of the original paper, but leads in practice to the same result.

**Table 2: Description of the datasets used in our experiments**

Dataset	Users	Items	Scale	Ratings > 3
movielens1M (ml1M)	6,038	3,533	1-5	575,281
movielens10M (ml10M)	69,816	10,472	0.5-5	5,885,448
movielens20M (ml20M)	138,362	22,884	0.5-5	12,195,566
AmazonMovies (AM)	57,430	171,356	1-5	3,263,050
DBLP	18,889	203,030	5	692,752
Gowalla (GW)	20,270	135,540	5	1,107,467

*MovieLens.* MovieLens [22] is a group of anonymous datasets containing movie ratings collected on-line between 1995 and 2015 by GroupLens Research [40]. The datasets (before binarization) contain movie ratings on a 0.5-5 scale by users who have at least performed 20 ratings. We use 3 versions of the dataset, movielens1M (ml1M), movielens10M (ml10M) and movielens20M (ml20M), containing between 575,281 and 12,195,566 positive ratings (i.e. higher than 3).

*AmazonMovies.* AmazonMovies [36] (AM) is a dataset of movie reviews from Amazon collected between 1997 and 2012. We restrain our study to users with at least 20 ratings (before binarization) to avoid users with not enough data (this problem, the *cold start problem*, is generally treated separately [28]). After binarization, the dataset contains 57,430 users; 171,356 items; and 3,263,050 ratings.

*DBLP.* DBLP [46] is a dataset of co-authorship from the DBLP computer science bibliography. In this dataset, both the user set and the item set are subsets of the author set. If two authors have published at least one paper together, they are linked, which is expressed in our case by both of them rating each other with a rating of 5. As with AM, we only consider users with at least 20 ratings. The resulting dataset contains 18,889 users, 203,030 items; and 692,752 ratings.

*Gowalla.* Gowalla [16] (GW) is a location-based social network. As DBLP, both the user set and the item set are subsets of the set of the users of the social network. The undirected friendship link from  $u$  to  $v$  is represented by  $u$  rating  $v$  with a 5. As previously, only the users with at least 20 ratings are considered. The resulting dataset contains 20,270 users, 135,540 items; and 1,107,467 ratings.

### 3.2 Baseline algorithms and competitors

We apply GoldFinger to four existing KNN algorithms: Brute Force (as a reference point), NNDescent [19], Hyrec [8] and LSH [23]. We compare the performance and results of each of these algorithms in their native form (*native* for short) and when accelerated with GoldFinger (*GolFi* for short). For completeness, we also consider two direct competitors, *random sampling* and *b-bit minwise hashing* [31].

*Random sampling.* An intuitive idea to lower Jaccard’s index cost is to sample each user profile by keeping a fixed number of items per profile.

*b-bit minwise hashing (MinHash).* A standard technique to approximate Jaccard’s index values between sets is the *MinHash* [10] algorithm. MinHash creates multiple independent permutations on the IDs of an item universe, and keeps for each profile the item with the smallest ID, after each permutation. The Jaccard’s index between two profiles can be estimated by counting the proportion

of minimal IDs that are equal in the compacted representation of the two profiles. *MinHash* was extended to keep only the lowest  $b$  bits of each minimal element [31]. This approach, called *b-bit minwise hashing* (MinHash for short), creates very compact binary summaries of profiles, comparable to our SHFs, from which a Jaccard’s index can be estimated.

*Brute force.* The Brute Force algorithm computes the similarities between every pair of profiles, performing a constant number of similarity computations equal to  $\frac{n \times (n-1)}{2}$ . While this is computationally intensive, this algorithm produces an exact KNN graph.

*NNDescent.* NNDescent [19] constructs an approximate KNN graph (or ANN) by relying on a local search and by limiting the number of similarities computations.

NNDescent starts from an initial random graph, which is then iteratively refined to converge to an ANN graph. During each iteration, for each user  $u$ , NNDescent compares all the pairs  $(u_i, u_j)$  among the neighbors of  $u$ , and updates the neighborhoods of  $u_i$  and  $u_j$  accordingly. NNDescent includes a number of optimizations: it exploits the order on user IDs, and maintains update flags to avoid computing several times the same similarities. It also reverses the current KNN approximation to increase the space search among neighbors. The algorithm stops either when the number of updates during one iteration is below the value  $\delta \times k \times n$ , with a fixed  $\delta$ , or after a fixed number of iterations.

*Hyrec.* Hyrec [8] uses a strategy similar to that of NNDescent, exploiting the fact that a neighbor of a neighbor is likely to be a neighbor. As NNDescent, Hyrec starts with a random graph which is then refined. Hyrec primarily differs from NNDescent in its iteration strategy. At each iteration, for each user  $u$ , Hyrec compares all the neighbors’ neighbors of  $u$  with  $u$ , rather than comparing  $u$ ’s neighbors between themselves. Hyrec also does not reverse the current KNN graph. As NNDescent, it stops when the number of changes is below the value  $\delta \times k \times n$ , with a fixed  $\delta$ , or after a fixed number of iterations.

*LSH.* Locality-Sensitive-Hashing (LSH) [23] reduces the number of similarity computations by hashing each user into several buckets. Neighbors are then selected among users found in the same buckets. To ensure that similar users tend to be hashed into the same buckets, LSH uses min-wise independent permutations of the item set as its hash functions, similarly to the MinHash algorithm [10].

### 3.3 Experimental settings

We set  $k$  to 30 (the neighborhood size). The parameter  $\delta$  of Hyrec and NNDescent is set to 0.001, and their maximum number of iterations to 30. The number of hash functions for LSH is 10. GoldFinger uses 1024 bits long SHFs computed with Jenkins’ hash function [24].

*Evaluation metrics.* We measure the effect of GoldFinger on Brute Force, Hyrec, NNDescent and LSH along with two main metrics: (i) their computation *time* (measured from the start of the algorithm, once the dataset has been prepared), and (ii) the *quality* of the resulting KNN (Sec. 2.1). When applying GoldFinger to recommendation, we also measure the *recall* obtained by the recommender. Throughout our experiments, we use a 5-fold cross-validation, and average results on the 5 resulting runs.

**Table 3: Preprocessing time of each dataset for the native approach, b-bit minwise hashing (MinHash) & GoldFinger.**

Dataset	Native	MinHash	GoldFinger	speedup (×)
ml1M	0.37s	6.24s	0.31s	20.1
ml10M	3.90s	203s	3.24s	62.7
ml20M	8.71s	820s	7.06s	116.1
AM	3.40s	3250s	1.92s	1692.7
DBLP	0.42s	944s	0.29s	3255.2
GW	0.47s	594s	0.40s	1485.0

GoldFinger is orders of magnitude faster than MinHash, whose overhead is prohibitive.

*Implementation details and hardware.* We have implemented Brute Force, Hyrec, NNDescent and LSH (with and without GoldFinger) in Java 1.8. Our experiments run on a 64-bit Linux server with two Intel Xeon E5420@2.50GHz, totaling 8 hardware threads, 32GB of memory, and a HDD of 750GB. Unless stated otherwise, we use all 8 threads. Our code is available online<sup>5</sup>.

## 4 EVALUATION RESULTS

We first compare GoldFinger against Random sampling and MinHash, before discussing in detail its impact on Brute Force, Hyrec, NNDescent, and LSH.

### 4.1 Random sampling and MinHash

*Random sampling.* When applied to Brute Force on ml10M, a random sampling of 75 items per profile decreases the computation time from 2028s to 1626s, while delivering a quality of 0.94. By contrast, GoldFinger (with 1024 bits) is substantially faster (producing a graph in 606s, a 62.7% improvement compared to random sampling) while producing a similar graph (also delivering a quality of 0.94). Similar results are obtained on other datasets or on Hyrec, (with temporal gains by GoldFinger ranging from 54.9% to 72.2% compared to sampling, and producing better graphs across the board, except in one case in which GoldFinger’s quality is slightly lower, by 0.04). Together these results show GoldFinger consistently outperforms random sampling.

*MinHash.* We use  $b = 4$  and 256 permutations for BBHM (a configuration that provides the best trade-off between time and KNN quality). MinHash decreases the computation time on ml10M from 2028s to 1028s with a Brute Force approach, while delivering a quality of 0.93. GoldFinger (with 1024 bits) is both substantially faster (producing a graph in 606s, a 41.1% improvement compared to MinHash), and better (delivering a quality of 0.94). Similar results are obtained on other datasets, or when using LSH (with an average temporal gain of 34.8% compared to MinHash, temporal gains ranging from 2.7% to 58.8%, and while producing better graphs across the board). When used on AmazonMovies with NNDescent and Hyrec, MinHash causes these two algorithms to collapse, yielding qualities below 0.11. Together these results show GoldFinger outperforms MinHash by a wide margin. Furthermore, computing MinHash summaries is extremely costly (as it requires creating a large number of permutations on the entire item set), which renders the approach self-defeating in our context. Table 3 summarizes

<sup>5</sup><https://gitlab.inria.fr/oruas/SamplingKNN>

**Table 4: Computation time and KNN quality with native algorithms (*nat.*) and GoldFinger (*GolFi*).**

	algorithms	comp. time (s)			KNN quality		
		<i>nat.</i>	<i>GolFi</i>	<i>gain%</i>	<i>nat.</i>	<i>GolFi</i>	<i>loss</i>
ml1M	Brute Force	19.0	4.0	78.9	1.00	0.93	0.07
	Hyrec	14.4	4.4	69.4	0.98	0.92	0.06
	NNDescent	19.0	11.0	42.1	1.00	0.93	0.07
	LSH	9.5	<b>3.0</b>	68.4	0.98	0.92	0.06
ml10M	Brute Force	2028	606	70.1	1.00	0.94	0.06
	Hyrec	314	<b>110</b>	65.0	0.96	0.90	0.06
	NNDescent	374	147	60.7	1.00	0.93	0.07
	LSH	689	255	63.0	0.99	0.94	0.06
ml20M	Brute Force	8393	2616	68.8	1.00	0.92	0.08
	Hyrec	842	<b>289</b>	65.7	0.95	0.88	0.07
	NNDescent	919	383	58.3	0.99	0.92	0.07
	LSH	2859	1060	62.9	0.99	0.93	0.06
AM	Brute Force	1862	435	76.6	1.00	0.96	0.04
	Hyrec	235	<b>62</b>	73.6	0.82	0.93	-0.11
	NNDescent	324	91	71.9	0.98	0.95	0.03
	LSH	144	141	2.1	0.98	0.96	0.02
DBLP	Brute Force	100	46	54.0	1.0	0.82	0.18
	Hyrec	46	27	41.3	0.86	0.81	0.05
	NNDescent	31	<b>24</b>	22.6	0.98	0.82	0.16
	LSH	40	38	2.6	0.87	0.86	0.01
Gowalla	Brute Force	160	54	66.3	1.0	0.78	0.22
	Hyrec	39	<b>22</b>	43.6	0.95	0.78	0.17
	NNDescent	45	26	42.2	1.0	0.79	0.21
	LSH	30	27	3.7	0.87	0.82	0.05

GoldFinger yields the shortest computation times across all datasets (in bold), yielding gains (*gain*) of up to 78.9% against native algorithms. The loss in quality is at worst moderate, ranging from 0.22 to an improvement of 0.11.

the time required to load and construct the internal representation of each dataset when using a native (explicit) approach, GoldFinger (using Jenkins’ hash function [24]), and MinHash. Whereas GoldFinger is slightly faster than a native approach (as it does not need to create extensive in-memory objects to store the dataset), MinHash is one to 3 orders of magnitude slower than GoldFinger (1692 times slower on AmazonMovies for instance). This kind of overhead makes it impractical for environments with limited resources, and we, therefore, do not consider MinHash in the rest of our evaluation.

### 4.2 Brute Force, Hyrec, NNDescent, and LSH

The impact of GoldFinger (*GolFi*) on these four algorithms is summarized in Table 4 in terms of execution time and KNN quality. The columns marked *nat.* indicate the results with the native algorithms, while those marked *GolFi* contain those with GoldFinger. The columns in italics show the gain in computation time brought by GoldFinger (*gain %*), and the loss in quality (*loss*). The fastest time for each dataset is shown in bold. Excluding LSH for space reasons, the same results are shown in Figures 6 (time) and 7 (quality).

Overall, GoldFinger delivers the fastest computation times across all datasets, for a small loss in quality ranging from 0.22 (with Brute



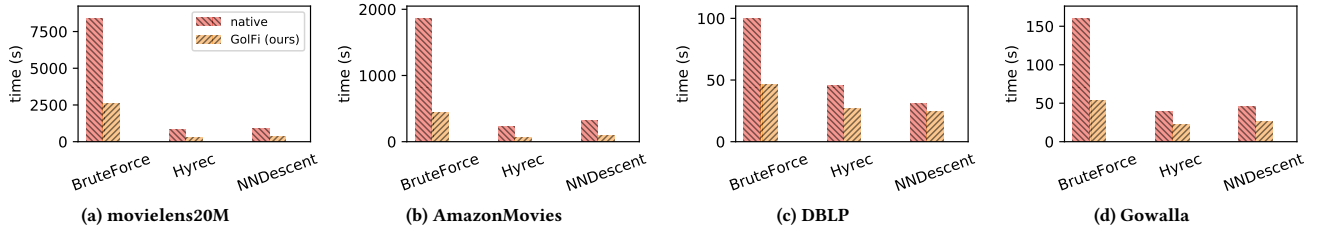


Figure 6: Execution time using a 1024 bits SHF (lower is better).

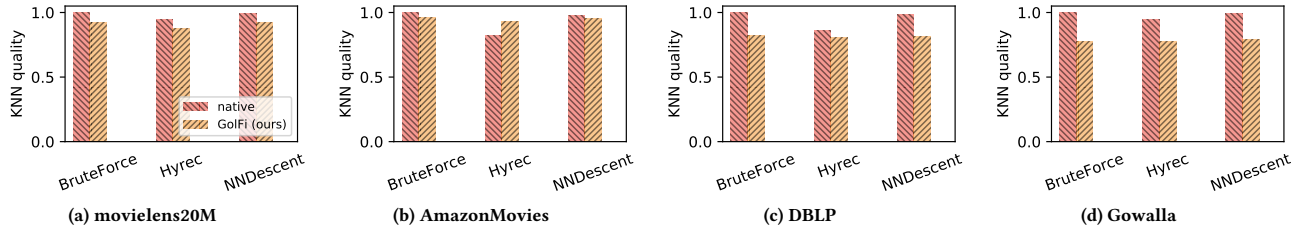


Figure 7: KNN quality using a 1024 bits SHF (higher is better).

Force on Gowalla) to an improvement of 0.11 (Hyrec on Amazon-Movies). Excluding LSH on AmazonMovies, DBLP, and Gowalla for the moment, GoldFinger is able to reduce computation time substantially, from 42.1% (NNDescent on ml1M) to 78.9% (Brute Force on ml1M), corresponding to speedups of 1.72 and 4.74 respectively. Experiments show that the most important parameter is the distribution of ratings, rather than the sparsity or the number of items/users. The more ratings are concentrated onto a small subset of items, reducing collisions, the higher the accuracy provided by GoldFinger. Ratings appear to be more concentrated of few items in movie datasets which explain the differences between the datasets.

GoldFinger only has a limited effect on the execution time of LSH on the AmazonMovies, DBLP and Gowalla datasets. This lack of impact can be explained by the characteristics of LSH and the datasets. LSH must first create user buckets using permutations on the item universe, an operation that is proportional to the number of items. Because AmazonMovies, DBLP and Gowalla are comparatively very sparse (Table 2), the buckets created by LSH tend to contain few users. As a result, the overall computation time is dominated by the bucket creation, limiting the impact of GoldFinger.

In spite of these results, GoldFinger consistently outperforms native LSH on these datasets, for instance, taking 62s (with Hyrec) instead of 141s with LSH on AmazonMovies (a speedup of  $\times 2.27$ ), for a comparable quality.

### 4.3 Memory and cache accesses

By compacting profiles, GoldFinger reduces the amount of memory needed to process a dataset. To gauge this effect, we use the `perf`<sup>6</sup> command line tool to profile the memory accesses of GoldFinger. `perf` uses hardware counters to measure accesses to the cache hierarchy (L1, LLC, and physical memory). To eliminate accesses performed during the dataset preparation, we subtract the values

Table 5: L1 stores and L1 loads with the native algorithms (*nat.*) and GoldFinger (*GolFi*) on ml10M.

algo	L1 stores ( $\times 10^{12}$ )			L1 loads ( $\times 10^{12}$ )		
	nat.	GolFi	gain%	nat.	GolFi	gain%
Brute Force	2.82	0.34	87.9	8.26	1.08	86.9
Hyrec	0.35	0.08	77.1	1.14	0.28	75.4
NNDescent	0.57	0.16	71.9	1.93	0.59	69.4
LSH	0.84	0.85	-1.19	2.96	2.90	2.03

GoldFinger drastically reduces the number of L1 accesses, yielding reductions (*gain*) of up to 87.9%.

returned by `perf` when only preparing the dataset from the values obtained on a full execution.

Table 5 summarizes the measures obtained on Brute Force, Hyrec, NNDescent and LSH on ml10M, both without (*native*) and with GoldFinger (*GolFi*). We only show L1 accesses for space reasons, since LLC and RAM accesses are negligible in comparison. Except on LSH, GoldFinger significantly reduces the number of L1 cache loads and stores, confirming the benefits of GoldFinger in terms of memory footprint. For LSH, L1 accesses are almost not impacted by GoldFinger. Again, we conjecture this is because memory accesses are dominated by the creation of buckets.

### 4.4 GoldFinger in action: recommendations

We evaluate the applicability of GoldFinger in the context of a concrete application, namely a recommender. Item recommendation is one of the main applications of KNN graphs, and consists in providing every user with a list of items she is likely to rate positively. To do so, we compute for each user  $u$  and each item  $i$  not known to  $u$  that is present in  $u$ 's KNN neighborhood a score  $score(u, i)$ , using a weighted average of the ratings given by other users in  $u$ 's KNN:

<sup>6</sup>[https://perf.wiki.kernel.org/index.php/Main\\_Page](https://perf.wiki.kernel.org/index.php/Main_Page)

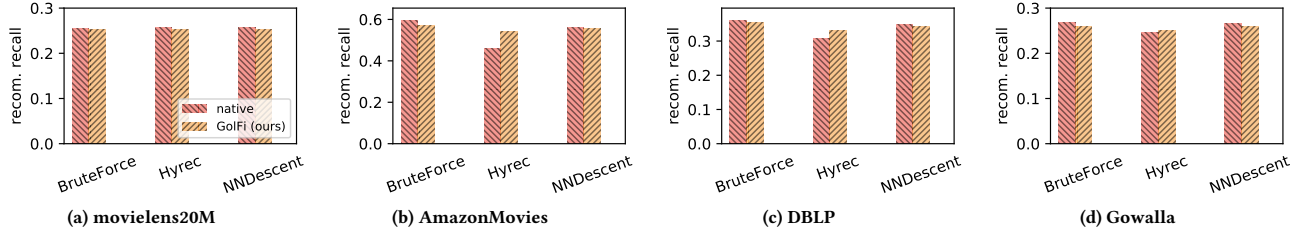


Figure 8: Recommendation recall using a 1024 bits SHF (higher is better). GoldFinger's (GolFi) recall loss is negligible.

$$\text{score}(u, i) = \frac{\sum_{v \in \widehat{\text{knn}}(u)} r(u, i) \times \text{sim}(u, v)}{\sum_{v \in \widehat{\text{knn}}(u)} \text{sim}(u, v)}.$$

Using the KNN graphs computed in the previous sections, we recommend 30 items to each user in every dataset. Since we use a 5-fold cross-validation, we use the 1/5 of each dataset not used in an experiment as our testing set, and consider a recommendation successful if the user has produced a positive rating for the recommended item in the testing set. We evaluate the quality recommendation using *recall*, i.e. the number of successful recommendations divided by the number of positively rated items hidden in the testing set.

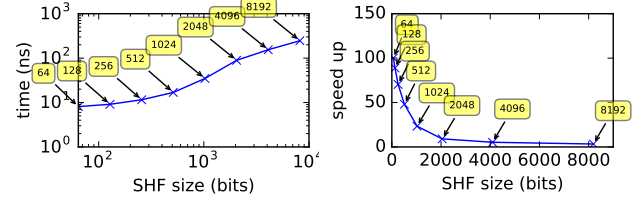
Figure 8 shows the recall of the recommendation made with the native algorithms and with their GoldFinger counterparts on all datasets. The impact on recall spans from a 4.16% drop (brute force on AmazonMovies) to a gain of +16.98% (Hyrec on AmazonMovies) with a gain of +0.84% on average. These results clearly show that the small drop in KNN quality caused by GoldFinger has close to no impact on the outcome of the recommender, confirming the practical relevance of GoldFinger.

## 5 SENSITIVITY ANALYSIS

The size of SHFs ( $b$ ) determines the number of collisions occurring when computing SHFs, and when intersecting them. It thus affects the obtained KNN quality. Shorter SHFs also deliver higher speedups, resulting in an inherent trade-off between execution time and quality. In this section we focus on ml10M.

### 5.1 Impact on the similarity computation time

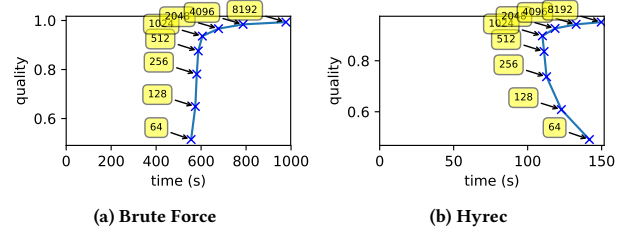
SHFs aim at drastically decreasing the cost of individual similarity computations. To assess this effect, Figure 9 shows the average computation time of one similarity computation when using SHFs (Eq. 4) and its corresponding speed-up. The measures were obtained by computing with a multithreaded program the similarities between two sets of  $5 \times 10^4$  users, sampled randomly from ml10M. The first set of users is divided into several parts, one for each thread. On each thread, each user of the part of the first set is compared to every user of the second set. The total time required is divided by the total number of similarities computed,  $2.5 \times 10^9$ , and then averaged over 4 runs. The computation time is linear in the size of the SHF. Computation time spans from 8 nanoseconds to 250 nanoseconds using SHF, against 800 nanoseconds with real profiles. The other datasets show similar results.



(a) Time of one similarity computation

(b) Speed-up

Figure 9: Effect of the size of the SHF on the similarity computation time, on ml10M.



(a) Brute Force

(b) Hyrec

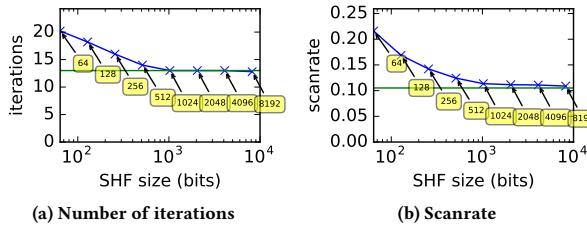
Figure 10: Relation between the execution time and the quality in function of the size of SHF.

### 5.2 Impact on the execution of the algorithm

Figure 10 shows how the overall execution time and the quality of Brute Force and Hyrec evolve when we increase the size of the SHFs. (LSH presents a similar behavior to that of Brute Force, and NNDescent to that of Hyrec.)

As expected, larger SHFs cause Brute Force to take longer to compute, while delivering a better KNN quality (Fig. 10a). The overall computation time does not exactly follow that of individual similarity computations (Figure 9a), as the algorithm involves additional bookkeeping work, such as maintaining the KNN graph and iterating over the dataset.

The KNN quality of Hyrec shows a similar trend, increasing with the size of SHFs. The computation time of Hyrec presents, however, an unexpected pattern: it first *decreases* when SHFs grow from 64 to 1024 bits, before increasing again from 1024 to 4096 bits (Figure 10b). This difference is due to the different nature of the two approaches. The Brute Force algorithm computes a fixed number of similarities, which is independent of the distribution of similarities between users. By contrast, Hyrec adopts a greedy approach: the number of similarities computed depends on the iterations performed by the algorithm, and these iterations are highly dependent on the



**Figure 11: Effect of compression on the convergence of Hyrec on ml10M. GoldFinger converges to the native approach when the size of SHF augments.**

distribution of similarity values between pairs of users (what we have termed the *similarity topology* of the dataset), a phenomenon we return to in the following section.

### 5.3 Impact on estimated similarity values

Still, most of the pairs of users have both their exact and approximated similarities below 0.1: 92% for 1024 bits and 94% for 4096 bits. This confirms our initial intuition (Section 2): two users with low similarity are likely to get a low approximation using GoldFinger. The large majority of the pairs of users in the KNN (as directed edges) do not see their similarity changed much by the use of SHFs. The pairs that experience a large variation between their real and their estimated similarity are too few in numbers to have a decisive impact on the quality of the resulting KNN graph.

This explains why the Brute Force algorithm experiments a decrease in execution time along with a small drop in quality with GoldFinger. Hyrec and NNDescent, however, iterate recursively on node neighborhoods and are therefore more sensitive to the overall distribution of similarity values. The recursive effect is the reason why Hyrec and NNDescent’s execution time first decreases as SHFs grow in size (as mentioned earlier, in Fig. 10).

To shed more light on this effect, Figure 11 shows the number of iterations and the corresponding scanrate performed by Hyrec for SHF sizes varying from 64 to 8192 bits. The scanrate is the number of similarity computations executed by Hyrec+GoldFinger divided by the number of pairs of users. The green horizontal line represents the results when using native Hyrec. As expected, the behavior of the GoldFinger version converges to that of native Hyrec as the size of the SHFs increases. Interestingly, short SHFs ( $< 1024$  bits) cause Hyrec to require more iterations to converge, leading to a higher scanrate. When this occurs, the performance gain on individual similarity computations (Figure 9) does not compensate for this higher scanrate, explaining the counter-intuitive pattern observed in Figure 10b.

## 6 RELATED WORK

For small datasets, KNNs can be solved efficiently using specialized data structures [7, 33, 38]. These solutions, unfortunately, do not scale, as computing an exact KNN efficiently remains an open problem. Most practical approaches, therefore, compute an approximation of the KNN graph (ANN), as we do.

A first way to accelerate the computation time is to decrease the number of comparisons between users, taking the risk to miss

some neighbors. *Recursive Lanczos Bisection* [15] computes an ANN graph using a divide-and-conquer method, while *NNDescent* [19] and *Hyrec* [8] rely on local search, i.e. they assume that a neighbor of a neighbor is likely to be a neighbor, and thus drastically decrease the scan rate, delivering substantial speedups. *KIFF* [9] computes similarities only when users share an item. KIFF works particularly well on sparse datasets but has more difficulties with denser datasets such as the ones we studied. *Locality Sensitive Hashing* (LSH) [23] allows fast ANN graph computations by hashing users into buckets. The neighbors are selected only between the users of the same buckets. Several hash functions have been introduced, for different metrics [10, 11, 14]. All of the above works can be combined with our approach—as we have demonstrated in the case of *NNDescent*, *Hyrec*, and *LSH*—and are thus complementary to our contribution.

Another strategy to accelerate a KNN graph’s construction consists in compacting users’ profiles, in order to obtain a fast approximation of the similarity metric. Keeping only a fraction of the profiles speeds-up Jaccard computation [26] but the resulting approach is not as fast as GoldFinger. *Minwise hashing* [4, 31] approximates Jaccard’s index by only keeping a small subset of items for each user. It is space efficient but has a prohibitive preprocessing.

Cui *et al.* [18] use a bit array to represent profiles: each feature has its value rounded to either 0 or 1, and stored in one bit. The approach is not scalable for the dataset we study. Closer to our work, Gorai *et al.* [21] use Bloom filters to encode the profiles and then estimates Jaccard’s index by using a bitwise AND. Despite providing privacy, the resulting loss in precision is prohibitive. *Sketches* [13, 17, 41] are other compacted datastructures, which have been used for instance to find frequent items in data streams [3]. Unfortunately sketches are not optimized for set intersection.

Privacy has today become a major concern for many information systems. Multiple metrics with different semantics have been proposed to characterized privacy protection such as *k-anonymity* [43], *l-diversity* [35], *t-closeness* [30] or *differential privacy* [20]. An increasing number of works currently seek to add those properties to existing machine learning techniques [2, 37, 45]. Blip [2] for instance provides differential privacy to users’ profiles encoded into Bloom filters by injecting additional noise. By contrast, our approach naturally provides *k-anonymity* and *l-diversity*, for free.

## 7 CONCLUSION

We have proposed *GoldFinger*, a new *compact* and *fast-to-compute* representation of datasets which accelerates the computation of Jaccard’s index. We have illustrated the effectiveness of GoldFinger on KNN graph construction. As a side effect, GoldFinger also protects users’ privacy for free.

Our extensive evaluation shows that GoldFinger is able to drastically accelerate the construction of KNN graphs against the native versions of prominent KNN construction algorithms such as NNDescent or LSH while incurring a small to moderate loss in quality, and close to no overhead in dataset preparation compared to the state of the art. We have also precisely characterized the privacy protection provided by GoldFinger in terms of *k-anonymity* and *l-diversity*.

## REFERENCES

- [1] Snips: The ai platform for voice-enabled devices. <https://snips.ai/>. last accessed April 19, 2017.
- [2] M. Alaggan, S. Gams, and A.-M. Kermarrec. Blip: Non-interactive differentially-private similarity computation on bloom filters. In *SSS*, 2012.
- [3] E. Anceaume, Y. Busnel, N. Rivetti, and B. Sericola. Identifying global icebergs in distributed streams. In *SRDS*, 2015.
- [4] Y. Bachrach and E. Porat. Sketching for big data recommender systems using fast pseudo-random fingerprints. In *ICALP*, 2013.
- [5] X. Bai, M. Bertier, R. Guerraoui, A.-M. Kermarrec, and V. Leroy. Gossiping personalized queries. In *EDBT*, 2010.
- [6] M. Bertier, D. Frey, R. Guerraoui, A.-M. Kermarrec, and V. Leroy. The gossip anonymous social network. In *Middleware*, 2010.
- [7] A. Beygelzimer, S. Kakade, and J. Langford. Cover trees for nearest neighbor. In *ICML*, 2006.
- [8] A. Boutet, D. Frey, R. Guerraoui, A.-M. Kermarrec, and R. Patra. Hyrec: leveraging browsers for scalable recommenders. In *Middleware*, 2014.
- [9] A. Boutet, A.-M. Kermarrec, N. Mittal, and F. Taiani. Being prepared in a sparse world: the case of knn graph construction. In *ICDE*, 2016.
- [10] A. Z. Broder. On the resemblance and containment of documents. In *Compression and Complexity of Sequences*, 1997.
- [11] A. Z. Broder, S. C. Glassman, M. S. Manasse, and G. Zweig. Syntactic clustering of the web. *Comp. Networks and ISDN Sys.*, 1997.
- [12] R. Bruno, D. Patricio, J. Simão, L. Veiga, and P. Ferreira. Runtime object lifetime profiler for latency sensitive big data applications. In *EuroSys*, 2019.
- [13] M. Charikar, K. Chen, and M. Farach-Colton. Finding frequent items in data streams. *Automata, languages and programming*, 2002.
- [14] M. S. Charikar. Similarity estimation techniques from rounding algorithms. In *STOC*, 2002.
- [15] J. Chen, H.-r. Fang, and Y. Saad. Fast approximate knn graph construction for high dimensional data via recursive lanczos bisection. *J. of ML Research*, 2009.
- [16] E. Cho, S. A. Myers, and J. Leskovec. Friendship and mobility: User movement in location-based social networks. In *KDD*, 2011.
- [17] G. Cormode and S. Muthukrishnan. An improved data stream summary: the count-min sketch and its applications. *J. of Algorithms*, 55(1), 2005.
- [18] B. Cui, H. T. Shen, J. Shen, and K.-L. Tan. Exploring bit-difference for approximate knn search in high-dimensional databases. In *AusDM*, 2005.
- [19] W. Dong, C. Moses, and K. Li. Efficient k-nearest neighbor graph construction for generic similarity measures. In *WWW*, 2011.
- [20] C. Dwork. Differential privacy: A survey of results. In *TAMC*, 2008.
- [21] M. Gorai, K. Sridharan, T. Aditya, R. Mukkamala, and S. Nukavarapu. Employing bloom filters for privacy preserving distributed collaborative knn classification. In *WICT*, 2011.
- [22] F. M. Harper and J. A. Konstan. The movielens datasets: History and context. *ACM Trans. Interact. Intell. Syst.*, 2015.
- [23] P. Indyk and R. Motwani. Approximate nearest neighbors: towards removing the curse of dimensionality. In *STOC*, 1998.
- [24] B. Jenkins. Hash functions. *Dr Dobbs Journal*, 1997.
- [25] K. Keeton, C. B. Morrey III, C. A. Soules, and A. Veitch. Lazybase: freshness vs. performance in information management. *SIGOPS Op. Sys. Review*, 2010.
- [26] A. Kermarrec, O. Ruas, and F. Taiani. Nobody cares if you liked star wars: KNN graph construction on the cheap. In *Euro-Par*, 2018.
- [27] A. Labrinidis and N. Roussopoulos. Exploring the tradeoff between performance and data freshness in database-driven web servers. *The VLDB Journal*, 2004.
- [28] X. N. Lam, T. Vu, T. D. Le, and A. D. Duong. Addressing cold-start problem in recommendation systems. In *IMCOM*, 2008.
- [29] J. J. Levandoski, M. Sarwat, A. Eldawy, and M. F. Mokbel. Lars: A location-aware recommender system. In *ICDE*, 2012.
- [30] N. Li, T. Li, and S. Venkatasubramanian. t-closeness: Privacy beyond k-anonymity and l-diversity. In *ICDE*, 2007.
- [31] P. Li and A. C. KÄnig. Theory and applications of b-bit minwise hashing. *CACM*, 2011.
- [32] G. Linden, B. Smith, and J. York. Amazon. com recommendations: Item-to-item collaborative filtering. *Internet Comp.*, 2003.
- [33] T. Liu, A. W. Moore, K. Yang, and A. G. Gray. An investigation of practical approximate nearest neighbor algorithms. In *Advances in neural information processing systems*, 2004.
- [34] A. Machanavajjhala, J. Gehrke, D. Kifer, and M. Venkatasubramanian. L-diversity: privacy beyond k-anonymity. In *ICDE*, 2006.
- [35] A. Machanavajjhala, J. Gehrke, D. Kifer, and M. Venkatasubramanian. l-diversity: Privacy beyond k-anonymity. In *ICDE*, 2006.
- [36] J. J. McAuley and J. Leskovec. From amateurs to connoisseurs: modeling the evolution of user expertise through online reviews. In *WWW*, 2013.
- [37] F. McSherry and I. Mironov. Differentially private recommender systems: building privacy into the net. In *KDD*, 2009.
- [38] A. W. Moore. The anchors hierarchy: Using the triangle inequality to survive high dimensional data. In *UAI*, 2000.
- [39] N. Nodarakis, S. Sioutas, D. Tsoumakos, G. Tzimas, and E. Pitoura. Rapid aknn query processing for fast classification of multidimensional data in the cloud. *CoRR*, abs/1402.7063, 2014.
- [40] P. Resnick, N. Iacovou, M. Suchak, P. Bergstrom, and J. Riedl. Grouplens: an open architecture for collaborative filtering of netnews. In *CSCW*, 1994.
- [41] P. Roy, A. Khan, and G. Alonso. Augmented sketch: Faster and more accurate stream processing. In *SIGMOD*, 2016.
- [42] P. Samarati. Protecting respondents identities in microdata release. *IEEE Trans. on Knowledge and Data Eng.*, 2001.
- [43] L. Sweeney. k-anonymity: A model for protecting privacy. *International J. of Uncertainty, Fuzziness and Knowledge-Based Sys.*, 2002.
- [44] C. J. van Rijsbergen. *Information retrieval*. Butterworth, 1979.
- [45] K. Vu, R. Zheng, and J. Gao. Efficient algorithms for k-anonymous location privacy in participatory sensing. In *INFOCOM*, 2012.
- [46] J. Yang and J. Leskovec. Defining and evaluating network communities based on ground-truth. *CoRR*, abs/1205.6233, 2012.