# LSF-Join: Locality Sensitive Filtering for Distributed All-Pairs Set Similarity Under Skew

Cyrus Rashtchian UCSD

crashtchian@eng.ucsd.edu

Aneesh Sharma Google

aneesh@google.com

David P. Woodruff CMU

dwoodruf@cs.cmu.edu

March 9, 2020

#### Abstract

All-pairs set similarity is a widely used data mining task, even for large and high-dimensional datasets. Traditionally, similarity search has focused on discovering very similar pairs, for which a variety of efficient algorithms are known. However, recent work highlights the importance of finding pairs of sets with relatively small intersection sizes. For example, in a recommender system, two users may be alike even though their interests only overlap on a small percentage of items. In such systems, some dimensions are often highly skewed because they are very popular. Together these two properties render previous approaches infeasible for large input sizes. To address this problem, we present a new distributed algorithm, LSF-Join, for approximate all-pairs set similarity. The core of our algorithm is a randomized selection procedure based on Locality Sensitive Filtering. Our method deviates from prior approximate algorithms, which are based on Locality Sensitive Hashing. Theoretically, we show that LSF-Join efficiently finds most close pairs, even for small similarity thresholds and for skewed input sets. We prove guarantees on the communication, work, and maximum load of LSF-Join, and we also experimentally demonstrate its accuracy on multiple graphs.

### 1 Introduction

Similarity search is a widely used primitive in data mining applications, and all-pairs similarity in particular is a common data mining operation [1, 7, 21, 34]. Motivated by recommender systems and social networks, we design algorithms for computing all-pairs set similarity (a.k.a., a set similarity join). In particular, we consider the similarity of nodes in terms of a bipartite graph. We wish to determine similar pairs of nodes from one side of the graph. For each node v on the right, we consider its neighborhood  $\Gamma(v)$  on the left. Equivalently, we can think of  $\Gamma(v)$  as a set of the neighbors of v in the graph. Using this representation, many graph-based similarity problems can be formulated as finding pairs of nodes with significantly overlapping neighborhoods. We focus on the cosine similarity between pairs  $\Gamma(v)$  and  $\Gamma(u)$  represented as high-dimensional vectors.

Although set similarity search has received a lot of attention in the literature, there are three aspects of modern systems that have not been adequately addressed yet. Concretely, we aim to develop algorithms that come with provable guarantees and that handle the following three criteria:

- 1. **Distributed and Scalable.** The algorithm should work well in a distributed environment like MapReduce, and should scale to large graphs using a large number of processors.
- 2. Low Similarity. The algorithm should output most pairs of sets with relatively low normalized set similarity, such as a setting of cosine similarity  $\tau$  taking values  $0.1 \le \tau \le 0.5$ .
- 3. **Extreme Skew.** The algorithm should provably work well even when the dimensions (degrees on the left) are highly irregular and skewed.

The motivation for these criteria comes from recommender systems and social networks. For the first criteria, we consider graphs with a large number of vertices. For the second, we wish to find pairs of nodes that are semantically similar without having a large cosine value. This situation is common in collaborative filtering and user similarity [27], where two users may be alike even though they overlap on a small number of items (e.g., songs, movies, or citations). Figure 1 depicts the close pair histogram of a real graph, where most similar pairs have low cosine similarity. For the third criteria, skewness has come to recent attention as an important property [5, 23, 36], and it can be thought of as power-law type behavior for degrees on the left. In contrast, most other prior work assumes that the graph has uniformly small degrees on the left [22, 27, 28]. This smoothness assumption is reasonable in settings when the graph is curated by manual actions (e.g., Twitter follow graph). However, this is too restrictive in some settings, such as a graph of documents and entities, where entities can legitimately have high degrees, and throwing away these entities may remove a substantial source of information. Another illustration of this phenomenon can be observed even on human-curated graphs, e.g., the Twitter follow graph, where computing similarities among consumers (instead of producers, as in [27]) runs into a similar issue.

Previous work fails to handle all three of the above criteria. When finding low similarity items (e.g., cosine similarity < 0.5), standard techniques like Locality-Sensitive Hashing [16, 29, 36] are no longer effective (because the number of hashing iterations is too large). Recently, there have been several proposals for addressing this, and the closest one to ours is the wedge-sampling algorithm from [27]. However, the approach in [27] has one severe shortcoming: it requires that each dimension has a relatively low frequency (i.e., the bipartite graph has small left degrees).

In this work, we address this gap by presenting a new distributed algorithm LSF-Join for approximate all-pairs similarity that can scale to large graphs with high skewness. As a main contribution, we provide theoretical guarantees on our algorithm, showing that it achieves very high accuracy. We also provide guarantees on the communication, work, and maximum load in a distributed environment with a very large number of processors.

Our approach uses Locality Sensitive Filtering (LSF) [11]. This is a variant of the ideas used for Locality Sensitive Hashing (LSH). The main difference between LSF and LSH is that the LSF constructs a single group of surviving elements based on a hash function (for each iteration). In contrast, LSH constructs a whole hash table, each time, for a large number of iterations. While the hashing and sampling ideas are similar, the benefit of LSF is in its computation and communication costs. Specifically, our LSF scheme will have the property that if an element v survives in k' out of k total hash functions, then the computation scales with k' and not k. For low similarity elements, k' is usually substantially smaller than k, resulting in a lower overall cost (for example k' will be

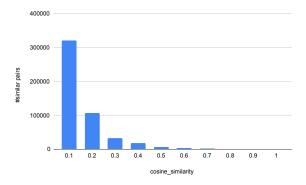


Figure 1: Histogram of the similar pairs at varying cosine similarity thresholds  $\tau$  for a citation network. The majority of pairs are concentrated at cosine similarity  $\tau \approx 0.1$ .

sublinear, while k is linear, in the input size). We also provide an efficient way to execute this filtering step on a per-node basis.

Our LSF procedure can also be a viewed as a pre-processing step before applying any all-pairs similarity algorithm (even one needing a smaller problem size and a graph without skew). The reason is that the survival procedure outputs a number of smaller subsets of the original dataset, each with a different, smaller set of dimensions, along with a guarantee that no dimension has a high degree. The procedure also ensures that similar pairs are preserved with high probability. Then, after performing this filtering, we may use other steps to improve the computation time. For example, applying a hashing technique may reduce the effective dimensionality without affecting the similarity structure.

#### Problem Set-up

The input consists of a bipartite graph G with a set of M vertices on the left and N vertices on the right. We denote that graph as G = (U, V, E), and we refer to U as the set of *dimensions*, and to V as the set of *nodes*. Given a parameter  $\tau > 0$ , we want to output all similar pairs of nodes (v, v') from V such that

$$\frac{|\Gamma(v) \cap \Gamma(v')|}{\sqrt{|\Gamma(v)| \cdot |\Gamma(v')|}} \ge \tau.$$

This problem also encapsulates other objectives, such as finding top-k results per node. Note that we could equivalently identify each node v with the set of its neighbors  $\Gamma(v) \subseteq U$ , and hence, this problem is the same as the set similarity join problem with input  $\{\Gamma(v) \mid v \in V\}$  and threshold  $\tau$  for cosine similarity. We describe our algorithm in a MapReduce-like framework, and we analyze it in the massively parallel computation model [8, 19], which captures the theoretical properties of MapReduce-inspired models (e.g., [26, 18]). We have p processors, in a shared-nothing distributed environment. The input data starts arbitrarily partitioned among the processors. Associated to each node v on the right is a vector  $\Gamma(v) \in \{0,1\}^M$  which is an indicator vector for the  $|\Gamma(v)|$  neighbors of v on the left. We would like to achieve the twin properties of load-balanced servers and low communication cost.

#### **Our Contributions**

The main contribution of our work is a new randomized, distributed algorithm, LSF-Join, which provably finds almost all pairs of sets with cosine similarity above a given threshold  $\tau$ . Our algorithm will satisfy all three of the criteria mentioned above (scalability, low similarity, and skewness). A key component of LSF-Join is a new randomized LSF scheme, which we call the survival procedure. The goal of this procedure is to find subsets of the dataset that are likely to contain similar pairs. In other words, it acts as a filtering step. Our LSF procedure comes with many favorable empirical and theoretical properties. First, we can execute it in nearly-linear time, which allows it to scale to very large datasets. Second, we exhibit an efficient way to implement it in a distributed setting with a large number of processors, using only a single round of communication for the whole LSF-Join algorithm. Third, the survival procedure leads to sub-quadratic local work, even when the dimensions are highly skewed and the similarity threshold is relatively low. To achieve these properties, we demonstrate how to implement the filtering using efficient, pairwise independent hash functions, and we show that even in this setting, the algorithm has good provable guarantees on the accuracy and running time. We also present a number of theoretical optimizations that better illuminate the behavior of the algorithm on datasets with different structural properties. Finally, we empirically validate our results by testing LSF-Join on multiple graphs.

#### Related Work

Many filtering-based similarity join algorithms provide exact algorithms and rely on heuristics to improve the running time [4, 6, 15, 22, 30, 32, 33, 34]. We primarily review prior work that is relevant to our setting and provides theoretical guarantees.

One related work uses LSF for set similarity search and join on skewed data [23]. Their *data* dependent method leads to a sequential algorithm based on the frequency of dimensions, improving a prior LSF-based algorithm [11]. Unfortunately, it seems impossible to adapt their method to the one-round distributed setting. Another relevant result is the wedge-sampling approach in [27]. They provide a distributed algorithm for low-similarity joins on large graphs. However, their algorithm assumes that the dataset is *not* skewed.

In the massively-parallel computation model [8, 9], multi-round algorithms have been developed that build off of LSH for approximate similarity joins, achieving output-optimal guarantees on the maximum load [17, 24]. However, it can be prohibitively expensive to use multiple rounds in modern shared-nothing clusters with a huge number of processors. In particular, the previous work achieves good guarantees only when the number of nodes N and number of processors p satisfy  $N \geq p^{1+c}$  for a constant c > 0. We focus on one-round algorithms, and we allow the possibility of  $p = \Theta(N)$ , which may be common in very large computing environments. Algorithms using LSH work well when  $\tau$  is large enough, such as  $0.6 \leq \tau < 1.0$ . However, for smaller  $\tau$ , LSH-based distributed algorithms require too much computation and/or communication due to the large number of repetitions [12, 31, 27, 35]. Prior work has also studied finding extremely close pairs [2, 1, 10] or finding pairs of sets with constant-size intersection [14]. These results do not apply to our setting because we aim to find pairs of large-cardinality sets with cosine similarity  $\tau$ 

in the range  $0.1 \le \tau \le 0.5$ , and we allow for the intersection size to be large in magnitude.

Finally, there are also conditional lower bounds showing that provably sub-quadratic time algorithms for all pairs set similarity (even approximate) may not exist in general [3, 25].

# 2 The LSF-Join Algorithm

We start with a high-level overview of our set similarity join algorithm, LSF-Join, which is based on a novel and effective LSF scheme. Let G = (U, V, E) be the input graph with |U| = M dimensions on the left, and |V| = N nodes on the right. For convenience, we refer to the vertices V and their indices [N] interchangeably, where we use [N] to denote the set  $\{1, 2, \ldots, N\}$ .

The LSF-Join algorithm uses k independent repetitions of our filtering scheme (where  $k \approx N$  achieves the best tradeoff). In the i-th repetition we create a set  $S_i \subseteq [N]$  of survivors of the set [N] of vertices on the right. We will define the LSF procedure shortly, which will determine the subsets  $\{S_i\}_{i=1}^k$  in a data-independent fashion. During the communication phase, the survival sets will be distributed in their entirety across the processors. In particular, if there are p processors, then each processor will handle roughly k/p different repetitions. During the local computation, the processors will locally compute all similar pairs in  $S_i$  for  $i \in [k]$  and output these pairs in aggregate (in a distributed fashion). As part of the theoretical analysis, we show that the size of each  $S_i$  is concentrated around its mean, and therefore, our algorithm has balanced load across the processors. To achieve high recall of similar pairs, we will need to execute the LSF-Join algorithm  $O(\log N)$  times independently, so that the failure probability will be polynomially small. Fortunately, this only increases the communication and computation by a  $O(\log N)$  factor. We execute the iterations in parallel, and LSF-Join requires only one round of communication.

# 2.1 Constructing the Survival Sets $S_i$

We now describe our LSF scheme, which boils down to describing how to construct the  $S_i$  survival sets. We have two main parameters of interest:  $\alpha \in (0, 1/2]$  denotes the survival probability of a single dimension (on the left), and k denotes the number of repetitions. The simplest way to describe our LSF survival procedure goes via uniform random sampling. We refer to this straightforward scheme as the *Naive-Filter* method, and we describe it first. Then, we explain how to improve this method by using a pairwise independent filtering scheme, which will be much more efficient in practice. We refer to the improved LSF scheme as the *Fast-Filter* method. Later, we also show that Fast-Filter enjoys many of the same theoretical guarantees of Naive-Filter, with much lower computational cost.

Naive-Filter. For the naive version of our filtering scheme, consider a repetition number  $i \in [k]$ . We choose a uniformly random set  $U_i \subseteq U$  of vertices on left by choosing each node  $u \in U$  to be in  $U_i$  with probability  $\alpha$  independently. Then, we filter vertices v on the right depending on whether their neighborhood is completely contained in  $U_i$  or not (that is, whether  $\Gamma(v) \subseteq U_i$  or not). The i-th survival set  $S_i$  will be the set of vertices  $v \in V$  such that  $\Gamma(v) \subseteq U_i$ . We repeat this process independently for each i = 1, 2, ..., k, to derive k filtered sets of vertices  $S_1, ..., S_k$ . Notice that

for each i, the probability that v survives in  $S_i$  is exactly  $\alpha^{|\Gamma(v)|}$ , where  $|\Gamma(v)|$  is the number of neighbors of v on the left.

The intuition behind using this filtering method for set similarity search is that similar pairs are relatively likely to survive in the same set. Indeed, the chance that both v and v' survive in  $S_i$  is equal to  $\alpha^{|\Gamma(v)\cup\Gamma(v')|}$ . When the cosine similarity is large, we must have that  $|\Gamma(v)\cap\Gamma(v')|$  is large and also that  $|\Gamma(v)\cup\Gamma(v')|$  is much smaller than  $|\Gamma(v)|+|\Gamma(v')|$ . In other words, v and v' are more likely to survive together if they are similar, and less likely if they are very different. For example, consider the case where  $d=|\Gamma(v)|=|\Gamma(v')|$  is a large constant. Then, pairs with cosine similarity at least  $\tau$  will survive together with probability  $\alpha^{(2-\tau)d}$ . At the other extreme, disjoint pairs only survive together with probability  $\alpha^{2d}$ .

The main drawback of the Naive-Filter method is that it takes too much time to determine all indices i such that  $v \in S_i$ . Consider the set of v's neighbors  $\Gamma(v)$ . We need to determine whether  $\Gamma(v) \subseteq U_i$  for every  $i \in [k]$ . Hence, it requires at least  $O(\alpha|\Gamma(v)|k)$  work to compute the indices where v survives, that is, the set  $\{i : v \in S_i\}$ . We will need to set  $k \gg N$ , and hence, the work of Naive-Filter is linear in N or worse for each node v. To improve upon this, our Fast-Filter method will have work proportional to  $|\{i : v \in S_i\}|$ , and we show that this is often considerably smaller than k.

### 2.2 The Fast-Filter Method

The key idea behind our fast filtering method is to develop a pairwise independent filtering scheme that approximates the uniform sampling of the survival sets. We then devise a way to efficiently compute the survival sets on a per-node basis, by using fast matrix operations. More precisely, for each node v on the right, Fast-Filter will determine the indices  $I_v \subseteq [k]$  of survival sets in which v survives (that is, we have  $I_v = \{i : v \in S_i\}$ ). We develop a way to compute  $I_v$  independent for each vertex v by using Gaussian elimination on binary matrices. The Fast-Filter method only requires a small amount of shared randomness between the processors.

To describe the Fast-Filter method, it will be convenient to assume that  $\log_2(1/\alpha)$  and  $\log_2 k$  are both integers. We now explain the pairwise independent filtering scheme. For each node  $u \in U$  on the left, we sample a random  $\log_2(1/\alpha) \times \log_2 k$  binary matrix  $A'_u$  and a  $\log_2 1/\alpha$ -length bit-string  $b'_u$ . We identify each of the k repetitions  $i \in [k]$  with binary vectors in the  $\log_2(k)$ -dimensional vector space over GF(2), the finite field with two elements. In other words, we use the binary representation of i to associate i with a length  $\log_2 k$  bit-string, and we perform matrix and vector operations modulo two. We abuse notation and use i for both the integer and the bit-string, where context will distinguish the two.

To determine whether a node  $v \in [N]$  survives in  $S_i$ , we perform the following operation. We first stack the matrices  $A'_u$  on top of each other for each of v's neighbors  $u \in \Gamma(v)$ . This forms a  $|\Gamma(v)| \cdot \log_2(1/\alpha) \times \log_2 k$  matrix  $A^v$ . We also stack the vectors  $b'_u$  on top of each other, forming a length  $|\Gamma(v)| \cdot \log_2(1/\alpha)$  bit-string  $b^v$ . Finally, we define  $S_i$  by setting  $v \in S_i$  if and only if  $A^vi + b^v = 0$ , where 0 denotes the all-zeros vector. We say that v survives the i-th repetition if  $A^vi + b^v = 0$ . Then  $I_v = \{i : v \in S_i\}$  is the set of indices  $I_v \subseteq [k]$  in which v survives.

In a one-round distributed setting, the processors can effectively pre-compute the submatrices

### Algorithm 1 Efficient LSF for a Single Node

- 1: **function** Fast-Filter( $G, v, \alpha, k$ )
- 2: Compute  $A^v$  and  $b^v$  using the shared random seed
- 3: Determine the solution space of  $A^{v}i + b^{v} = 0$
- 4: Let  $I_v \leftarrow \{i : A^v i + b^v = 0\}$
- 6: end function

### **Algorithm 2** Approximate Cosine Similarity Join

- 1: Repeat the following procedure  $O(\log N)$  times in parallel:
- 2: **function** LSF-Join(  $G = (U, V, E), \tau, \alpha, k$  )
- 3: **For** each vertex  $v \in V$  do in parallel:
- 4: FAST-FILTER $(G, v, \alpha, k)$  to determine sets containing v
- 5: Partition the sets  $S_1, \ldots, S_k$  across processors
- 6: Locally compute all pairs in each  $S_i$  with similarity  $\geq \tau$
- 7: Output all close pairs in a distributed fashion
- 8: end function

 $A'_u$  and the subvectors  $b'_u$  using a shared seed. In particular, these may be computed on the fly, as opposed to stored up front, by using a shared random seed and by using an efficient hash function to compute the elements of  $A'_u$  and  $b'_u$  only when processing v such that  $u \in \Gamma(v)$ . By doing so, the processors will use the same values of  $A'_u$  and  $b'_u$  as one another, leading to consistent survival sets, without incurring any extra rounds of communication.

To gain intuition about this filtering procedure, let  $d = |\Gamma(v)|$  denote the number of v's neighbors. Node v will survive in  $S_i$  if i satisfies  $A^v i + b^v = 0$ . This consists of  $d \cdot \log_2(1/\alpha)$  linear equations that i must satisfy. As the matrix  $A^v$  and the vector  $b^v$  are chosen uniformly at random, it is easy to check that v survives in  $S_i$  with probability  $\alpha^{|\Gamma(v)|} = \alpha^d$ , and hence, their expected sizes satisfy

$$\mathbb{E}[|S_i|] = \alpha^d N$$
 and  $\mathbb{E}[|I_v|] = \alpha^d k$ 

over a random  $A^v$  and  $b^v$ .

Theoretically, the main appeal of Fast-Filter is that it is pairwise independent in the following sense. For any two distinct repetitions i and i', the bit-strings for i and i' differ in at least one bit. Therefore, we see that  $A^vi + b^v = 0$  is satisfied or not independently of  $A^vi' + b^v = 0$ , over the random choice of  $A^v$  and  $b^v$ . While this is only true for pairs of repetitions, this level of independence will suffice for our theoretical analysis. Furthermore, we show that we can determine the survival sets containing v in time proportional to the number  $|I_v|$  of such sets, which is often much less than the total number k of possible sets.

We now explain how to efficiently compute the survival sets on a per-node basis. For a fixed node  $v \in [N]$ , the Fast-Filter method determines the repetitions i that v survives in, or in other words, the set  $I_v = \{i : v \in S_i\}$ . This is equivalent to finding all length  $\log_2(k)$  bit-strings i that are solutions to  $A^v i + b^v = 0$ . The processor can form  $A^v$  and  $b^v$  in O(d) time, where  $d = |\Gamma(v)|$ ,

assuming the unit cost RAM model on words of  $O(\log_2(N))$  bits. Then, we can use Gaussian elimination over bit-strings to very quickly find all  $i \in [k]$  that satisfy  $A^v i + b^v = 0$ . To understand the complexity of this, first note that  $A^v$  has  $\log_2 k$  columns. Moreover, without loss of generality, we see that  $A^v$  has at most  $\log_2 k$  rows, as otherwise there exists no solution. Therefore, Gaussian elimination takes  $O(\log^3 k)$  time to write  $A^v$  in upper triangular form (and correspondingly rewrite  $b^v$ ) so that all solutions to  $A^u i = b^u$  can be enumerated in time proportional to the number of solutions to this equation. The expected total work is

$$O(N \log^3 k + \alpha^d k N).$$

This can be parallelized for each node v independently.

We prove guarantees about Fast-Filter in Theorem 2. The pseudo-code for Fast-Filter appears as Algorithm 1. The main difference between the two filtering methods is how the random survival sets are chosen. For the sake of this discussion, we set k = N, which is reasonable in practice, and we continue to let  $d = |\Gamma(v)|$ . In the Fast-Filter method, we use a random linear map over GF(2) with enough independent randomness to decide for each repetition, whether or not a node survives not. By using Gaussian elimination, we are able to compute  $I_v = \{i : v \in S_i\}$  in time proportional to  $|I_v| \leq N$ . In particular, the amount of work for v is  $O(\log^3 N + \alpha^d N)$  in expectation, because  $\mathbb{E}[|I_v|] = \alpha^d N$  when k = N.

The pseudo-code for LSF-Join appears as Algorithm 2. We assume that the vertices v start partitioned arbitrarily across p processors. For each vertex v in parallel, we use Fast-Filter determine the indices  $I_v$  of the sets in which v survives. As detailed above, we can do so consistently by using a shared random seed for Fast-Filter. During the communication phase, we randomly distribute the sets  $S_1, \ldots, S_k$  across p processors, so that each processor handles k/p sets in expectation. Then, during local computation, we compare all pairs in  $S_i$  for each  $i \in [k]$  in parallel. We use  $O(\log N)$  independent iterations of the algorithm in parallel to find all close pairs with high probability (e.g., recall close to one). Finally, we output all pairs with cosine similarity at least  $\tau$  in a distributed fashion.

One way of processing each  $S_i$  set is to compare all pairs in this set. Specifically, for all pairs of nodes  $v, v' \in S_i$ , explicitly compute  $|\Gamma(v) \cap \Gamma(v')|$  and check if it is at least  $\tau \sqrt{|\Gamma(v)| \cdot |\Gamma(v')|}$ . One can assume the lists  $\Gamma(v)$  and  $\Gamma(v')$  are sorted arrays of d' integers, where  $d' = \max\{|\Gamma(v)|, |\Gamma(v')|\}$ . Thus, one can compute  $|\Gamma(u) \cap \Gamma(v)|$  by merging these sorted lists in O(d') time, assuming words of length  $O(\log_2(N))$  can be manipulated in constant time in the unit cost RAM model.

Letting  $d_i$  be the maximum of  $|\Gamma(v)|$  over  $v \in S_i$ , the time to locally compare all pairs in set  $S_i$  is  $O(|S_i|^2 d_i)$ . We can also bound the average amount of work across p processors to handle all sets  $S_1, \ldots, S_k$ . This can be bounded by

$$O\left(\sum_{i=1}^{k} |S_i|^2 \cdot d_i \cdot \frac{k}{p}\right).$$

We call this the *brute-force all-pairs* algorithm.

### 2.2.1 Setting the Parameters

Let  $\bar{d}$  denote the average degree on the right in the input graph. Ideally, these parameters should satisfy

$$\alpha^{(2-\tau)\bar{d}} \cdot k = 2,\tag{1}$$

or in other words,  $\alpha = (2/k)^{1/((2-\tau)\bar{d})}$ , where 2 could be replaced with a larger constant for improved recall. If it is possible to approximately satisfy (1) with  $\log_2(1/\alpha)$  being an integer, then running  $O(\log N)$  independent iterations of the algorithm with these parameters will work very well. For example, this is the case when  $(1/2)^{\bar{d}} = 1/N^c$  for constant  $c \approx 1$ . However, for large average degree  $\bar{d}$ , the parameter  $\alpha$  may exceed 1/2. To approximate  $\alpha > 1/2$ , we can subsample the matrices  $A^v$  and vectors  $b^v$  to increase the effective collision probability. More precisely, consider  $d = |\Gamma(v)|$ . If we wish to survive in a repetition with probability  $\alpha^d$ , then we can solve for  $d^*$  in the equality  $\alpha^d = (1/2)^{d^*}$ , and we subsample the d rows in  $A^v$  and  $b^v$  down to  $d^*$ . This effectively constructs survival sets  $S_i$  as in Naive-Filter with  $\alpha$  probability of each neighbor surviving. In the theoretical results, we will assume that  $\alpha$  and k satisfy (1). In the experiments, we either set  $\alpha$  to be 1/2, or we use the matrix subsampling approach; we also vary the number of independent iterations to improve recall (where we use  $\beta$  to denote the number of iterations).

# 3 Theoretical Guarantees

We assume on the graph G = (U, V, E) is right-regular with nodes in V having degree d for simplicity. In practice, we can repeat the algorithm for different small ranges of d. First, notice that

$$\Pr[v \in S_i] = \Pr[A^v i + b^v = 0] = \frac{1}{2^{d \log_2 1/\alpha}} = \alpha^d$$
 (2)

Now consider two nodes  $u, v \in [N]$ . Then both u and v are in  $S_i$  if and only if the following event occurs. Let  $A^{u,v}$  be the matrix obtained by stacking  $A^u$  on top of  $A^v$ , and  $b^{u,v}$  be the vector obtained by stacking  $b^u$  on top of  $b^v$ . Note that for each  $w \in \Gamma(u) \cap \Gamma(v)$ , the rows of  $A_w$  occur twice in  $A^{u,v}$  and the coordinates of  $b_w$  occur twice in  $b^{u,v}$ . Thus, it suffices to retain only one copy of  $A_w$  and  $b_w$  in  $A^{u,v}$  for each  $w \in \Gamma(u) \cap \Gamma(v)$ , and by doing so we reduce the number of rows of  $A^{u,v}$  and entries of  $b^{u,v}$  to at most  $|\Gamma(u) \cup \Gamma(v)| \cdot d \log_2 1/\alpha$ . Consequently,

$$\Pr[u \in S_i \text{ and } v \in S_i] = \Pr[A^{u,v}i + b^{u,v} = 0] = \alpha^{|\Gamma(u) \cup \Gamma(v)|}$$
(3)

Notice that on one extreme if  $\Gamma(u)$  and  $\Gamma(v)$  are disjoint, then (3) evaluates to  $\alpha^{2d}$ . On the other hand, if  $|\Gamma(u) \cap \Gamma(v)| \ge \tau \cdot d$ , then  $|\Gamma(u) \cup \Gamma(v)| \le (2-\tau)d$ , and then (3) evaluates to  $\alpha^{(2-\tau)d}$ .

The discrepancy in (2) and (3) is exactly what we exploit in our LSF scheme; namely, we use the fact that similar pairs are more likely to survive together in a repetition than dissimilar pairs. We first justify the setting of  $\alpha$  in (1).

**Lemma 1.** Let u, v be such that  $|\Gamma(u) \cap \Gamma(v)| \ge \tau d$ . The expected number of repetitions i for which both  $u \in S_i$  and  $v \in S_i$  is at least 2.

*Proof.* As shown in (3), the probability both u and v survive in a single repetition is  $\alpha^{|\Gamma(u)\cup\Gamma(v)|} \geq \alpha^{(2-\tau)d}$ , and therefore the expected number of repetitions for which both  $u \in S_i$  and  $v \in S_i$  is at least  $k \cdot \alpha^{(2-\tau)d}$ , which by (1) is at least 2.

**Lemma 2.** The expected load per processor is  $\alpha^d Nk/p$ , and the expected total communication is  $\alpha^d kN$ .

Proof. There are k repetitions, each concerning one  $S_i$  survival set. Each node  $v \in [N]$  survives in  $S_i$  with probability  $\alpha^d$  independently. The expected size of  $S_i$  is  $\mathbf{E}|S_i| = \alpha^d N$ . Each processor handles k/N repetitions, leading to  $\alpha^d N k/p$  expected load. The total communication is  $\sum_{i=1}^k |S_i|$ , which has expectation  $\alpha^d k N$ .

**Lemma 3.** Using brute-force all-pairs locally, the expected work per machine is  $(\alpha^d N)^2 k/p$ .

*Proof.* Each repetition has expected size  $\alpha^d N$ , leading to work  $\alpha^{2d} N^2$ . Each processor handles k/p repetitions, implying  $\alpha^{2d} N^2 k/p$  work per processor in expectation.

Combining the lemmas and plugging in  $\alpha$  gives us the following.

**Theorem 1.** Setting  $\alpha = (2/k)^{1/((2-\tau)d)}$ , the survival procedure has total communication is

$$O(Nk^{1-1/(2-\tau)})$$

and local work

$$O(N^2k^{1-2/(2-\tau)}/p)$$

in expectation.

As an example, we compare to hash-join when p=N, which has total communication  $N^{3/2}$  and local work N. We set  $k=N^{\frac{2-\tau}{2-2\tau}}$ , and by Theorem 1, the expected total communication is  $Nk^{-\tau/(2-2\tau)}=N^{3/2}$ . The local work per processor is  $Nk^{-\tau/(2-\tau)}=N^{1-\frac{\tau}{2-2\tau}}$ . Since  $\tau>0$ , the work is always sublinear, thus improving over hash-join while using the same amount of total communication. As we will see in the theorem below, it is crucial that we use the family of pairwise independent hash functions above for generating our randomness.

**Theorem 2.** The expected total time the nodes in [N] need to generate the  $S_i$  is

$$O(N\log^3 k + \alpha^d k N + |E|),$$

and the expected total time and communication that the nodes in [N] need to send the sets  $\Gamma(v)$  for each  $v \in S_i$  for each i is

$$O(N\log^3 k + \alpha^d k N \cdot d\log N + |E|).$$

Proof. Each node  $u \in [N]$  needs to figure out the repetitions i that it survives in. It can form  $A^u$  and  $b^u$  in O(d) time assuming the unit cost RAM model on word of  $O(\log_2(N))$  bits. Note u then needs to figure out which  $i \in [N]$  satisfy  $A^u \cdot i + b^u = 0$ . To do so, in can just solve this equation using Gaussian elimination. Note that  $A^u$  has at most  $\log_2 k$  rows, and has  $\log_2 k$  columns. Therefore Gaussian elimination takes at most  $O(\log^3 k)$  time to write  $A^u$  in upper triangular form and corresponding  $b^u$  so that all solutions to the equation  $A^u x = b^u$  can be enumerated in time proportional to the number of solutions to this equation. Thus, the expected time per processor is  $O(\log^3 N + \alpha^d N)$ , where we have used (2) to bound the expected number of repetitions that u survives in by  $k \cdot \alpha^d$ . Thus, the total expected time to form all of the  $S_i$ , for i = 1, 2, ..., k, is  $O(N \log^3 k + \alpha^d k N)$ . Note that  $O(\alpha^d k N \cdot d \log N)$  is the total expected amount of communication.

While correct in expectation, since the randomness uses across the repetitions is not independent, namely, we use the same matrices  $A_w$  and vectors  $b_w$  for each node  $w \in [M]$ , it is important to show that the variance of the number of repetitions i for which both  $u \in S_i$  and  $v \in S_i$  is small. This enables one to show the probability there is at least one repetition i for which both u and v survive is a large enough constant, which can be amplified to any larger constant by independently repeating a constant number of times.

**Lemma 4.** Let u, v be such that  $|\Gamma(u) \cap \Gamma(v)| \geq \tau d$ . With probability at least 1/2, there is a repetition i with both  $u \in S_i$  and  $v \in S_i$ .

Proof. Let  $X_i$  be an indicator random variable which is 1 if u and v survive the i-th repetition, and is 0 otherwise. Let  $X = \sum_{i=1}^k X_i$  be the number of repetitions for which both u and v survive. By Lemma 1,  $\mathbf{E}[X] \geq 2$ . It is well-known that the hash function family  $f(x) = Ax + b \mod 2$ , where A and b range over all possible binary matrices and vectors, respectively, is a pairwise independent family. It follows that  $X_1, X_2, \ldots, X_k$  are pairwise independent random variables, and consequently  $\mathbf{Var}[X] = \sum_{i=1}^k \mathbf{Var}[X_i]$ . As  $X_i \in \{0,1\}$ , we have  $\mathbf{Var}[X_i] \leq \mathbf{E}[X_i]$ , and hence,  $\mathbf{Var}[X] \leq \mathbf{E}[X]$ . By Chebyshev's inequality,

$$\Pr[X=0] \leq \Pr\left[|X - \mathbf{E}[X]| \geq \mathbf{E}[X]\right] \leq \frac{\mathbf{Var}[X]}{(\mathbf{E}[X])^2} \leq \frac{1}{\mathbf{E}[X]} \leq \frac{1}{2}.$$

Efficiently Amplifying Recall. At this point, we have shown that one iteration of LSF-Join will find a constant fraction of close pairs. To amplify the recall, we run  $\beta = O(\log N)$  copies of LSF-Join in parallel. We emphasize that this is a more efficient way to achieve a high probability result, better than simply increasing the number of repetitions k in a single LSF-Join execution. Intuitively, this is because the repetitions are only guaranteed to be pairwise independent. Theoretically,  $O(\log(1/\delta))$  independent copies leads to a failure probability of  $1 - \delta$  by a Chernoff bound. But, if we only increased the number of repetitions, then by Chebyshev's inequality, we would need to use  $O(k/\delta)$  repetitions for the same success probability  $1 - \delta$ . The latter requires  $O(1/\delta)$  times the amount of communication/computation, while the former is only a  $O(\log(1/\delta))$  factor. Setting  $\delta = 1/N^3$  leads to a failure probability of 1 - 1/N after taking a union bound over the  $O(N^2)$  possible pairs.

# 4 Optimizations

In this section, we present several extensions of the LSF-Join algorithm and analysis, such as considering the number of close pairs, using hashing to reduce dimensionality, combining LSF-Join with hash-join, and lowering the communication cost when the similarity graph is a matching.

### 4.1 Processing Time as a Function of the Profile

While Theorem 1 gives us a worst-case tradeoff between computation and communication, we can better understand this tradeoff by parameterizing the total amount of work of the servers by a data-dependent quantity  $\Phi$ , introduced below, which may give a better overall running time in certain cases.

Supposing that  $k \geq p$ , the processors receive multiple sets to process. We choose a random hash function  $H:[k] \to [p]$  so that processor j receives all sets  $S_i$  for which H(i) = j. When  $k \geq p$ , each processor handles k/p sets  $S_i$  in expectation.

The processor handling the set  $S_i$  receives  $S_i$  together with the neighborhood  $\Gamma(u)$  for each  $u \in S_i$ , and is responsible for outputting all pairs  $u, v \in S_i$  for which  $|\Gamma(u) \cap \Gamma(v)| \ge \tau d$ .

To bound the total amount of computation, we introduce a data-dependent quantity  $\Phi$ . Note that the  $S_i$  are independent and identically distributed, so we can fix a particular i. We define the profile  $\Phi$  of a dataset as follows:

$$\Phi = N \cdot \alpha^d + \sum_{u \neq v \in [N]} \alpha^{|\Gamma(u) \cup \Gamma(v)|}.$$

**Lemma 5.**  $\mathbf{E}[|S_i|] = N \cdot \alpha^d$  and  $\mathbf{E}[|S_i|^2] \leq \Phi$ .

*Proof.* Let  $|S_i| = \sum_u X_u$ , where  $X_u$  is an indicator that node u survives repetition i. Then  $|S_i| = \sum_u X_u$ , and so  $\mathbf{E}[|S_i|] = N \cdot \alpha^d$  by (2). For the second moment,

$$\mathbf{E}[|S_i|^2] = \sum_{u,v} \mathbf{E}[X_u X_v] \le \sum_u \mathbf{E}[X_u^2] + \sum_{u \ne v} \mathbf{E}[X_u X_v].$$

Plugging  $\alpha^{|\Gamma(u)\cup\Gamma(v)|}$  from (3) for  $\mathbf{E}[X_uX_v]$  and using the definition of  $\Phi$  proves the lemma.

We are interested in bounding the overall time for all nodes in [p] to process the sets  $S_i$ .

**Theorem 3.** The total work of the nodes in [p] to process the sets  $S_1, \ldots, S_k$ , assuming that we use the brute-force all-pairs algorithm is  $O(dk\Phi)$ . The average work per processor is  $O(\frac{dk}{p}\Phi)$ .

*Proof.* After receiving the  $S_i$ , the total time for all processors to execute their *brute-force all-pairs* algorithm is  $O(dk\Phi)$ , which allows for outputting the similar pairs. The theorem follows.

### 4.2 Hashing to Speed Up Processing

Recall that the processor responsible for finding all similar pairs in  $S_i$  receives the set  $\Gamma(u)$  of neighbors of each node  $u \in S_i$ . In the case when the neighborhoods are all of comparable size,

sat size d, we can think of  $\Gamma(u)$  as a vector  $\chi_u \in \{0,1\}^M$  with exactly d ones in it; here  $\chi_u$  is the characteristic vector of the neighbors of u. We can first hash the vector  $\chi_u$  down to  $s = d/(z\tau)$  dimensions, for a parameter z > 0. To do this, we use the CountMin map [13], which can be viewed as a random matrix  $S \in \{0,1\}^{s \times M}$  with a single non-zero per column, and this non-zero is chosen uniformly at random and independently for each of the M columns of S. We replace  $\chi_u$  with  $S \cdot \chi_u$ . If an entry of  $S \cdot \chi_u$  is larger than 1, we replace it with 1, and let the resulting vector be denoted  $\gamma_u$ , which is in  $\{0,1\}^s$ . Note that we can compute all of the  $\gamma_u$  for a given repetition i using  $O(|S_i|d)$  time, assuming arithmetic operations on  $O(\log M)$  bit words can be performed in constant time.

While  $\langle \chi_u, \chi_v \rangle = |\Gamma(u) \cap \Gamma(v)|$  for two nodes  $u, v \in [N]$ , it could be that  $\langle \gamma_u, \gamma_v \rangle \neq |\Gamma(u) \cap \Gamma(v)|$ . We now quantify this.

**Lemma 6.** For any two nodes  $u, v \in [N]$ , it holds that with probability at least 1 - 2/z,

$$|\langle \gamma_u, \gamma_v \rangle - \langle \chi_u, \chi_v \rangle| \le d\tau/2.$$

Proof. Note that  $\langle \gamma_u, \gamma_v \rangle \leq |\Gamma(u) \cap \Gamma(v)|$  since each node  $w \in \Gamma(u) \cap \Gamma(v)$  is hashed to a bucket by CountMin, which will be a coordinate that is set to 1 in both  $\gamma(u)$  and  $\gamma(v)$ . Also the probability that w hashes to a bucket containing a  $\hat{w} \in \Gamma(u) \cap \Gamma(v)$  with  $\hat{w} \neq w$  is at most  $d/(d/(z\tau)) = z\tau$ , and the expected number of w with this property is at most  $dz\tau$ . By a Markov bound, the number of such w is at most  $d\tau/2$  with probability at least 1 - 2/z, as desired.

By the previous lemma we can replace the original dimension-M vectors  $\chi_u$  with the potentially much smaller dimension- $d/(z\tau)$ -vectors  $\gamma_u$  with a small price in accuracy and success probability.

### 4.3 Combining LSF-Join with Hash-Join

The LSH-based approach of Hu et. al [17] suggests (in our framework) an alternate strategy of sub-partitioning the survival sets, using a hash-join to distribute the brute-force all-pairs algorithm. Here we analyze this combined approach and plot the tradeoffs. We show that this strategy does not provide any benefit in the communication vs. computation tradeoff, perhaps surprisingly.

The combined strategy, using p processors, starts by using  $k = p^c$  repetitions for a parameter c < 1, and this is followed by a hash join on each survival set. More precisely, we first construct k sets  $S_1, \ldots, S_k$  using the Fast-Filter survival procedure. Then, for each set  $S_j$ , we will process all pairs in  $S_j \times S_j$  using  $p/k = p^{1-c}$  machines. This can be implemented in one round, because all we need to do is estimate the size of each set  $S_j$  approximately, that is,  $|S_j| \approx \alpha^d N$ . Then, we can implement the hash-join in a distributed fashion.

We first review the guarantees of the standard hash-join.

**Lemma 7.** For N vectors and p machines, a hash-join has expected total communication  $N\sqrt{p}$  and expected  $N^2/p$  work per machine.

We use this bound to compute the communication and work, when using a hash-join to process each survival set.

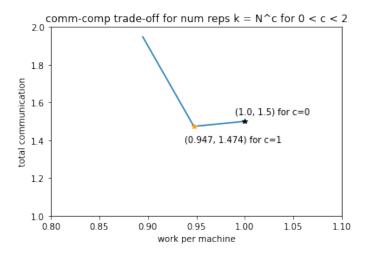


Figure 2: For  $\tau = 0.1$  and p = N, a comparison of LSF-Join and the combined approach of using hash-join to distribute the brute-force all-pairs. We plot the exponent of N for the different settings of  $k = N^c$  repetitions for 0 < c < 2.

**Theorem 4.** The combined approach has expected total communication  $N^{1+\frac{c(1-\tau)}{2-\tau}}p^{\frac{1-c}{2}}$  and expected  $N^2/p^{1+\frac{c\tau}{2-\tau}}$  work per processor.

*Proof.* When  $k=N^c$ , we have that  $\alpha^d=k^{\frac{-c}{2-\tau}}$ , and hence, we have  $|S_j|=\alpha^dN=Np^{-\frac{c}{2-\tau}}$  in expectation. We use Lemma 7 to analyze the hash-join for each of the  $p^c$  groups of  $p^{1-c}$  processors. Each group handles  $N'=Np^{-\frac{c}{2-\tau}}$  inputs, and therefore the communication of the group is  $N'\cdot p^{1/2-c/2}$ , which is  $N^{1-\frac{c}{2-\tau}}p^{\frac{1-c}{2}}$ . Multiplying by  $N^c$ , the exponent of N becomes

$$1 + c - \frac{c}{2 - \tau} = 1 + \frac{c(1 - \tau)}{2 - \tau},$$

which gives the claimed communication bound. For the per processor work, we have that this is the claimed bound:

$$(N')^2/p^{1-c} = N^2 p^{-1+c-\frac{2c}{2-\tau}} = N^2 p^{-1-\frac{c\tau}{2-\tau}}.$$

Figure 2 demonstrates that the combination approach is never better than the original LSF-Join approach. For a comparison, we consider p=N processors, and hence, the number of repetitions will be  $k=N^c$  for 0 < c < 2. Then, when  $c \ge 1$ , the survival procedure has expected total communication  $N^{1+\frac{c(1-\tau)}{2-\tau}}$ , and it has expected  $N^{1-\frac{c\tau}{2-\tau}}$  work per processor. And, when  $c \le 1$  we have that the combined approach has expected total communication  $N^{\frac{3}{2}-\frac{c\tau}{2(2-\tau)}}$ , and it has expected  $N^{1-\frac{c\tau}{2-\tau}}$  work per processor. Notice that c=1 corresponds to standard LSF-Join, and c=0 corresponds to using a hash-join on the whole dataset.

# 4.4 When the Similarity Graph is a Matching

Recall that to recover all close pairs with high probability, we need to iterate the LSF-Join algorithm  $O(\log N)$  times, because each time finds a constant fraction of close pairs. We exhibit an improvement using multiple communication steps when the similar pairs are structured. An important application of all-pairs similarity is constructing the similarity graph. In our setting, the *similarity graph* connects all pairs  $v, v' \in V$  such that their cosine similarity is at least  $\tau$ . The structure that we consider is when the similarity graph happens to be a *matching*, containing exactly N/2 disjoint pairs v, v' with similarity at least  $\tau$ .

The key idea is that each iteration decreases the number of input nodes by a constant fraction. We will remove these nodes (or at least one endpoint from each close pair) from consideration, and then repeat the procedure using the remaining nodes. We observe that this method can also be extended to near-matchings (e.g., small disjoint cliques). Similarly, our result is not specific to LSF-Join, and the technique would work for any LSF similarity join method.

We state our result using the r-th iterated log function  $\log^{(r)} N$ , where  $\log^{(1)} N = \log N$ , and  $\log^{(r)} N = \log(\log^{(r-1)} N)$  for  $r \ge 2$ . Then, we show:

**Theorem 5.** Using 2r-1 communication steps, we can find all but a negligible fraction of close pairs when the similarity graph is a matching. The total communication and computation is  $O(\log^{(r)} N)$  times the cost of one execution of LSF-Join.

*Proof.* For r = 1, we simply run LSF-Join  $O(\log N)$  times independently in a single communication step, where each time finds a constant fraction of close pairs. For  $2r - 1 \ge 3$  communication steps, we will use r rounds of LSF-Join, and we will remove all found pairs between subsequent rounds (each round will take two communication steps, except for the last, which takes one).

In the first round, we run LSF-Join  $T_r = O(\log^{(r)} N)$  times. Then, the expected number of pairs that are *not* found will be  $O(N/2^{T_r})$ , where  $2^{T_r} = \text{poly}(\log^{(r-1)} N)$ . In the next round, with r-1 rounds remaining, we will only consider the remaining pairs, and we will iterate LSF-Join  $T_{r-1}$  times. We repeat this process until no more rounds remain, and output the close pairs from all rounds.

We can implement each round of the above algorithm using at most two communication steps. We do so by marking the found pairs between rounds using a single extra communication step. More formally, the input pairs start partitioned across p processors. We denote the input partition as  $V = V_1 \cup \cdots \cup V_p$ . After finding some fraction of close pairs, processor i must be notified of which nodes in  $V_i$  are no longer active. Whenever processor j finds a close pair (v, v'), it sends the index of v to processor i such that  $v \in V_i$  (and similarly for  $v' \in V_{i'}$ ), where i is known to processor j because processor i must have sent v to processor j in LSF-Join. We reduce the total input set from V to V', where V' denotes the remaining nodes after removing the found pairs.

To analyze this procedure, notice that the dominant contribution to the total communication and computation is the first round. This is because the subsequent rounds have a geometrically decreasing number of input nodes. The first round uses  $T_r = O(\log^{(r)} N)$  iterations of LSF-Join, which shows that overall communication and computation is  $O(\log^{(r)} N)$  times the cost of one iteration.

### 4.5 Hashing to Improve Recall

Not only is hashing helpful in order to reduce the description size of the neighborhood sets, as described in Section 4.2, hashing can also be used to increase the number of similar pairs surviving a repetition, and thus the recall. Before, a node pair (u, v) survives a repetition with probability  $\alpha^{|\Gamma(u)\cup\Gamma(v)|}$ . Hashing can, however, make  $|\Gamma(u)\cup\Gamma(v)|$  smaller due to collisions. Suppose we hash the characteristic vector  $\chi_u \in \{0,1\}^M$  of the neighborhood of a node u down to d/C dimensions for some parameter  $C \geq 1$ , obtaining the vector  $\gamma_u \in \{0,1\}^{d/C}$ , as in Section 4.2. We could, for example, set  $C = z\tau$  as in Section 4.2.

**Lemma 8.** Thinking of  $\gamma_u$  and  $\gamma_v$  as characteristic vectors of sets, and letting  $t = |\Gamma(u) \cup \Gamma(v)|$ , we have

$$\mathbf{E}[|\gamma_u \cup \gamma_v|] = (d/C)(1 - (1 - C/d)^t) < t.$$

Proof. Let  $X_i = 1$  be an indicator random variable for the event that i-th bin is non-empty when throwing t balls into d/C bins. If the bin is empty, then let  $X_i = 0$ . Then  $\mathbf{E}[X_i] = 1 - (1 - C/d)^t$ , and so  $\mathbf{E}[X] = d/C - (d/C)(1 - C/d)^t = d/C(1 - (1 - C/d)^t) \le d/C$ , where  $X = |\gamma_u \cup \gamma_v|$  is the total number of non-empty bins.

By Lemma 8, the expected size of the union of the neighborhoods drops after hashing. This is useful, as the survival probability of the node pair (u,v) in a repetition after hashing is now  $\alpha^{|\gamma(u)\cup\gamma(v)|}$ , which by the previous lemma is larger than before since  $|\gamma(u)\cup\gamma(v)|\leq |\Gamma(u)\cup\Gamma(v)|$ , and this inequality is strict in expectation. Note, however, that the communication and work per machine increase in expectation, but this tradeoff may be beneficial.

# 5 Experimental Results

In this section, we complement the theoretical analysis presented earlier with experiments that measure the recall and efficiency of LSF-Join on three real world graphs from the SNAP repository [20]: WikiVote, PhysicsCitation, and Epinions. In accordance with our motivation, we also run LSF-Join on an extremely skewed synthetic graph, on which the WHIMP algorithm fails.

### Experimental Setup

We compare LSF-Join against the state of the art WHIMP algorithm from [27], and hence our setup is close to the one for WHIMP. In this vein, we transform our graphs into bipartite graphs, either by orienting edges from left to right (for directed graphs), or by duplicating nodes on either side (for undirected ones). This is in accordance with the setup of the left side denoting sets and the right side denoting nodes that is described in the introduction. Also, we pre-filter each bipartite graph to have a narrow degree range on the right (the left degrees can still be O(n)) to minimize variance in cosine similarity values due to degree mismatch. This makes the experiments cleaner, and the algorithm itself can run over all degrees in a doubling manner. We use sparse matrix multiplication for computing all-pairs similarity after computing the survivor sets  $S_i$  for each bucket i, as it is quite fast in practice and consumes  $d \cdot O(|S_i|)$  memory on each server. Finally, even though we

Dataset	N	M	Communication Cost		Recall	
			LSF-Join	WHIMP <sup>†</sup>	LSF-Join	WHIMP
WikiVote	7K	104K	710MB $(\sum_{i}  S_{i}  = 71M, \beta = 30)$	60MB	100%	100%
Citation	34K	421K	410MB $(\sum_{i}  S_{i}  = 41M, \beta = 1)$	50MB	100%	100%
Epinions	60K	500K	6GB $(\sum_{i}  S_{i}  = 573M, \beta = 1)$	$60 \mathrm{MB}$	100%	100%
Synthetic	10M	200M	160GB $(\sum_{i}  S_{i}  = 8B, \beta = 50)$	Failed	90%	_

Table 1: Summary of the performance of LSF-Join and WHIMP on the four datasets, in terms of communication cost and recall (precision for WHIMP was also high). We note that LSF-Join was run at the minimum number of independent iterations  $\beta$  to achieve high recall for  $\tau = 0.1$ .

computed a theoretically optimal value of  $\alpha$  earlier, in practice, a smaller choice of  $\alpha$  often suffices in combination with repeating the Fast-Filter method for  $\beta \geq 1$  independent iterations.

For each of the graphs, we run LSF-Join on the graph on a distributed MapReduce platform internal to Google, and compare the output similar pairs against a ground truth set generated from a sample of the data. The ground truth set is generated by doing an exact all-pairs computation for a small subset of nodes chosen at random. Using this ground truth, we can measure the efficacy of the algorithm, and the measure we focus on for the evaluation is the recall of similar pairs<sup>1</sup>. Specifically, let the set of true similar pairs in the ground truth with similarity at least  $\tau$  be denoted by S. Furthermore, let the set of similar pairs on the same set of nodes that are returned by the algorithm be  $\hat{S}$ . Then, the recall  $R = \frac{|\hat{S} \cap S|}{|S|}$ . For a fixed value of  $\tau$ , we can measure the change in recall as the number  $\beta$  of independent iterations varies (with fixed  $\alpha$  and k = N). We run our experiments at a value of  $\beta$  that achieves high recall (which is a strategy that carries across datasets), and the results are summarized in Table 1 for ease of comparison. There is a synthetic dataset included in the table, which is described later. The communication cost for LSF-Join is dependent on the number of survivors, which in turn depends on the choice of  $\beta$ . We do ignore a subtlety here in that the communication cost will actually often be much less than the number of survivors, since multiple independent repetitions will produce many copies of the same node and hence we can only send one of those copies to a processor.

We reiterate that our experimental comparison is only against the WHIMP algorithm as the WHIMP paper demonstrated that commonly used LSH-based techniques are provably worse. Since WHIMP is only applicable in the scenario where there are no high degree left nodes, our three public

<sup>\*</sup> The communication cost of LSF-Join depends on the number of survivors, which we note along with the value of  $\beta$ .

<sup>†</sup> WHIMP communication cost is dominated by shuffling SimHash sketches. We use 8K bits for SimHash, as suggested in [27].

<sup>&</sup>lt;sup>1</sup>The precision is dependent on the method used to compute all-pairs similarity in a bucket, and since we use sparse matrix multiplication, for us this is 100%.

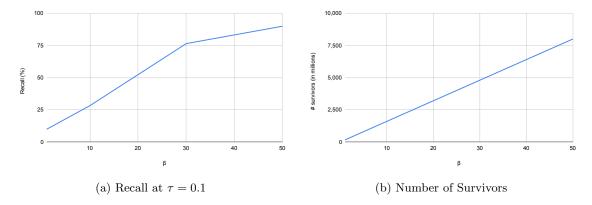


Figure 3: Recall and number of survivors as  $\beta$  increases for the synthetic skewed graph.

graphs are those for which this assumption holds in order to be able to do a comparison. Since the WHIMP algorithm has output-optimal communication complexity, we expect WHIMP to have lower communication cost than LSF-Join, as WHIMP's communication cost is dominated by the number of edges in the graph. This is indeed seen to be the case from Table 1. However, LSF-Join trades-off higher communication cost with the benefit of load balancing across individual servers. WHIMP does not do any load balancing in the worst case, which can render it inapplicable for a broad class of graphs, as we shall see in the next section. Indeed, the WHIMP job failed for our synthetic graph.

### 5.1 Synthetic Graph With Extreme Skew

To illustrate a case that WHIMP fails to address, we present results on a synthetic graph that contains the core element of skeweness that we set out to address in this work. We anticipate that the same results will hold for several real world settings, but a synthetic graph is sufficient for comparison with WHIMP. Indeed, the motivation for this randomly generated synthetic graph comes from user behavior where even though users consume almost the same amount of content (say, videos) online, the content being consumed sees a power-law distribution (e.g., some videos are vastly more popular than others). A simplified setting of the same phenomenon can be captured in the following random bipartite graph construction: we build an  $N \times N$  bipartite graph G(U, V, E), where each right node has degree d. Each right node  $v \in V$  chooses to connect to left nodes as follows: first pick d/2 nodes at random (without replacement) from a small set of hot nodes  $H \subset U$ , and pick d/2 nodes at random (again, without replacement) from the rest of  $U \setminus H$ . If  $|H| = \gamma \cdot d$ , and  $|H| \ll N$ , this results in right nodes having pairwise cosine similarity that scale with  $1/\gamma$  while the hot dimensions have degree O(n) for constant  $\gamma$ . In this setting, we expect wedge sampling-based methods to fail since the hot dimensions have large neighborhoods.

We constructed such a synthetic random bipartite graph with the following parameters: N = 10 million, d = 20, and  $\gamma = 10$ . Then, we repeated the same experiment as the one described above for the real world graphs. This time, we noted that WHIMP failed as the maximum degree for

left nodes was around 500K. We were able to run our procedure though, and the recall and the communication cost of the Fast-Filter procedure is shown in Table 1. The recall of the Fast-Filter procedure is shown in Fig 3a, and the number of survivors in Fig 3b. Note that, as before, we are able to achieve high recall even on this graph with a heavily skewed degree distribution, with reasonable communication cost.

### 6 Conclusion

We present a new distributed algorithm, LSF-Join, for approximate all-pairs set similarity search. The key idea of the algorithm is the use of a novel LSF scheme. We exhibit an efficient version of this scheme that runs in nearly linear time, utilizing pairwise independent hash functions. We show that LSF-Join effectively finds low similarity pairs in high-dimensional datasets with extreme skew. Theoretically, we provide guarantees on the accuracy, communication, and work of LSF-Join. Our algorithm improves over hash-join and LSH-based methods. Experimentally, we show that LSF-Join achieves high accuracy on real and synthetic graphs, even for a low similarity threshold. Moreover, our algorithm succeeds for a graph with extreme skew, whereas prior approaches fail.

**Acknowledgments.** Part of this work was done while D. Woodruff was visiting Google Mountain View. D. Woodruff also acknowledges support from the National Science Foundation Grant No. CCF-1815840.

# References

- [1] Foto N Afrati, Anish Das Sarma, David Menestrina, Aditya Parameswaran, and Jeffrey D Ullman. Fuzzy Joins using MapReduce. In *ICDE*. IEEE, 2012.
- [2] Foto N. Afrati, Anish Das Sarma, Anand Rajaraman, Pokey Rule, Semih Salihoglu, and Jeffrey D. Ullman. Anchor-Points Algorithms for Hamming and Edit Distances Using MapReduce. In ICDT, 2014.
- [3] Thomas Dybdahl Ahle, Rasmus Pagh, Ilya Razenshteyn, and Francesco Silvestri. On the complexity of inner product similarity join. In *Proceedings of the 35th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems*, pages 151–164. ACM, 2016.
- [4] Maha Ahmed Alabduljalil, Xun Tang, and Tao Yang. Optimizing parallel algorithms for all pairs similarity search. In *Proceedings of the sixth ACM international conference on Web search and data mining*, pages 203–212, 2013.
- [5] Nikolaus Augsten and Michael H Böhlen. Similarity joins in relational database systems. Synthesis Lectures on Data Management, 5(5):1–124, 2013.
- [6] Ranieri Baraglia, Gianmarco De Francisci Morales, and Claudio Lucchese. Document similarity self-join with mapreduce. In 2010 IEEE International Conference on data mining, pages 731– 736. IEEE, 2010.

- [7] Roberto J Bayardo, Yiming Ma, and Ramakrishnan Srikant. Scaling up All Pairs Similarity Search. In WWW. ACM, 2007.
- [8] Paul Beame, Paraschos Koutris, and Dan Suciu. Communication steps for parallel query processing. In *PODS*, pages 273–284. ACM, 2013.
- [9] Paul Beame, Paraschos Koutris, and Dan Suciu. Skew in Parallel Query Processing. In PODS. ACM, 2014.
- [10] Paul Beame and Cyrus Rashtchian. Massively-parallel similarity join, edge-isoperimetry, and distance correlations on the hypercube. In *Proceedings of the Twenty-Eighth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 289–306. Society for Industrial and Applied Mathematics, 2017.
- [11] Tobias Christiani. A Framework for Similarity Search with Space-Time Tradeoffs Using Locality-Sensitive Filtering. In *Proceedings of the Twenty-Eighth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 31–46. Society for Industrial and Applied Mathematics, 2017.
- [12] Tobias Christiani, Rasmus Pagh, and Johan Sivertsen. Scalable and robust set similarity join. In 2018 IEEE 34th International Conference on Data Engineering (ICDE), pages 1240–1243. IEEE, 2018.
- [13] Graham Cormode and S. Muthukrishnan. An improved data stream summary: the count-min sketch and its applications. *J. Algorithms*, 55(1):58–75, 2005.
- [14] Dong Deng, Yufei Tao, and Guoliang Li. Overlap set similarity joins with theoretical guarantees. In *Proceedings of the 2018 International Conference on Management of Data*, pages 905–920. ACM, 2018.
- [15] Fabian Fier, Nikolaus Augsten, Panagiotis Bouros, Ulf Leser, and Johann-Christoph Freytag. Set similarity joins on mapreduce: an experimental survey. *Proceedings of the VLDB Endowment*, 11(10):1110–1122, 2018.
- [16] Sariel Har-Peled, Piotr Indyk, and Rajeev Motwani. Approximate nearest neighbor: Towards removing the curse of dimensionality. *Theory of Computing*, 8(1):321–350, 2012.
- [17] Xiao Hu, Ke Yi, and Yufei Tao. Output-Optimal Massively Parallel Algorithms for Similarity Joins. *ACM Trans. Database Syst.*, 44(2):6:1–6:36, April 2019.
- [18] Howard Karloff, Siddharth Suri, and Sergei Vassilvitskii. A model of computation for mapreduce. In *Proceedings of the twenty-first annual ACM-SIAM symposium on Discrete Algorithms*, pages 938–948. SIAM, 2010.
- [19] Paraschos Koutris, Semih Salihoglu, and Dan Suciu. Algorithmic aspects of parallel data processing. Foundations and Trends® in Databases, 8(4):239–370, 2018.

- [20] Jure Leskovec and Andrej Krevl. SNAP Datasets: Stanford large network dataset collection. http://snap.stanford.edu/data, June 2014.
- [21] Jure Leskovec, Anand Rajaraman, and Jeffrey David Ullman. *Mining of Massive Datasets*. Cambridge University Press, 2014.
- [22] Willi Mann, Nikolaus Augsten, and Panagiotis Bouros. An empirical evaluation of set similarity join techniques. *Proceedings of the VLDB Endowment*, 9(9):636–647, 2016.
- [23] Samuel McCauley, Jesper W Mikkelsen, and Rasmus Pagh. Set similarity search for skewed data. In *Proc. of the 37th Symp. on Principles of Database Systems (PODS)*, pages 63–74. ACM, 2018.
- [24] Samuel McCauley and Francesco Silvestri. Adaptive mapreduce similarity joins. In *Proc. 5th ACM SIGMOD Workshop on Algorithms and Systems for MapReduce and Beyond*, page 4. ACM, 2018.
- [25] Rasmus Pagh, Nina Mesing Stausholm, and Mikkel Thorup. Hardness of bichromatic closest pair with jaccard similarity. In 27th Annual European Symposium on Algorithms (ESA 2019), 2019.
- [26] Anish Das Sarma, Foto N. Afrati, Semih Salihoglu, and Jeffrey D. Ullman. Upper and Lower Bounds on the Cost of a Map-reduce Computation. *Proc. VLDB Endow.*, 6(4):277–288, February 2013.
- [27] Aneesh Sharma, C Seshadhri, and Ashish Goel. When hashes met wedges: A distributed algorithm for finding high similarity vectors. In *Proceedings of the 26th International Conference on World Wide Web (WWW)*, pages 431–440, 2017.
- [28] Yasin N. Silva, Jason Reed, Kyle Brown, Adelbert Wadsworth, and Chuitian Rong. An experimental survey of mapreduce-based similarity joins. In Laurent Amsaleg, Michael E. Houle, and Erich Schubert, editors, Similarity Search and Applications: 9th International Conference, SISAP 2016, Tokyo, Japan, October 24-26, 2016, Proceedings, pages 181–195, Cham, 2016. Springer International Publishing.
- [29] N. Sundaram, A. Turmukhametova, N. Satish, T. Mostak, P. Indyk, S. Madden, and P. Dubey. Streaming Similarity Search Over One Billion Tweets Using Parallel Locality-Sensitive Hashing. PVLDB, 6(14):1930–1941, 2013.
- [30] Rares Vernica, Michael J Carey, and Chen Li. Efficient Parallel Set-similarity Joins using MapReduce. In SIGMOD, pages 495–506. ACM, 2010.
- [31] Hongya Wang, Jiao Cao, LihChyun Shu, and Davood Rafiei. Locality Sensitive Hashing Revisited: Filling the Gap Between Theory and Algorithm Analysis. In *CIKM*, pages 1969–1978, New York, NY, USA, 2013. ACM.

- [32] Jiannan Wang, Guoliang Li, and Jianhua Feng. Can we beat the prefix filtering?: an adaptive framework for similarity join and search. In *Proceedings of the 2012 ACM SIGMOD International Conference on Management of Data*, pages 85–96. ACM, 2012.
- [33] Xubo Wang, Lu Qin, Xuemin Lin, Ying Zhang, and Lijun Chang. Leveraging set relations in exact set similarity join. *Proceedings of the VLDB Endowment*, 10(9):925–936, 2017.
- [34] Chuan Xiao, Wei Wang, Xuemin Lin, Jeffrey Xu Yu, and Guoren Wang. Efficient Similarity Joins for Near-duplicate Detection. ACM Transactions on Database Systems, 36(3):15, 2011.
- [35] Chenyun Yu, Sarana Nutanong, Hangyu Li, Cong Wang, and Xingliang Yuan. A generic method for accelerating lsh-based similarity join processing. *IEEE Transactions on Knowledge and Data Engineering*, 29(4):712–726, 2016.
- [36] Erkang Zhu, Fatemeh Nargesian, Ken Q Pu, and Renée J Miller. Lsh ensemble: internet-scale domain search. *Proceedings of the VLDB Endowment*, 9(12):1185–1196, 2016.