# Fast Similarity Sketching[*]

Søren Dahlgaard, Jakob Bæk Tejs Houen, Mathias Bæk Tejs Langhede, and Mikkel Thorup

University of Copenhagen
soren.dahlgaard@gmail.com,jakob@tejs.dk,mathias@tejs.dk,
mikkel2thorup@gmail.com

## Abstract

We consider the *Similarity Sketching* problem: Given a universe $[u] = \{0, \ldots, u-1\}$ we want a random function $S$ mapping subsets $A \subseteq [u]$ into vectors $S(A)$ of size $t$, such that the Jaccard similarity $J(A, B) = |A \cap B|/|A \cup B|$ between sets $A$ and $B$ is preserved. More precisely, define $X_i = [S(A)[i] = S(B)[i]]$ and $X = \sum_{i \in [t]} X_i$. We want $\mathrm{E}[X_i] = J(A, B)$, and we want $X$ to be strongly concentrated around $\mathrm{E}[X] = t \cdot J(A, B)$ (i.e. Chernoff-style bounds). This is a fundamental problem which has found numerous applications in data mining, large-scale classification, computer vision, similarity search, etc. via the classic MinHash algorithm. The vectors $S(A)$ are also called *sketches*. Strong concentration is critical, for often we want to sketch many sets $B_1, \ldots, B_n$ so that we later, for a query set $A$, can find (one of) the most similar $B_i$. It is then critical that no $B_i$ looks much more similar to $A$ due to errors in the sketch.

The seminal $t \times MinHash$ algorithm uses $t$ random hash functions $h_1, \ldots, h_t$, and stores $(\min_{a \in A} h_1(A), \ldots, \min_{a \in A} h_t(A))$ as the sketch of $A$. The main drawback of MinHash is, however, its $O(t \cdot |A|)$ running time, and finding a sketch with similar properties and faster running time has been the subject of several papers. Addressing this, Li et al. [NIPS'12] introduced *one permutation hashing (OPH)*, which creates a sketch of size $t$ in $O(t + |A|)$ time, but with the drawback that possibly some of the $t$ entries are "empty" when $|A| = O(t)$. One could argue that sketching is not necessary in this case, however the desire in most applications is to have *one* sketching procedure that works for sets of all sizes. Therefore, filling out these empty entries is the subject of several follow-up papers initiated by Shrivastava and Li [ICML'14]. However, these "densification" schemes fail to provide good concentration bounds exactly in the case $|A| = O(t)$, where they are needed.

In this paper we present a new sketch which obtains essentially the best of both worlds. That is, a fast $O(t \log t + |A|)$ expected running time while getting the same strong concentration bounds as $t \times \mathrm{MinHash}$. Our new sketch can be seen as a mix between sampling with replacement and sampling without replacement. We demonstrate the power of our new sketch by considering popular applications in large-scale classification with linear Support Vector Machines (SVM) as introduced by Li et al. [NIPS'11] as well as approximate similarity search using the Locality Sensitive Hashing (LSH) framework of Indyk and Motwani [STOC'98].

---

# 1   Introduction

In this paper we consider the following problem which we call the *similarity sketching* problem. Given a large key universe $[u] = \{0, \ldots, u-1\}$ and positive integer $t$ we want a random function $S$ mapping subsets $A \subseteq [u]$ into vectors (which we will call sketches) $S(A)$ of size $t$, such that similarity is preserved. More precisely, we consider the *Jaccard similarity* $J(A, B) = |A \cap B|/|A \cup B|$ between sets $A$ and $B$. Define $X_i = [S(A)[i] = S(B)[i]]$ for each $i \in [t]$, where $S(A)[i]$ denotes the $i$th entry of the vector $S(A)$ and $[x]$ is the Iverson bracket notation with $[x] = 1$ when $x$ is true and 0 otherwise. We want $\mathrm{E}[X_i] = J(A, B)$ for each $i \in t$. Moreover, we want $X = \sum_{i \in [t]} X_i$ to be strongly concentrated around $\mathrm{E}[X] = t \cdot J(A, B)$. That is, the sketches can be used to estimate $J(A, B)$ by doing a pair-wise comparison of corresponding entries. We will call this the *alignment property* of the similarity sketch. Strong concentration, with error probabilities dropping exponentially in $t$ like with Chernoff-bounds, is critical. Often we want to sketch many sets $B_1, \ldots, B_n$ so that we later, for a query set $A$, can find (one of) the most similar $B_i$. It is then critical that no $B_i$ looks much more similar to $A$ due to errors in the sketch.

The standard solution to the similary sketching problem is $t \times$MinHash algorithm[1]. The algorithm works as follows: Let $h_0, \ldots, h_{t-1} : [u] \to [0, 1]$ be random hash functions and define $S(A) = (\min_{a \in A} h_0(a), \ldots, \min_{a \in A} h_{t-1}(a))$. This corresponds to sampling $t$ elements from $A$ with replacement and thus has all the above desired properties.

MinHash was originally introduced by Broder et al. [4, 6, 5] for the AltaVista search engine and has since been used as a standard tool in many applications including duplicate detection [6, 12], all-pairs similarity [3], large-scale learning [17], computer vision [22], and similarity search [14]. The application [17] to Support Vector Machines (SVMs) in Machine Learning is an instructive example of the use of the alignment property. The basic idea is that we want the similarity between sets to be grow with a dot-product between associated vectors. To get a bit vector, [17] suggests that for each coordinate, we hash to get a single bit and based on this bit, replace the coordinate with two bits: `01` or `10`. Now, for the dot-product, if we had a match in a coordinate, we get a match in both bits, adding two to the dot-product. If we did not have a match, then with probability $1/2$, either both or no bits will match. Dissimilarity is thus essentially halved in the reduction. The more important thing is that more similar sets are expected to get higher dot-products, and this is the main point for the SVM applications. Mathematically, a cleaner alternative is to use the hash bit to replace the the coordinate with $-1/\sqrt{t}$ or $1/\sqrt{t}$. Now, in expectation, the dot-product is exactly the Jaccard similarity.

The main drawback of $t \times$MinHash is the $O(t \cdot |A|)$ running time that we pay because we have to find the minimum in $A$ with $t$ independent hash functions. To appreciate the scale of the different parameters, let us just consider the classic application from [4, 6, 18] for text similarity. For each text, they look at the set of $w$-shingles[2] where a $w$-shingle is a tuple of $w$ consecutive words, e.g., [4, 6] use $w = 4$. The similarity between texts are the similarity between their sets of $w$-shingles. With roughly $10^5$ common english words, the number of possible 4-shingles is $u \approx 10^{20}$. Also, note that the number of distinct 4-shingles in a large text can be close to the text size even though the number of distinct words is much smaller, so the sets we consider can be quite large.

While the sketch size $t$ is typically meant to be much smaller than the set size, it can still be sizeable, e.g., [17] suggests using $t = 500$ and [15] suggests using $t = 4000$. From a more theoretical perspective, if we use Locallity Sensitive Hashing (LSH) for set similarity among $n$

---

[1] `https://en.wikipedia.org/wiki/MinHash`

[2] $w$-shingles are also referred to as $w$-grams, e.g., when used for finding similarity between DNA sequences (see `wikipedia.org/wiki/N-gram`).

sets, then $t = \Omega(n^\rho)$ where $\rho$ is a constant parameter. We shall return to this application later.

If we do not care about alignment, then an alternative to $t\times$MinHash sketching is the Bottom-$t$ sketch described in [4, 8, 26]. The idea is to use a single hash function and just store the $t$ smallest hash values from $A$ (so $1\times$MinHash=Bottom-1). We can find the bottom-$t$ sketch of $A$ in $O(|A|)$ time. However, with $t > 1$, there is no alignment between the $t$ sketch values from different sets, and the alignment is needed for applications in LSH as well as for Support Vector Machines (SVM) that we will also discuss later.

Bachrach and Porat [2] suggested a more efficient way of computing $t\times$MinHash values with $t$ different hash functions. They use $t$ different polynomial hash functions that are related, yet pairwise independent, so that they can systematically maintain the MinHash for all $t$ polynomials in $O(\log t)$ time per element of $A$. There are two issues with this approach: It is specialized to work with polynomials and MinHash is known to have constant bias unless the polynomials considered have super-constant degree [20], and this bias does not decay with independent repetitions. Also, because the experiments are only pairwise independent, the concentration is only limited by Chebyshev's inequality and thus nowhere near the Chernoff bounds we want for many applications.

Another direction introduced by Li et al. [16] is *one permutation hashing* (OPH) which works by hashing the elements of $A$ into $t$ buckets and performing a MinHash in each bucket using the same hash function. While this procedure gives $O(t + |A|)$ sketch creation time it also may create empty buckets and thus only obtains a sketch with $t' \leqslant t$ entries when $|A| = o(t \log t)$. One may argue that sketching is not even needed in this case. However, a common goal in applications of similarity sketching is to have *one* sketching procedure which works for all set size – one data structure that works for an entire collection of data sets of different sizes in the case of approximate similarity search. It is thus very desirable that the similarity sketch works well independently of the size of the input set.

Motivated by this, several follow-up papers [25, 24, 23] have tried to give different schemes for filling out the empty entries of the OPH sketch ("densifying" the sketch). These papers all consider different ways of copying from the full entries of the sketch into the empty ones. Due to this approach, however, these densification schemes all fail to give good concentration guarantees when $|A|$ is small, which is exactly the cases in which OPH gives many empty bins and densification is needed. This is because of the fundamental problem that unfortunate collisions in the first round cannot be resolved in the second round when copying from the full bins. To understand this consider the following extreme example: Let $A$ be a set with two elements. Then with probability $1/t$ these two elements end in the same bin where only one survives (the one with the smallest hash value). Densification will then just copy the surviver to all $t$ entries. Such issues may lead to very poor similarity estimation and this is illustrated with experiments in Figure 1. A further issue is that the state-of-the-art densification scheme of Shrivastava [23] has a running time of $O(t^2)$ when $|A| = O(t)$.

## 1.1 Our contribution

In this paper we obtain a sketch which essentially obtains the best of both worlds. That is, strong concentration guarantees for similarity estimation as well as a fast expected sketch creation time of $O(t \log t + |A|)$. On top of that, our new sketching algorithm is simple and easy to implement.

Our new sketch can be seen as a mixture between sampling with and without replacement and in many cases outperforms $t\times$MinHash. An example of this can be seen in the toy example of Figure 1, where the "without replacement"-part of our sketch gives better concentration compared to MinHash. Our sketch can be employed in places where $t\times$MinHash is currently employed to improve the running time from the "quadratic" $O(t \cdot |A|)$ to the "near-linear" $O(t \log t + |A|)$. In
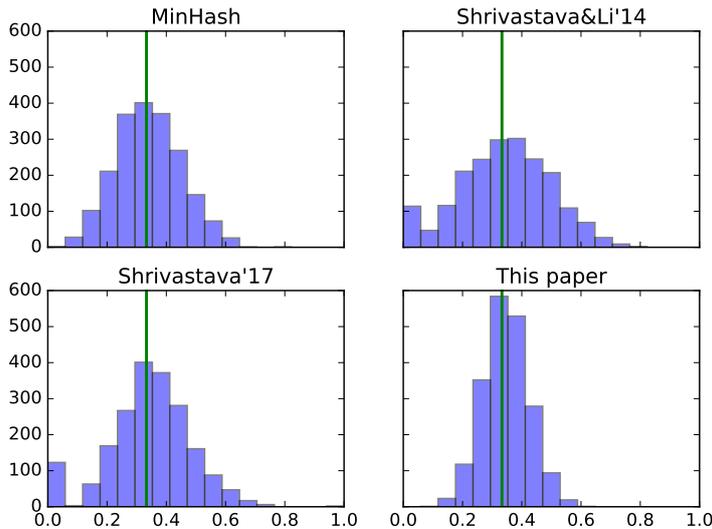
Figure 1: Experimental evaluation of similarity estimation of the sets $A = \{1, 2\}$ and $B = \{2, 3\}$ with different similarity sketches and $t = 16$. Each experiment is repeated 2000 times and the $y$-axis reports the frequency of each estimate. The green line indicates the actual similarity. The two methods based on OPH perform poorly as each set has a probability of $1/t$ to be a single-element sketch. Our new method outperforms $t \times$MinHash as it has an element of "without replacement".

this paper we focus on two popular applications, which are large-scale learning with linear SVM and approximate similarity search with LSH. Before discussing these applications, we describe our result more formally.

**Theorem 1.** [3] *Let* $[u] = \{0, 1, 2, \ldots, u - 1\}$ *be a set of keys and let* $t$ *be a positive integer. There exists an algorithm that given a set* $A \subseteq [u]$ *in expected time* $O(|A| + t \log t)$ *creates a size-t vector* $v(A)$ *of non-negative real numbers with the following properties.*

*For two sets* $A, B \subseteq [u]$ *with Jaccard similarity* $J(A, B) = J$ *it holds that* $v(A \cup B)_i = \min\{v(A)_i, v(B)_i\}$ *for each index* $i \in [t]$. *For* $i \in [t]$ *let* $X_i = 1$ *if* $v(A)_i = V(B)_i$ *and* $0$ *otherwise. Let* $I \subseteq [t]$ *be a subset of* $k$ *indices. Let* $j \in [t] \backslash I$, *then,*

$$\min\left(J, \frac{tJ - \sum_{i \in I} X_i}{t - k}\right) \leqslant \Pr\left[X_j = 1 \,|\, \sigma\left((X_i)_{i \in I}\right)\right] \leqslant \max\left(J, \frac{tJ - \sum_{i \in I} X_i}{t - k}\right).$$

If we sampled with replacement then the probability that $X_j = 1$ when conditioning on $(X_i)_{i \in I}$ is exactly $J$, and if we sampled without replacement then the probability that $X_j = 1$ when conditioning on $(X_i)_{i \in I}$ is exactly $\frac{tJ - \sum_{i \in I} X_i}{t - k}$. Thus the theorem shows qualitatively that our new sketch falls between sampling with replacement and sampling without replacement.

Now a nice implication of the theorem is that we get classic Chernoff concentration bounds.

**Corollary 1.** *Use the same setting as Theorem 1 and let* $X = \frac{1}{t} \sum_{i \in [t]} X_i$. *Then* $E[X] = J$ *and*

---

[3]This powerful theorem was not in our original conference version [10].

3

*for every $\delta > 0$ it holds that:*

$$\Pr[X \geqslant J(1+\delta)] \leqslant \left(\frac{e^\delta}{(1+\delta)^{1+\delta}}\right)^{tJ},$$

$$\Pr[X \leqslant J(1-\delta)] \leqslant \left(\frac{e^{-\delta}}{(1-\delta)^{1-\delta}}\right)^{tJ}.$$

The above theorems assume access to fully random hashing (this is also assumed for the standard MinHash sketch). However, we will also show how to implement our sketch using mixed tabulation hashing (introduced by Dahlgaard et al. [9]) which is practical and can be evaluated in $O(1)$ time. The concentration bounds get slightly weaker and more complicated, as they do with any realistic hashing scheme.

## 1.2 Speeding up Locality Sensitive Hashing on Sets

One of the most powerful applications of the $t \times$MinHash algorithm is for the *approximate set similarity search problem* using *Locality Sensitive Hashing (LSH)*. For this application we will crucially need both the alignment and strong probability bounds from Theorem 1. Our fundamental contribution will be to improve the time bound, reducing the effect of the set size from multiplicative to additive. This is very important when dealing with larger sets, e.g., the above mentioned examples where the sets are the 4-shingles of texts [4, 6, 18].

While our new similarity sketch forms the base of our improvement, we need several other ideas to address other challenges, particularly when it comes to high probability results. To describe our contribution, we first need to revisit the Locality Sensitive Hashing framework introduced by Indyk and Motwani [14]. It is a general framework based on the following type of data-independent hash families:

**Definition 1** (Locality sensitive hashing [14]). Let $(X, S)$ be a similarity space and let $\mathcal{H}$ be a family of hash functions $h : X \to R$. We say that $\mathcal{H}$ is $(s_1, s_2, p_1, p_2)$-sensitive if for any $x, y \in X$ and $h \in \mathcal{H}$ chosen uniformly random we have that

- If $S(x, y) \geqslant s_1$ then $\Pr[h(x) = h(y)] \geqslant p_1$.

- If $S(x, y) \leqslant s_2$ then $\Pr[h(x) = h(y)] \leqslant p_2$.

The Locality Sensitive Hashing framework takes a $(s_1, s_2, p_1, p_2)$-sensitive family $\mathcal{H}$ for a similarity space $(X, S)$, and uses it to solve the Approximate Similarity Search problem for a stored set $Y \subseteq X$ with parameters $0 < s_2 < s_1 < 1$. That is, given a query $q \in X$ it returns an element $a \in Y$ with $S(a, q) \geqslant s_2$, if there exists an element $b \in Y$ with $S(b, q) \geqslant s_1$ with constant probability. Note that the definition of the locality sensitive hash function $\mathcal{H}$ is oblivious to the concrete stored set $Y$.

In this paper we only focus on the Jaccard set similarity with stored set of sets $Y = \mathcal{F}$. The classical way of using the LSH framework with respect to the Jaccard similarity is using the seminal MinHash algorithm, which was originally introduced by Broder et al. [4, 6]. The MinHash algorithm is defined as follows: Given a family $\mathcal{H}$ of hash functions $h : U \to R$ we define a new family $\mathcal{H}^{\min}$ of hash functions $h^{\min} : \mathcal{P}(U) \to U$ where $h^{\min}(A) = \arg\min_{x \in A} h(x)$ for any $A \subseteq U$. If the domain $R$ is large enough, we can assume that there are no collisions, hence that $h^{\min}$ is well-defined. If $h \in \mathcal{H}$ is chosen uniformly at random we get that $\Pr[h^{\min}(A) = h^{\min}(B)] = J(A, B)$ for any $A, B \subseteq U$. This shows that $\mathcal{H}^{\min}$ is a $(j_1, j_2, j_1, j_2)$-sensitive family for any $0 < j_2 < j_1 < 1$.

The next idea from [14] is to create sketches using $K$ functions from $\mathcal{H}^{\min}$. For a set $A$, we get the sketch $S(A) = (h_0^{\min}(A), \ldots, h_{K-1}^{\min}(A))$ where $h_0, \ldots, h_{K-1} \in \mathcal{H}$. If the hash functions are independent, then $\Pr[S(A) = S(B)] = J(A,B)^K$, so $S$ is $(j_1, j_2, j_1^K, j_2^K)$-sensitive. Setting $K = \left\lceil \frac{\log(n)}{\log(1/j_2)} \right\rceil$, we get that $S$ is $(j_1, j_2, O(1/n^\rho), 1/n)$-sensitive where $\rho = \frac{\log(1/j_1)}{\log(1/j_2)}$. While the above construction may seem a bit ad-hoc, O'Donnell et al. [19] have shown it to be optimal in the sense that we cannot in general construct a $(j_1, j_2, o(1/n^\rho), 1/n)$-sensitive family of hash functions. With $n = |\mathcal{F}|$ stored sets, we only expect a constant number of false matches $S(A) = S(Q)$ where $A \in \mathcal{F}$ and $J(A,Q) < j_2$.

To get a positive match with constant probability, we use $L = \lceil n^\rho \rceil$ sketch functions $S_0, \ldots, S_{L-1}$. If we have a set $B$ with $J(Q,B) \geqslant j_2$, then with constant probability, there is an $i \in [L]$ such that $S_i(B) = S_i(Q)$. To create a data structure, for each $i$, we have a hash table that with each sketch value $s$ stores pointers to all sets $A \in \mathcal{F}$ with $S_i(A) = s$. When we query a set $Q$ we compute the sketch $S_i(Q)$ and lookup all the matches $A \in \mathcal{F}$ with $S_i(A) = S_i(Q)$ using the hash tables for each $i \in [L]$. We check the matches one by one to see if $J(A,Q) \leqslant j_2$, stopping if a good one is found. This data structure has constant 1-sided error probability. It uses $O(n \cdot L + \sum_{A \in \mathcal{F}} |A|) = O(n^{1+\rho} + \sum_{A \in \mathcal{F}} |A|)$ space and $O(L \cdot K \cdot |Q|) = O(n^\rho \log n \cdot |Q|)$ query time. More precisely, the query time is dominated by two parts:

(1) We have to compute $O(L \cdot K)$ hash values for similarity sketches which takes $O(L \cdot K \cdot |Q|)$ time using $O(L \cdot K) \times$ MinHash.

(2) In expectation the data structure returns $O(L)$ false positives which have to be filtered out. This takes $O(L \cdot |Q|)$ time.

Different techniques has been used to speed up the query time and mostly the focus has been on improving the dominant part (1). Andoni and Indyk [1] looked at the general LSH framework and limited the number of evaluations of locality sensitive hash functions. The idea is to create the sketches by combining smaller sketches together. More precisely, a much smaller collection of sketches of size $m = o(L)$ is created and then every $\binom{m}{t}$ combination of sketches is formed. This technique is also known as tensoring. In the context of Jaccard similarity it improved part (1) of the query time from $O(L \cdot K \cdot |Q|)$ to $O(L \cdot |Q|)$, matching the bound for part (2). Thus they got an overall query time of $O(L \cdot |Q|) = O(n^\rho \cdot |Q|)$.

**Contribution** Our original conference version [10] had overlooked [1] and focussed directly on improving the original bounds in (1) and (2). Concerning (1), we note that the analysis assumes that all the values in the $O(L \cdot K) \times$ MinHash are independent. This is not the case for our new similarity sketch from Theorem 1 with $t = O(L \cdot K)$, but it turns out that the probability bounds from Theorem 1 do suffice. This implements (1) in $O(t \log t + |Q|) = O(L \cdot K \cdot \log n + |Q|)$. For a better implementation, we can first create an intermediate sketch of size $O(\log^2(n))$ and then sample the $L$ sketches of size $K$ from this intermediate sketch. This improves part (1) of the query time to $O(K \cdot L + |Q|)$.

Improving (2) requires some different ideas, but already in [10], we improved part (2) of the query time to $O(L + |Q|)$.

Christiani [7] noticed that our construction also composes nicely with tensoring technique [1] so that we only need $L$ instead of $L \cdot K$ sketch values. As a result, part (1) of the query time is improved from $O(K \cdot L + |Q|)$ to $O(L + |Q|)$, matching the time for part (2). Hence the total query time becomes $O(L + |Q|) = O(L + |Q|) = O(n^\rho + |Q|)$, which is the natural target for constant error probability. This is the combined solution presented in the current full paper.

Often we want a smaller error probability $\varepsilon > 0$, e.g., $\varepsilon = 1/n$. The generic standard approach is to use $O(\log(1/\varepsilon))$ independent data structures, each failing with constant probability, and

then return the best solution found, if any. Indeed this is suggested by Motwani and Indyk in [14, p. 605] and [11, p. 327]. Since the data structures can use a common representation of the sets, we would get a space usage of $O(n^{1+\rho} \log(1/\varepsilon) + \sum_{A \in \mathcal{F}} |A|)$ and a query time of $O((n^\rho + |Q|) \log(1/\varepsilon))$.

Here we further improve the query time to $O(n^\rho \log(1/\varepsilon) + |Q|)$ while having the same space usage. Thus we preserve our optimal linear dependence on $|Q|$ even for high probability results. Adding it all up, we prove the following result for the approximate set similarity search problem using LSH:

**Theorem 2.** *We are given a family $\mathcal{F}$ of up to $n$ sets from a large universe $U$. Moreover, we are given two constant parameters $j_1$ and $j_2$ with $0 < j_2 < j_1 < 1$, as well as an error parameter $\varepsilon > 0$ which may be subconstant.*

*We present an Approximate Similarity Search data structure with 1-sided error probability $\varepsilon$: given a query set $Q \subseteq U$, if there is a set $B \in \mathcal{F}$ with $J(B, Q) \geqslant j_1$, the data structure returns a set $A \in \mathcal{F}$ with $J(A, Q) \geqslant j_2$ with probability $1 - \varepsilon$. Moreover, if $A$ is returned, it always has $J(A, Q) \geqslant j_2$. With $\rho = \frac{\log(1/j_1)}{\log(1/j_2)}$ our data structure uses $O\left(n^{1+\rho} \log(1/\varepsilon) + \sum_{A \in \mathcal{F}} |A|\right)$ space and it has query time $O\left(n^\rho \log(1/\varepsilon) + |Q|\right)$.*

## 1.3 Notation

For a real number $x$ and an integer $k$ we define $x^{\underline{k}} = x(x-1)(x-2)\dots(x-k+1)$. For an expression $P$ we let $[P]$ denote the variable that is 1 if $P$ is true and 0 otherwise. For a non-negative integer $n$ we let $[n]$ denote the set $[n] = \{0, 1, 2, \dots, n-1\}$.

## 2 Fast Similarity Sketching

In this section we present our new sketching algorithm, which takes a set $A \subseteq [u]$ as input and produces a sketch $S(A, t)$ of size $t$. When $t$ is clear from the context we may write just $S(A)$.

Our new similarity sketch is simple to describe: Let $h_0, \dots, h_{2t-1}$ be random hash functions such that for $i \in [t]$ we have $h_i : [u] \to [t] \times [i, i+1)$ and for $i \in \{t, \dots, 2t-1\}$ we have $h_i : [u] \to \{i - t\} \times [i, i+1)$. For each hash function $h_i$ we say that the output is split into a bin, $b_i$, and a value, $v_i$. That is, for $i \in [2t]$ and $a \in [u]$ we have $h_i(a) = (b_i(a), v_i(a))$, where $b_i(a)$ and $v_i(a)$ are restricted as described above. We may then define the $j$th entry of the sketch $S(A)$ as follows:

$$S(A)[j] = \min\{v_i(a) \mid a \in A, i \in [2t], b_i(a) = j\} \ . \tag{1}$$

In particular, the hash functions $h_t, \dots, h_{2t-1}$ ensure that each entry of $S(A)$ is well-defined. Furthermore, since we have $v_i(a) < v_j(b)$ for any $a, b \in [u]$ and $0 \leqslant i < j < 2t$ we can efficiently implement the sketch defined in (1) using the Similarity-Sketch procedure in Algorithm 1. A bin $S[b]$ gets "filled" the first time it gets assigned a value $< \infty$ in line 5. The algorithm terminates when all bins are filled, and if this happens in the first round with $i = 0$, then the sketch created is identical to that of *one permutation hashing* [16].

We will start our analysis of $S(A)$ by bounding the running time of Algorithm 1.

**Lemma 1.** *Let $A \subseteq [u]$ be some set and let $t$ be a positive integer. Then the expected running time of Algorithm 1 is $O(t \log t + |A|)$.*

*Proof.* We always have a trivial worst-case upper bound of $O(t \cdot |A|)$ and this is $O(t \log t)$ if $|A| = O(\log t)$. Otherwise, we may assume $|A| \geqslant 2 \log t$. Fix $i$ to be the smallest value such that

---

**Algorithm 1:** `Similarity-Sketch`

    **input** : $A$, $t$, $h_0, \ldots, h_{2t-1}$
    **output:** The sketch $S(A, t)$

1  $S \leftarrow \infty^t$
2  $c \leftarrow 0$
3  **for** $i = 0, \cdots, 2t - 1$ **do**
4      **for** $a \in A$ **do**
5         $b, v \leftarrow h_i(a)$
6         **if** $S[b] = \infty$ **then**
7            $c \leftarrow c + 1$
8         $S[b] \leftarrow \min(S[b], v)$
9      **if** $c = t$ **then**
10        **return** $S$

---

$|A| \cdot i \geqslant 2 \cdot t \log t$. Then $i \leqslant t$, and then the probability that a given bin is empty after evaluating $h_0, \ldots, h_{i-1}$ is at most

$$(1 - 1/t)^{|A| \cdot i} \leqslant (1 - 1/t)^{2 \cdot t \log t} \leqslant 1/t^2 \ .$$

It follows that the probability of any bin being empty is at most $1/t$ and thus the expected running time is $O(|A| \cdot i + \frac{|A| \cdot t}{t}) = O(t \log t + |A|)$. $\qquad\square$

Next, we will prove several properties of the sketch. The first is an observation that the sketch of the union of two sets can be computed solely from the sketches of the two sets.

**Fact 1.** *Let $A, B$ be two sets and let $t$ be a positive integer. Then*

$$S(A \cup B, t)[i] = \min(S(A, t)[i], S(B, t)[i]) \ .$$

The main technical lemma regarding the sketch is Lemma 2 below. The lemma show that our sketch can qualitatively be seen as a mixture between sampling with replacement and sampling without replacement. We will use this lemma to show that we get an unbiased estimator as well as Chernoff-style concentration bounds.

**Lemma 2.** *Let $A, B$ be sets with Jaccard similarity $J(A, B) = J$, let $t$ be a positive integer For each $i \in [t]$ let $X_i = [S(A, t)[i] = S(B, t)[i]]$. Let $I \subseteq [t]$ be a subset of $k$ indices. Let $j \in [t] \backslash I$ then*

$$\min\left(J, \frac{tJ - \sum_{i \in I} X_i}{t - k}\right) \leqslant \Pr\left[X_j = 1 \mid \sigma\left((X_i)_{i \in I}\right)\right] \leqslant \max\left(J, \frac{tJ - \sum_{i \in I} X_i}{t - k}\right) \ .$$

*Proof.* Define $\mathcal{T} = (T_0, T_1, \ldots, T_{2t-1})$ in the following way. Let $T_0 = b_0 (A \cup B)$ and for $i \geqslant 1$ let $T_i = b_i (A \cup B) \backslash (T_0 \cup \ldots \cup T_{i-1})$. Assume in the following that $\mathcal{T}$ is fixed. It clearly suffices to prove this theorem for all possible choices of $\mathcal{T}$. Let $n = |A \cup B|$, then $nJ = |A \cap B|$.

We will prove the claim when the set $I$ is chosen uniformly random among the subsets of $[t]$ of size $k$, and where $j$ is chosen uniformly random from $[t] \backslash I$. Because of symmetry this will suffice.

The probability that $j \in T_h$ is $\frac{|T_h| - |T_h \cap I|}{t - k}$. Conditioned on $j \in T_h$ the probability that $X_j = 1$ is exactly $\frac{nJ - \sum_{i \in T_h \cap I} X_i}{n - |T_h \cap I|}$. So the probability that $X_j = 1$ is

$$p = \mathrm{E}\left[\sum_{h \in [2t]} \frac{|T_h| - |T_h \cap I|}{t - k} \cdot \frac{nJ - \sum_{i \in T_h \cap I} X_i}{n - |T_h \cap I|}\right]$$

7

If we fix $|T_h \cap I|$ then $T_h \cap I$ is a uniformly random subset of $I$ of size $|T_h \cap I|$. This implies that

$$E\left[\sum_{i \in T_h \cap I} X_i \,\middle|\, \sigma\left((X_i)_{i \in I}, |T_h \cap I|\right)\right] = |T_h \cap I| \frac{\sum_{i \in I} X_i}{k}$$

and that

$$\sum_{h \in [2t]} E\left[\frac{|T_h| - |T_h \cap I|}{t-k} \cdot \frac{nJ - \sum_{i \in T_h \cap I} X_i}{n - |T_h \cap I|} \,\middle|\, \sigma\left((X_i)_{i \in I}, |T_h \cap I|\right)\right]$$

$$= \sum_{h \in [2t]} \frac{|T_h| - |T_h \cap I|}{t-k} \cdot \frac{nJ - |T_h \cap I| \frac{\sum_{i \in I} X_i}{k}}{n - |T_h \cap I|}$$

If $J \leqslant \frac{tJ - \sum_{i \in I} X_i}{t-k}$ then $J \geqslant \frac{\sum_{i \in I} X_i}{k}$ which implies that

$$J \leqslant \frac{nJ - |T_h \cap I| \frac{\sum_{i \in I} X_i}{k}}{n - |T_h \cap I|} \leqslant \frac{|T_h| J - |T_h \cap I| \frac{\sum_{i \in I} X_i}{k}}{|T_h| - |T_h \cap I|}$$

and inserting these estimates give us

$$J = \sum_{h \in [2t]} \frac{|T_h| - |T_h \cap I|}{t-k} \cdot J$$

$$\leqslant \sum_{h \in [2t]} \frac{|T_h| - |T_h \cap I|}{t-k} \cdot \frac{nJ - |T_h \cap I| \frac{\sum_{i \in I} X_i}{k}}{n - |T_h \cap I|}$$

$$\leqslant \sum_{h \in [2t]} \frac{|T_h| - |T_h \cap I|}{t-k} \cdot \frac{|T_h| J - |T_h \cap I| \frac{\sum_{i \in I} X_i}{k}}{|T_h| - |T_h \cap I|}$$

$$= \frac{tJ - \sum_{i \in I} X_i}{t-k}$$

Similar calculations shows that if $J \geqslant \frac{tJ - \sum_{i \in I} X_i}{t-k}$ then

$$J \geqslant \sum_{h \in [2t]} \frac{|T_h| - |T_h \cap I|}{t-k} \cdot \frac{nJ - |T_h \cap I| \frac{\sum_{i \in I} X_i}{k}}{n - |T_h \cap I|} \geqslant \frac{tJ - \sum_{i \in I} X_i}{t-k}$$

which finishes the proof. $\qquad \square$

As a corollary we immediately get that the estimator is unbiased.

**Lemma 3.** *Let $A, B$ be sets with Jaccard similarity $J(A, B) = J$ and let $t$ be a positive integer. Let $X_i = [S(A, t)[i] = S(B, t)[i]]$ and let $X = \sum_{i \in [t]} X_i$. Then $E[X] = tJ$.*

*Proof.* This follows directly by applying Lemma 2 with $k = 0$. $\qquad \square$

We also get nice bounds on the moments even when conditioning on a subset of the indices. This lemma will be important in Section 3 where we will use it to prove our result of improving LSH.

**Lemma 4.** *Let $A, B$ be sets with Jaccard similarity $J(A, B) = J$, let $t$ be a positive integer. Let $I \subseteq [t]$ be a set of $k$ indices and $K \subseteq [t] \backslash I$ another disjoint set of $m$ indices. Then*

$$\left(\frac{tJ - k}{t-k}\right)^m \leqslant E\left[\prod_{j \in K} X_j \,\middle|\, \sigma\left((X_i)_{i \in I}\right)\right] \leqslant \left(\frac{tJ}{t-k}\right)^m.$$

8

*Proof.* The proof for the lower bound is completely analogous to proof for the upper bound so we only show the arguments for the upper bound.

First we note that if $k \geq t(1 - J)$ then $\frac{tJ}{t-k} \geq 1$ and the result is trivially true. So in the rest of the proof we assume that $k \leq t(1 - J)$

We will prove that for any subset $K' \subseteq K$ and $h \in K \backslash K'$ then the following is true

$$\mathrm{E}\left[ X_h \,\middle|\, \sigma\left((X_i)_{i \in I}\right) \wedge \bigwedge_{j \in K'} (X_j = 1) \right] \leq \frac{tJ}{t - k}$$

This is easily seen by using Lemma 2

$$\mathrm{E}\left[ X_h \,\middle|\, \sigma\left((X_i)_{i \in I}\right) \wedge \bigwedge_{j \in K'} (X_j = 1) \right] \leq \max\left( J, \frac{tJ - |K'| - \sum_{i \in I} X_i}{t - |K'| - k} \right) \leq \frac{tJ}{t - k}$$

Now we enumerate the elements of $K = \{v_0, \ldots, v_{m-1}\}$ and see that

$$\mathrm{E}\left[ \prod_{j \in [m]} X_{v_j} \,\middle|\, \sigma\left((X_i)_{i \in I}\right) \right] = \prod_{j \in [m]} \mathrm{E}\left[ X_{v_j} \,\middle|\, \sigma\left((X_i)_{i \in I}\right) \wedge \bigwedge_{h \in [j]} (X_{v_h} = 1) \right] \leq \left( \frac{tJ}{t - k} \right)^m$$

$\square$

Finally, we also get Chernoff-style concentration bounds as follows.

**Lemma 5.** *Let $A, B$ be sets with Jaccard similarity $J(A, B) = J$ and let $t$ be a positive integer. Let $X_i = [S(A, t)[i] = S(B, t)[i]]$ and let $X = \sum_{i \in [t]} X_i$. Then for $\delta > 0$*

$$\Pr[X \geq J(1 + \delta)] \leq \left( \frac{e^\delta}{(1 + \delta)^{1+\delta}} \right)^t,$$

$$\Pr[X \leq J(1 - \delta)] \leq \left( \frac{e^{-\delta}}{(1 - \delta)^{1-\delta}} \right)^t.$$

*Proof.* The upper bound follows from Lemma 4 with $k = 0$ and [21, Corollary 1] since Chernoff bounds are derived by bounding $\mathrm{E}\left[e^{\lambda X}\right]$ for some $\lambda > 0$.

The lower bounds follows from considering $Y = \sum_{i \in [t]} Y_i$ where $Y_i = 1 - X_i$ and $Y = t - X$. Since $Y_i = [S(A \cup B, t)[i] = S((A \cup B) \backslash (A \cap B), t)[i]]$ we can use the same argument as for the upper bound, see [21, Page 4]. $\square$

**Practical implementation** In Algorithm 1 we used $2t$ fully random hash functions to implement our new similarity sketch. We now briefly describe how to avoid this requirement by instead using just *one* Mixed Tabulation hash function as introduced by Dahlgaard et al. [9]. We do not present the entire details, but refer instead to the theorems of [9] which can be used directly in a black-box fashion. We note that Dahlgaard et al. [9] did address *one permutation hashing* [16] which is identical to our similarity sketch if all bins are filled in the first round with $i = 0$.

In tabulation-based hashing we view each key, $x \in [u]$, as a vector $(x_0, \ldots, x_{c-1})$ of $c$ characters, where each $x_i \in [u^{1/c}] = \Sigma$, and $\Sigma$ is called the *alphabet size*. The space will be $O(c|\Sigma|)$ and hash values are computed in in $O(c)$ time. As in Dahlgaard et al. [9], we need $|\Sigma| \geq \delta \cdot t \log t$ for some sufficiently large constant $\delta$. The output of tabulation hashing is a bit-string that we can easily split into two parts, one for the bin and one for the value.

In our similarity sketch from Algorithm 1, we use $2t$ independent hash fuctions $h_i$. Here we view the index $i$ as an extra most significant character from $[2t] \subseteq \Sigma$. We then use a single mixed tabulation hash function $H$ taking indexed keys from $\Sigma^{c'}$ for $c' = c + 1$. With with original key $a = (a_0, \ldots, a_{c-1})$, we compute $h_i(a)$ in Line 5 of Algorithm 1 as $H(i, a_0, \ldots, a_{c-1})$, but fixing the bin to $i - t$ if $i \geqslant t$,

We now consider two cases:

- Suppose $|A| \leqslant |\Sigma|/2$ and let $i \leqslant 2t$ be an integer such $i|A| \leqslant |\Sigma|/2$. Then we have at most $|\Sigma|/2$ indexed keys in $[i] \times A$. It now follows from [9, Theorem 1], that w.h.p., the indexed keys from $[i] \times A$ hash fully randomly, just like in our original analysis. If $2t|A| \leqslant |\Sigma|/2$, we pick $i = 2t$, and then this implies full randomness over all indexed keys. Otherwise, we pick $i = \min\{t, |\Sigma|/(2|A|)]\}$. Then $i|A| \geqslant |\Sigma|/(4|A|)$. Now, as in the proof of Lemma 1, the probability that a given bin is empty after $i$ rounds in Algorithm 1 is at most

$$(1 - 1/t)^{|A| \cdot i} \leqslant (1 - 1/t)^{(\delta t \log t)/4} < 1/t^{\delta/4}.$$

For a large enough $\delta$, we conclude that all bins are filled, w.h.p., hence that we have full randomness over all indexed keys considered before termination.

- Suppose now that $|A| > |\Sigma|/2$. In particular this implies that we have a subset $A'$ of $|\Sigma|/2$ keys. As in the previous case, $[1] \times A'$ get hashed fully randomly, filling all bins, w.h.p., but this must then also be the case for $[1] \times A$. Thus, w.h.p., Algorithm 1 terminates at the end of the first round, meaning that it behaves like *one permutation hashing* [16]. In this case both correctness and running time follows immediately from [9, Theorem 2].

We note that concentration bounds from [9] for mixed tabulation have been later been improved by Houen and Thorup in [13, §1.6]. This does not, however, change the above description of how we would apply mixed tabulation in our similarity sketch.

# 3   Speeding up LSH

We are now ready to prove Theorem 2. We want to solve the approximate similarity search problem with parameters $0 < j_2 < j_1 < 1$ on a family, $\mathcal{F}$, of $n$ sets from a large universe $U$. We show a solution that uses $O\left(n^{1+\rho} \log(1/\varepsilon) + \sum_{A \in \mathcal{F}} |A|\right)$ space and $O\left(n^\rho \log(1/\varepsilon) + |Q|\right)$ query time.

## 3.1   Creating the sketches

We will create a data structure similar to the LSH structure as described in Section 1.2. Our data structure will have parameters $L, K, M$. For each $A \in \mathcal{F}$ (and query $Q$) we will create $2M \cdot L$ sketches $\left(S'_{i,j}(A)\right)_{i \in [2M], j \in [L]}$ of size $K$ such that for any two sets $A, B \subseteq U$, $i \in [2M]$ and $j \in [L]$ we have the following properties

- $\Pr\left[S'_{i,j}(A) = S'_{i,j}(B)\right] = O(J(A, B)^K)$   .

- If $J(A, B) \geqslant j_1$ then $\Pr\left[S'_{i,j}(A) = S'_{i,j}(B)\right] = \Theta(J(A, B)^K)$   .

We will then combine these sketches into $M \cdot L^2$ new sketches of size $2K$ by using the tensoring technique. Specifically for every $i \in [M]$ and $j, k \in [L]$ we define $S_{i,j \cdot L + k}(A) = (S'_{2i,j}(A), S'_{2i+1,k}(A))$ which is clearly well-defined. By using standard hashing techniques the new sketches can be calculated in $O(M \cdot L^2)$ time. Then for any two sets $A, B \subseteq U$, $i \in [M]$ and $j \in [L^2]$ we have the following properties
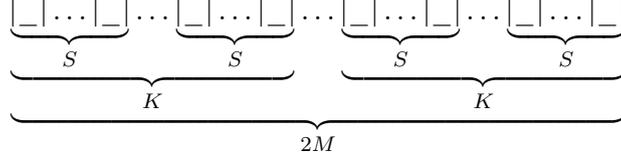
Figure 2: The intermediate sketch $S(A)$ is first partitioned into $2M$ segments which corresponds to the $2M$ subexperiments. Each of these segments then partitioned further into $K$ blocks of size $S$, which corresponds to the $K$ entries in the $L$ sketches in each of the subexperiments.

- $\Pr[S_{i,j}(A) = S_{i,j}(B)] = O(J(A,B)^{2K})$ .

- If $J(A,B) \geqslant j_1$ then $\Pr[S_{i,j}(A) = S_{i,j}(B)] = \Theta(J(A,B)^{2K})$ .

By setting $K = \left\lceil \frac{\log(n)}{2\log(1/j_2)} \right\rceil$, $L = 6\left\lceil (1/j_1)^K \right\rceil$ and $M = \Theta(\log(1/\varepsilon))$ and using the analysis of [14] immediately give us

$$O\left(ML^2 + \sum_{A\in\mathcal{F}} |A|\right) = O\left(n^{1+\rho}\log(1/\varepsilon) + \sum_{A\in\mathcal{F}} |A|\right)$$

space usage and

$$O\left(ML^2 + T(\log(1/\varepsilon)n^{\rho/2}, \log(n), |Q|)\right) = O\left(n^{\rho}\log(1/\varepsilon) + T(n^{\rho/2}\log(1/\varepsilon), \log(n), |Q|)\right)$$

query time, where $T(x,y,z)$ is the time it takes to calculate $x$ sketches of size $y$ for a set of size $z$, and $\rho = \frac{\log(1/j_1)}{\log(1/j_2)}$.

**Remark 1.** The parameters $K$ and $L$ differs from the usual analysis of the LSH framework. We have that $j_2^K \leqslant \frac{1}{\sqrt{n}}$ instead of $j_2^K \leqslant \frac{1}{n}$, and $n^{\rho/2} = L \geqslant 6j_1^K$ instead of $n^{\rho} = L \geqslant j_1^K$. The use of tensoring entails that most of the analysis will be centered around the sketches $\left(S'_{i,j}(A)\right)_{i\in[2M],j\in[L]}$. If we used the usual parameters then there would be $M\left\lceil\sqrt{L}\right\rceil$ sketches of size $\left\lceil\frac{K}{2}\right\rceil$. This would clutter the notation making the analysis harder to understand.

In order to create the $2M \cdot L$ sketches $\left(S'_{i,j}(A)\right)_{i\in[2M],j\in[L]}$ described above we first create $2M$ tables $T_0, \ldots, T_{2M-1}$ of size $L \times K$ such that for each $i \in [2M]$, $j \in [L]$ and $k \in [K]$ we have that $T_i[j,k]$ is an uniformly random integer chosen from the set $\{KSi + Sj, \ldots, KSi + S(j+1) - 1\}$ where $S$ is parameter to be chosen later. The tables are independent of each other. In every table the rows are chosen independently. Every row in every table is filled using a source of 2-independence. Now for a given $A \subseteq U$ we do the following

1. Let $S(A)$ be a size $2MKS$ similarity sketch of Section 2.

2. For each $i \in [2M]$, $j \in [L]$ and $k \in [K]$ let $S'_{i,j}(A)[k] = S(A)\left[T_i[j,k]\right]$.

See Figure 2 for an intuition about how the intermediate sketch $S(A)$ is partitioned. It takes $O(MKS\log MKS + |A|)$ time to create the sketch $S(A)$ by Lemma 1, and it takes $O(MLK)$ time to calculate $S'$. Thus the time needed to create sketches $\left(S'_{i,j}(A)\right)_{i\in[2M],j\in[L]}$ for any $A \subseteq U$ is $O\left(MLK + MKS\log MKS + |A|\right)$. We will set $S = \Theta(K/j_1)$. Now the total running time for creating the sketches $(S_{i,j}(A))_{i\in[2M],j\in[L^2]}$ becomes

$$O\left(ML^2 + MLK + MKS\log(MKS) + |A|\right) = O\left(n^{\rho}\log(1/\varepsilon) + |A|\right)$$

11

thus the running time of a query is $O(n^\rho \log(1/\varepsilon) + |A|)$.

The main issue is to prove that the intermediate sketch has the desired properties despite the entries not being independent. We call the vector of sketches $\left(S'_{i,j}(A)\right)_{j\in[L]}$ a subexperiment for each $i \in [2M]$, and the vector of sketches $(S_{i,j}(A))_{j\in[L^2]}$ an experiment for each $i \in [M]$. It is clear that the experiment $(S_{i,j}(A))_{j\in[L^2]}$ is completely determined by the subexperiments $\left(S'_{2i,j}(A)\right)_{j\in[L]}$ and $\left(S'_{2i+1,j}(A)\right)_{j\in[L]}$ for every $i \in [M]$ by construction. We define $\mathcal{F}_{bad}(Q) = \{A \in \mathcal{F} \mid J(A,Q) \leqslant j_2\}$, which we will call the family of bad sets with respect to $Q$. We also define $\mathcal{F}_{good}(Q) = \{A \in \mathcal{F} \mid J(A,Q) \geqslant j_1\}$, which we will call the family of good sets with respect to $Q$.

Furthermore it will be useful to define the vector $\mathcal{M}'_Q(A)[i]$ which will be the number of matches that the set $A \in \mathcal{F}$ has in the $i \in [2M]$ subexperiment:

$$\mathcal{M}'_Q(A)[i] = \sum_{j\in[L]} \left[S'_{i,j}(A) = S'_{i,j}(Q)\right],$$

and the vector $\mathcal{M}_Q(A)[i]$ which will be the number of matches that the set $A \in \mathcal{F}$ has in the $i \in [M]$ experiment:

$$\mathcal{M}_Q(A)[i] = \sum_{j\in[L^2]} [S_{i,j}(A) = S_{i,j}(Q)] = \mathcal{M}'_Q[2i]\mathcal{M}'_Q[2i+1].$$

Each experiment will provide us with a lot of matches and amongst all those matches we need to find a set $A$ with $J(A,Q) \geqslant j_2$, so we need to filter away all the bad sets. To do this each experiment will choose one candidate set amongst the matched sets. This will give us $O(\log(1/\varepsilon))$ candidates which we are allowed to use extra time checking by our time budget. We will use two different techniques to choose the candidates depending on the number of matches. If an experiment has $O(L)$ matches then we can afford to check each set using a sketch of size $O(\max(\log n, \log(1/\varepsilon)))$ since by using the minwise $b$-bit hashing trick of Li et al.[17] this can be done in $O(L)$ time. If an experiment has $\omega(L)$ matches then we pick a random match which also can be done in $O(L)$ time.

We note that the experiments are conditionally independent given the intermediate sketch, $S(Q)$, that is, if we fix the intermediate sketch then the experiments are independent. The goal is to show that the intermediate sketch satisfies the following with probability at least $1 - \frac{\varepsilon}{3}$: After fixing the intermediate sketch at least a constant fraction of the experiments have the properties: (a) we expect $O\left(j_2^{2K}n\right)$ bad matches in expectation, and (b) the probability that a good set $A^* \in \mathcal{F}_{good}(Q)$ is matched is at least $\Omega(j_1^{2K})$. To show this we need a bit of notation, for a set $A \in \mathcal{F}$ we define

$$Y_i(A) = \frac{\sum_{j\in[S]} [S(A)[iS+j] = S(B)[iS+j]]}{S}$$

for every $i \in [2MK]$, which corresponds to the block of the subexperiments, and for every $l \in [2M]$ we define $Z_l(A) = \prod_{i\in[K]} Y_{lK+i}(A)$. From these definition we get that

$$\Pr\left[S_{i,j}(A) = S_{i,j}(Q) \mid Z_{2l}(A), Z_{2l+1}(A)\right] = Z_{2i}(A) \cdot Z_{2i+1}(A).$$

If the intermediate sketch is fixed then

$$\mathrm{E}\left[\sum_{A\in\mathcal{F}} \mathcal{M}_Q(A)[i] \,\middle|\, Z_{2i}, Z_{2i+1}\right] = L^2 \sum_{A\in\mathcal{F}_{bad}} Z_{2i}(A) \cdot Z_{2i+1}(A)$$

So to prove (a) we need to show that for most $i \in [M]$ then $\sum_{A\in\mathcal{F}_{bad}(Q)} Z_{2i}(A) \cdot Z_{2i+1}(A) = O\left(j_2^{2K}n\right)$, which formalized in the next Lemma.

12

**Lemma 6.** *Let $I \subseteq [M]$ be a set of $d$ indices, and $C$ a constant, then*

$$\Pr\left[\bigwedge_{i \in I}\left(\sum_{A \in \mathcal{F}_{bad}(Q)} Z_{2i}(A) \cdot Z_{2i+1}(A) \geqslant C j_2^{2K} n\right)\right] \leqslant \left(\frac{e}{C}\right)^d .$$

*Proof.* Using Markov's inequality we note that

$$\Pr\left[\bigwedge_{i \in I}\left(\sum_{A \in \mathcal{F}_{bad}} Z_{2i}(A) \cdot Z_{2i+1}(A) \geqslant C j_2^{2K} n\right)\right] \leqslant \Pr\left[\prod_{i \in I}\sum_{A \in \mathcal{F}_{bad}} Z_{2i}(A) \cdot Z_{2i+1}(A) \geqslant \left(C j_2^{2K} n\right)^d\right]$$

$$\leqslant \frac{\mathrm{E}\left[\prod_{i \in I}\sum_{A \in \mathcal{F}_{bad}} Z_{2i}(A) \cdot Z_{2i+1}(A)\right]}{\left(C j_2^{2K} n\right)^d} ,$$

so we just need to show that

$$\mathrm{E}\left[\prod_{i \in I}\sum_{A \in \mathcal{F}_{bad}} Z_{2i}(A) \cdot Z_{2i+1}(A)\right] \leqslant \left(e j_2^{2K} n\right)^d .$$

To do this we enumerate $I = \{v_0, \ldots, v_{d-1}\}$ and consider any $A_0, \ldots, A_{d-1} \in \mathcal{F}_{bad}(Q)$.

$$\mathrm{E}\left[\prod_{i \in [d]} Z_{2v_i}(A_i) \cdot Z_{2v_i+1}(A_i)\right]$$

$$= \prod_{i \in [d]} \mathrm{E}\left[Z_{2v_i}(A_i) \cdot Z_{2v_i+1}(A_i) \,\Big|\, \sigma\left(\left(Z_{2v_j}, Z_{2v_j+1}\right)_{j \in [i-1]}\right)\right]$$

$$\leqslant \prod_{i \in [d]} \left(\frac{2MKS \cdot J(A_i, Q)}{2MKS - (i-1)2K}\right)^{2K}$$

$$\leqslant \left(\left(\frac{2MKS \cdot j_2}{2MKS - 2MK}\right)^{2K}\right)^d$$

$$\leqslant \left(\left(\frac{S \cdot j_2}{S - 1}\right)^{2K}\right)^d$$

$$\leqslant \left(e j_2^{2K}\right)^d$$

where the first inequality follows by using Lemma 4, and the last inequality follows by the definition of $S$ and $K$. Using this we see that

$$\mathrm{E}\left[\prod_{i \in I}\sum_{A \in \mathcal{F}_{bad}} Z_{2i}(A) \cdot Z_{2i+1}(A)\right] \leqslant |\mathcal{F}_{bad}(Q)|^d \cdot \left(e j_2^{2K}\right)^d \leqslant \left(e j_2^{2K} n\right)^d$$

which proves the result. $\qquad \square$

Since $\Pr\left[S_{i,j}(A) = S_{i,j}(Q) \,|\, Z_{2l}(A), Z_{2l+1}(A)\right] = Z_{2i}(A) \cdot Z_{2i+1}(A)$, then to show (b): We need that for most $i \in [M]$ then $Z_{2i}(A) \cdot Z_{2i+1}(A) \geqslant \Omega(j_1^{2K})$, which is formalized in the next lemma.

**Lemma 7.** *For every set of indices $I \subseteq [M]$ of size $d$ and real number $\delta \in [0, 1]$, we have that*

$$\Pr\left[\bigwedge_{i \in I}\left(Z_{2i}(A^*) \leqslant (1 - \delta)j_1^K \vee Z_{2i+1}(A^*) \leqslant (1 - \delta)j_1^K\right)\right] \leqslant \left(6 e^{-\delta^2 j_1 S/(32K)}\right)^d .$$

First we need the following technical lemma.

**Lemma 8.** *Let $t, k, s$ be positive integers such that $t = ks$ and let $\alpha$ be a real number such that $\alpha t$ is a positive integer. Let $\mathcal{B}$ be the set of all $t$-tuples from $\{0,1\}^t$ for which the sum of entries is exactly $\alpha t$, i.e.:*

$$\mathcal{B} = \left\{ (b_0, \ldots, b_t) \in \{0,1\}^t \,\middle|\, \sum_{i \in [t]} b_i = \alpha t \right\}$$

*Let $(a_0, \ldots, a_t)$ be drawn uniformly random from $\mathcal{B}$, and let $Y_i = \sum_{j \in [s]} a_{is+j}$ for every $i \in [k]$. Then for any real number $\delta \in [0,1]$:*

$$\Pr\left[ \prod_{i \in [k]} Y_i \leqslant (1-\delta)(\alpha s)^k \right] \leqslant 2e^{-\delta^2 \alpha s/(10k)}$$

*Proof.* First note that we can assume that $2e^{-\delta^2 \alpha s/(10k)} \leqslant 1$ which implies that $5k \leqslant \delta^2 \alpha s$.

Let $l$ be a positive integer to be chosen later. If $\prod_{i \in [k]} Y_i \leqslant (1-\delta)(\alpha s)^k$, then one of two events must be true: either there exists $i \in [k]$ such that $Y_i \leqslant \frac{\alpha s}{2}$ , or we have that

$$\prod_{i \in [k]} (Y_i + l)^{\underline{l}} \leqslant \prod_{i \in [k]} (Y_i + l)^l$$

$$\leqslant \left( \prod_{i \in [k]} \left(1 + \frac{l}{Y_i}\right)^l \right) \left( \prod_{i \in [k]} Y_i \right)^l$$

$$\leqslant \left( \prod_{i \in [k]} \left(1 + \frac{l}{Y_i}\right)^l \right) (1-\delta)^l (\alpha s)^{kl}$$

$$\leqslant e^{\left( \sum_{i \in [k]} \frac{1}{Y_i} \right) l^2 - \delta l} (\alpha s)^{kl}$$

$$\leqslant e^{\frac{2k}{\alpha s} l^2 - \delta l} (\alpha s)^{kl}$$

For any $i \in [k]$ the probability that $Y_i \leqslant \frac{\alpha s}{2}$ is at most $e^{-\alpha s/8}$ by standard Chernoff bound. Hence the probability that there exists such an $i$ is at most $ke^{-\alpha s/8}$ by union bound. Now we note that

$$ke^{-\alpha s/8} = k \left( e^{-\delta^2 \alpha s/(10k)} \right)^{10k/(8\delta^2)} \leqslant k \left( e^{-\delta^2 \alpha s/(10k)} \right)^k \leqslant e^{-\delta^2 \alpha s/(10k)}$$

where the last inequality comes from the fact that $xa^x \leqslant a$ if $a \leqslant \frac{1}{2}$ and $x$ is a positive integer.

Now we want to bound the other event. By using Markov's inequality we have that

$$\Pr\left[ \prod_{i \in [k]} (Y_i + l)^{\underline{l}} \leqslant e^{\frac{2k}{\alpha s} l^2 - \delta l} (\alpha s)^{kl} \right] = \Pr\left[ \left( \prod_{i \in [k]} (Y_i + l)^{\underline{l}} \right)^{-1} \geqslant e^{-\frac{2k}{\alpha s} l^2 + \delta l} (\alpha s)^{-kl} \right]$$

$$\leqslant \mathrm{E}\left[ \left( \prod_{i \in [k]} (Y_i + l)^{\underline{l}} \right)^{-1} \right] \cdot e^{\frac{2k}{\alpha s} l^2 - \delta l} (\alpha s)^{kl}$$

Now we want to show that

$$\mathrm{E}\left[ \left( \prod_{i \in [k]} (Y_i + l)^{\underline{l}} \right)^{-1} \right] \leqslant (\alpha s)^{-kl}$$

14

First we define

$$\mathcal{C} = \left\{ (c_0, \ldots, c_{k-1}) \in \{0, \ldots, s\}^k \ \middle| \ \sum_{i \in [k]} c_i = \alpha t \right\}$$

and note that for $(c_0, \ldots, c_{k-1}) \in \mathcal{C}$ the probability that $Y_i = c_i$ for all $i \in [k]$ is exactly

$$\binom{t}{\alpha t}^{-1} \prod_{i \in [k]} \binom{s}{c_i}$$

Thus the expected value is

$$\binom{t}{\alpha t}^{-1} \sum_{(c_0, \ldots, c_{k-1}) \in \mathcal{C}} \prod_{i \in [k]} \left( \binom{s}{c_i} \frac{1}{(c_i + l)^{\underline{l}}} \right)$$

Now we note that

$$\binom{s}{c_i} \frac{1}{(c_i + l)^{\underline{l}}} = \frac{1}{(s + l)^{\underline{l}}} \binom{s + l}{c_i + l}$$

Hence we get that

$$\binom{t}{\alpha t}^{-1} \sum_{(c_0, \ldots, c_{k-1}) \in \mathcal{C}} \prod_{i \in [k]} \left( \binom{s}{c_i} \frac{1}{(c_i + l)^{\underline{l}}} \right)$$

$$= \binom{t}{\alpha t}^{-1} \left( \frac{1}{(s + l)^{\underline{l}}} \right)^k \sum_{(c_0, \ldots, c_{k-1}) \in \mathcal{C}} \prod_{i \in [k]} \binom{s + l}{c_i + l}$$

We see that we are in fact summing over all $t$-tuples $(c_0, \ldots, c_{k-1})$ where $c_i \in \{l, \ldots, l + s\}$ and $\sum_{i \in [k]} c_i = \alpha t + kl$. So we define

$$\mathcal{C}' = \left\{ (c_0, \ldots, c_{k-1}) \in \{0, \ldots, s + l\}^k \ \middle| \ \sum_{i \in [k]} c_i = \alpha t + lk \right\}$$

and get that

$$\binom{t}{\alpha t}^{-1} \left( \frac{1}{(s + l)^{\underline{l}}} \right)^k \sum_{(c_0, \ldots, c_{k-1}) \in \mathcal{C}} \prod_{i \in [k]} \binom{s + l}{c_i + l}$$

$$\leqslant \binom{t}{\alpha t}^{-1} \left( \frac{1}{(s + l)^{\underline{l}}} \right)^k \sum_{(c_0, \ldots, c_{k-1}) \in \mathcal{C}'} \prod_{i \in [k]} \binom{s + l}{c_i}$$

$$= \binom{t}{\alpha t}^{-1} \left( \frac{1}{(s + l)^{\underline{l}}} \right)^k \binom{t + kl}{\alpha t + kl}$$

$$= \frac{(t + kl)^{\underline{kl}}}{(\alpha t + kl)^{\underline{kl}}} \left( \frac{1}{(s + l)^{\underline{l}}} \right)^k$$

It is clear that

$$\left( k^l (s + l)^{\underline{l}} \right)^k \geqslant (ks + kl)^{\underline{kl}} = (t + kl)^{\underline{kl}}$$

hence using this we get that

$$\frac{(t + kl)^{\underline{kl}}}{(\alpha t + kl)^{\underline{kl}}} \left( \frac{1}{(s + l)^{\underline{l}}} \right)^k \leqslant \frac{k^{kl}}{(\alpha t + kl)^{\underline{kl}}} \leqslant \frac{k^{kl}}{(\alpha t)^{kl}} = \frac{1}{(\alpha s)^{kl}}$$

15

So the probability is bounded by

$$e^{\frac{2k}{\alpha s}l^2 - \delta l}$$

which has global minimum when $l = \frac{\delta \alpha s}{4k}$. Since $l$ is an integer we get that the probability is bounded by

$$e^{\frac{2k}{\alpha s}\left(\frac{\delta \alpha s}{4k} \pm \frac{1}{2}\right)^2 - \delta\left(\frac{\delta \alpha s}{4k} \pm \frac{1}{2}\right)} = e^{-\delta^2 \alpha s/(8k) + k/(2\alpha s)}$$

Now using that $5k \leqslant \delta^2 \alpha s \leqslant \alpha s$ we get that $k/(2\alpha s) \leqslant \delta^2 \alpha s/(50k)$. We see that $-\delta^2 \alpha s/(8k) + \delta^2 \alpha s/(50k) \leqslant -\delta^2 \alpha s/(10k)$ which finishes the proof. $\square$

We can now prove that most $Z_i(A^*) \geqslant (1-\delta)j_1^K$, which will almost prove Lemma 7.

**Lemma 9.** *For every set of indices $I \subseteq [2M]$ of size $d$ and real number $\delta \in [0,1]$, we have that*

$$\Pr\left[\bigwedge_{i \in I} Z_i(A^*) \leqslant (1-\delta)j_1^K\right] \leqslant \left(3e^{-\delta^2 j_1 S/(32K)}\right)^d$$

*Proof.* The first we do is to write the probability as an expectation

$$\Pr\left[\bigwedge_{l \in I} Z_l \leqslant (1-\delta)j_1^K\right] = \mathrm{E}\left[\prod_{l \in I}\left[Z_l \leqslant (1-\delta)j_1^K\right]\right]$$

For every $l \in I$ we define the random variable $\alpha_l$ by

$$\alpha_l = \frac{\sum_{i \in [K]} Y_{lK+i}}{K}$$

and let $P_l(\alpha)$ be the probability that $Z_l \leqslant (1-\delta)j_1^K$ given that $\alpha_l = \alpha$. We then get that

$$\mathrm{E}\left[\prod_{l \in I}\left[Z_l \leqslant (1-\delta)j_1^K\right]\right] = \mathrm{E}\left[\mathrm{E}\left[\prod_{l \in I}\left[Z_l \leqslant (1-\delta)j_1^K\right]\,\bigg|\,\sigma\left((\alpha_l)_{l \in I}\right)\right]\right]$$

$$= \mathrm{E}\left[\prod_{l \in I} P_l(\alpha_l)\right]$$

Now we will split the random variable $P_l(\alpha_l)$ as follows

$$P_l(\alpha_l) = \left[\alpha_l \leqslant \left(1 - \frac{\delta}{4K}\right)j_1\right]P_l(\alpha_l) + \left[\alpha_l > \left(1 - \frac{\delta}{4K}\right)j_1\right]P_l(\alpha_l)$$

If $\alpha_l > \left(1 - \frac{\delta}{4K}\right)j_1$ then $\alpha_l^K > \left(1 - \frac{\delta}{4}\right)j_1^K$ so $\left(1 - \frac{3}{4}\delta\right)\alpha_l^K > (1-\delta)j_1^K$. Using Lemma 8 we get that

$$\left[\alpha_l > \left(1 - \frac{\delta}{4K}\right)j_1\right]P_l(\alpha_l) \leqslant \left[\alpha_l > \left(1 - \frac{\delta}{4K}\right)j_1\right]2e^{-\delta^2 \alpha_l S 3^2/(4^2 \cdot 10K)}$$

$$\leqslant 2e^{-\delta^2 j_1 S 3^3/(4^3 \cdot 10K)}$$

$$\leqslant 2e^{-\delta^2 j_1 S/(32K)}$$

where the second inequality uses that $\alpha_l > \left(1 - \frac{\delta}{4K}\right)j_1 \geqslant \frac{3}{4}j_1$ and last inequality uses that $\frac{3^3}{4^3 \cdot 10} > \frac{1}{32}$.

Now using this we get that

$$\mathrm{E}\left[\prod_{l\in I}P_l(\alpha_l)\right]\leqslant\mathrm{E}\left[\prod_{l\in I}\left(\left[\alpha_l\leqslant\left(1-\frac{\delta}{4K}\right)j_1\right]+2e^{-\delta^2 j_1 s/(32k)}\right)\right]$$

Now for every subset $I'\subseteq I$ of size $d'$ we have that

$$\mathrm{E}\left[\prod_{l\in I'}\left[\alpha_l\leqslant\left(1-\frac{\delta}{4K}\right)j_1\right]\right]\leqslant\Pr\left[\frac{\sum_{l\in I'}\alpha_l}{d'}\leqslant\left(1-\frac{\delta}{4K}\right)j_1\right]$$

$$=\Pr\left[\frac{\sum_{l\in I'}\sum_{i\in[K]}Y_{lK+i}}{d't}\leqslant\left(1-\frac{\delta}{4K}\right)j_1\right]$$

$$\leqslant e^{-\delta^2 j_1 d'S/(32K)}$$

$$=\left(e^{-\delta^2 j_1 S/(32K)}\right)^{d'}$$

where the inequality uses that Lemma 5. Hence we get that

$$\Pr\left[\bigwedge_{l\in I}Z_l\leqslant(1-\delta)j_1^K\right]\leqslant\left(3e^{-\delta^2 j_1 S/(32K)}\right)^d$$

as we wanted. $\qquad\square$

We are now ready to prove Lemma 7 which follows easily from Lemma 9.

*Proof of Lemma 7.* By union bound and Lemma 9 we get that

$$\Pr\left[\bigwedge_{i\in I}\left(Z_{2i}(A^*)\leqslant(1-\delta)j_1^K\vee Z_{2i+1}(A^*)\leqslant(1-\delta)j_1^K\right)\right]$$

$$\leqslant\sum_{(v_i)_{i\in I}\in\{0,1\}^d}\Pr\left[\bigwedge_{i\in I}\left(Z_{2i+v_i}(A^*)\leqslant(1-\delta)j_1^K\right)\right]$$

$$\leqslant\sum_{(v_i)_{i\in I}\in\{0,1\}^d}\left(3e^{-\delta^2 j_1 S/(32K)}\right)^d$$

$$=\left(6e^{-\delta^2 j_1 S/(32K)}\right)^d$$

$\qquad\square$

Having Lemma 6 and Lemma 7 it becomes easy to prove the main theorem of this section.

**Theorem 3.** *Let $A\in\mathcal{F}_{good}(Q)$ and, $M_{good}\subseteq[M]$ be the set of experiments, such that, for every $i\in M_{good}$*

    *1. $\sum_{A\in\mathcal{F}_{bad}(Q)}Z_{2i}(A)\cdot Z_{2i+1}(A)\leqslant 4ej_2^{2K}n$   ,*

    *2. $Z_{2i}(A^*)\geqslant\frac{1}{2}j_1^K\wedge Z_{2i+1}(A^*)\geqslant\frac{1}{2}j_1^K$   .*

*Then $|M_{good}|\geqslant\frac{M}{3}$ with probability at least $1-\frac{\varepsilon}{2}$.*

*Proof.* By Lemma 6, Lemma 7, and a standard analysis of Chernoff bounds, we have Chernoff-type bounds on the number of experiments which does not satisfy either of the requirements. Now choosing $S$ large enough we get that $6e^{-\frac{1}{2}^2 j_1 S/(32K)}\leqslant\frac{1}{4}$, so choosing $M$ large enough at most $\frac{M}{3}$ experiments does not satisfy the first requirement with probability at least $1-\frac{\varepsilon}{4}$, and at most $\frac{M}{3}$ experiments does not satisfy the second requirement with probability at least $1-\frac{\varepsilon}{4}$. Now a union bound finishes the proof. $\qquad\square$

## 3.2 Filtering

Using Theorem 3, we get that good set $A^* \in \mathcal{F}_{good}(Q)$ is matched with probability at least $1 - \frac{3}{4}\varepsilon$. Now it is not enough to know that a good set $A^* \in \mathcal{F}_{good}(Q)$ is matched in one of the experiments with high probability, we also need to be able to find the good set amongst all the matched sets. To do this we will apply a two-tiered filtering algorithm. First we use a fast and imprecise filtering step to select a candidate $C_i \in \mathcal{F}$ amongst the matched set in the $i$'th experiment for every experiment $i \in [M]$. This will give us $M$ candidates $C_0, \ldots, C_{M-1}$ for which will use a slower and more precise filtering step to select the actual set. We will show that this two-tiered filtering algorithm finds a good set with high probability if such a set exists.

In the first filtering step we will use two different approaches depending on the number of matches in an experiment. If there is more than $2CL^2$ matches in a experiment then we will choose a random match. If there is less than $2CL^2$ matches then we will check every match using a sketch of size $\Theta(\max(\log n, \log(1/\varepsilon)))$, this can be done in constant time per element by using minwise $b$-bit hashing, and pick the first element that is above a certain threshold. If no element is above the threshold then we will for the sake of the analysis assume that the candidate picked is the empty set.

In the second filtering step we have a set of $M = O(\log(1/\varepsilon))$ candidates, which we will check using a sketch of size $\Theta(\log^2(n)\log(1/\varepsilon))$, again using minwise $b$-bit hashing this can be done in $\Theta(\log^2(n))$ time per element. This allows us distinguish between elements with similarity less than $j_2$ and elements with similarity at least $j_2'$. Again we pick the first element that is above a certain threshold.

We define

$$\mathcal{M}_Q(A)[i] = \sum_{j \in [L]} [S_{i,j}(A) = S_{i,j}(Q)]$$

to be the number of matches of $A \in \mathcal{F}$ in the $i$'th experiment where $i \in [M]$. Now let

$$I_Q(A^*) = \left\{ i \in M_{good} \,\middle|\, (\mathcal{M}_Q(A^*)[i] > 0) \wedge \left( \sum_{A \in \mathcal{F}_{bad}(Q)} \mathcal{M}_Q(A)[i] \leqslant CL^2 \right) \right\}$$

be the set of experiments where $A^* \in \mathcal{F}_{good}(Q)$ is matched and where there are not to many bad sets are matched. Using Theorem 3 and the analysis from Dahlgaard et al.[10] we get that $|I_Q(A^*)| = \Omega(\log(1/\varepsilon))$ with probability at least $1 - \frac{3}{4}\varepsilon$ by choosing $M = \Theta(\log(1/\varepsilon))$ and $C = \Theta(1)$ large enough.

We want to show that at least one of the experiments chooses a candidate with similarity at least $j_2'$, in particular we will show that

$$\Pr\left[ \bigwedge_{i \in I_Q(A^*)} \left( J(C_i, Q) < j_2' \right) \right] \leqslant \left( \frac{1}{2} \right)^{|I_Q(A^*)|}$$

Since we will filter in two different ways depending on the number of matches, we will split $I_Q(A^*)$ into two: The experiments $I_Q'(A^*)$ with many matches

$$I_Q'(A^*) = \left\{ i \in I_Q(A^*) \,\middle|\, \sum_{A \in \mathcal{F}} \mathcal{M}_Q(A)[i] > 2CL^2 \right\}$$

and the experiments $I_Q''(A^*)$ with few matches

$$I_Q''(A^*) = \left\{ i \in I_Q(A^*) \,\middle|\, \sum_{A \in \mathcal{F}} \mathcal{M}_Q(A)[i] \leqslant 2CL^2 \right\}$$

First we consider the case where the is at least one good experiment with few matches. As mentioned we will create a sketch $T(Q)$ of size $t = \Theta(\max(\log n, \log(1/\varepsilon)))$ using the minwise 1-bit hashing trick. We then get that

$$\Pr[T(Q)[i] = T(A)[i]] = J(A,Q) + (1 - J(A,Q))\frac{1}{2} = \frac{1 + J(A,Q)}{2}$$

for any $A \in \mathcal{F}$ and $i \in [t]$. We pick the threshold $\gamma = \frac{1 + \frac{j_1 + j_2'}{2}}{2}$. Using Hoeffding's inequality we get that for $A \in \mathcal{F}_{good}(Q)$ then

$$\Pr\left[\sum_{i \in [t]} [T(A)[i] = T(Q)[i]] \leq \gamma t\right] \leq \frac{\varepsilon}{8} \cdot \frac{1}{n}$$

and for $A \in \mathcal{F}_{bad}(Q)$ then

$$\Pr\left[\sum_{i \in [t]} [T(A)[i] = T(Q)[i]] \geq \gamma t\right] \leq \frac{\varepsilon}{8} \cdot \frac{1}{n}$$

So a union bound over all the sets in all the experiments, for which there are at most $O(CL^2 \log(1/\varepsilon))$, shows that the probability that none of the candidates has similarity at least $j_2'$ is at most $\frac{\varepsilon}{8}$. It takes $O(1)$ time to check each set so it takes $O(L^2)$ time to filter each experiment.

Now we assume that every good experiment has many matches. We note that since we pick a random element in every experiment independently then we get that

$$\Pr\left[\bigwedge_{i \in I_Q'(A^*)} \left(J(C_i, Q) < j_2'\right)\right] = \prod_{i \in I_Q'(A^*)} \Pr\left[J(C_i, Q) < j_2'\right]$$

Using Markov's inequality we get that

$$\Pr\left[J(C_i, Q) < j_2'\right] \leq \frac{1}{2}$$

for $i \in I'Q(A^*)$ since we know that $\sum_{A \in \mathcal{F}_{bad}(Q)} \mathcal{M}_Q(A)[i] \leq CL^2$. This shows that

$$\Pr\left[\bigwedge_{i \in I_Q'(A^*)} \left(J(C_i, Q) < j_2'\right)\right] \leq \left(\frac{1}{2}\right)^{|I_Q'(A^*)|}$$

Since every experiment has many matches this implies that the probability that none of the candidates has similarity at least $j_2'$ is at most $\frac{\varepsilon}{8}$. This shows that the probability that the first step of filtering fails is at most $\frac{\varepsilon}{8}$.

We have now shown that the probability, that none of the candidates has Jaccard similarity at least $j_2'$, is at most $\frac{7}{8}\varepsilon$. Now to find the correct candidate we will create a sketch $T'(Q)$ of size $t' = \Theta(\log^2(n)\log(1/\varepsilon))$ using the minwise 1-bit hashing trick. We pick the threshold $\gamma' = \frac{1 + \frac{j_2 + j_2'}{2}}{2}$. Using Hoeffding's inequality we get that for candidates $C$ with $J(C,Q) \geq j_2'$ then

$$\Pr\left[\sum_{i \in [t]} [T(C)[i] = T(Q)[i]] \leq \gamma' t'\right] \leq \frac{\varepsilon^2}{8}$$

and for candidates $C$ with $J(C, Q) \leqslant j_2$ then

$$\Pr\left[\sum_{i \in [t]} [T(C)[i] = T(Q)[i]] \geqslant \gamma' t'\right] \leqslant \frac{\varepsilon^2}{8}$$

A union bound over all $O(\log(1/\varepsilon))$ candidates shows that the probability, that the set chosen has similarity less than $j_2$, is at most $\frac{\varepsilon}{8}$. Checking each candidate takes $O(\log^2(n))$ time so we use a total of $O(\log^2(n) \log(1/\varepsilon))$ time for second filtering step.

Combining all the steps shows that the data structure finds a set $A \in \mathcal{F}$ with $J(A, Q) \geqslant j_2$ if there exists a set $B \in \mathcal{F}$ with $J(B, Q) \geqslant j_1$ with probability at least $1 - \varepsilon$. Now since we want the data structure to only have a 1-sided error then we calculate the exact Jaccard similarity of the set in $O(|Q|)$ time.

# References

[1] Alexandr Andoni and Piotr Indyk. Efficient algorithms for substring near neighbor problem. In *Proceedings of the Seventeenth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2006, Miami, Florida, USA, January 22-26, 2006*, pages 1203–1212. ACM Press, 2006.

[2] Yoram Bachrach and Ely Porat. Sketching for big data recommender systems using fast pseudo-random fingerprints. In Fedor V. Fomin, Rusins Freivalds, Marta Z. Kwiatkowska, and David Peleg, editors, *Automata, Languages, and Programming - 40th International Colloquium, ICALP 2013, Riga, Latvia, July 8-12, 2013, Proceedings, Part II*, volume 7966 of *Lecture Notes in Computer Science*, pages 459–471. Springer, 2013.

[3] Roberto J. Bayardo, Yiming Ma, and Ramakrishnan Srikant. Scaling up all pairs similarity search. In *Proceedings of the 16th International Conference on World Wide Web*, WWW '07, page 131–140, New York, NY, USA, 2007. Association for Computing Machinery.

[4] A. Broder. On the resemblance and containment of documents. In *Proceedings of the Compression and Complexity of Sequences 1997*, SEQUENCES '97, pages 21–29, Washington, DC, USA, 1997. IEEE Computer Society.

[5] Andrei Z. Broder, Moses Charikar, Alan M. Frieze, and Michael Mitzenmacher. Min-wise independent permutations. *J. Comput. Syst. Sci.*, 60(3):630–659, 2000. Announced at STOC'98.

[6] Andrei Z. Broder, Steven C. Glassman, Mark S. Manasse, and Geoffrey Zweig. Syntactic clustering of the web. *Computer Networks*, 29(8-13):1157–1166, 1997.

[7] Tobias Christiani. Fast locality-sensitive hashing for approximate near neighbor search. *CoRR*, abs/1708.07586, 2017.

[8] Edith Cohen and Haim Kaplan. Summarizing data using bottom-k sketches. In *Proc. 26th PODC*, pages 225–234, 2007.

[9] Søren Dahlgaard, Mathias Bæk Tejs Knudsen, Eva Rotenberg, and Mikkel Thorup. Hashing for statistics over k-partitions. In Venkatesan Guruswami, editor, *IEEE 56th Annual Symposium on Foundations of Computer Science, FOCS 2015, Berkeley, CA, USA, 17-20 October, 2015*, pages 1292–1310. IEEE Computer Society, 2015.

[10] Søren Dahlgaard, Mathias Bæk Tejs Knudsen, and Mikkel Thorup. Fast similarity sketching. In Chris Umans, editor, *58th IEEE Annual Symposium on Foundations of Computer Science, FOCS 2017, Berkeley, CA, USA, October 15-17, 2017*, pages 663–671. IEEE Computer Society, 2017.

[11] 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.

[12] Monika Henzinger. Finding near-duplicate web pages: A large-scale evaluation of algorithms. In *Proceedings of the 29th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '06, page 284–291, New York, NY, USA, 2006. Association for Computing Machinery.

[13] Jakob Bæk Tejs Houen and Mikkel Thorup. Understanding the moments of tabulation hashing via chaoses. In *49th International Colloquium on Automata, Languages, and Programming, ICALP 2022, July 4-8, 2022, Paris, France*, volume 229 of *LIPIcs*, pages 74:1–74:19. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2022.

[14] Piotr Indyk and Rajeev Motwani. Approximate nearest neighbors: Towards removing the curse of dimensionality. In Jeffrey Scott Vitter, editor, *Proceedings of the Thirtieth Annual ACM Symposium on the Theory of Computing, Dallas, Texas, USA, May 23-26, 1998*, pages 604–613. ACM, 1998.

[15] Ping Li. 0-bit consistent weighted sampling. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 665–674, 08 2015.

[16] Ping Li, Art B. Owen, and Cun-Hui Zhang. One permutation hashing. In Peter L. Bartlett, Fernando C. N. Pereira, Christopher J. C. Burges, Léon Bottou, and Kilian Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25: 26th Annual Conference on Neural Information Processing Systems 2012. Proceedings of a meeting held December 3-6, 2012, Lake Tahoe, Nevada, United States.*, pages 3122–3130, 2012.

[17] Ping Li, Anshumali Shrivastava, Joshua L. Moore, and Arnd Christian König. Hashing algorithms for large-scale learning. In John Shawe-Taylor, Richard S. Zemel, Peter L. Bartlett, Fernando C. N. Pereira, and Kilian Q. Weinberger, editors, *Advances in Neural Information Processing Systems 24: 25th Annual Conference on Neural Information Processing Systems 2011. Proceedings of a meeting held 12-14 December 2011, Granada, Spain.*, pages 2672–2680, 2011.

[18] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. *Introduction to information retrieval.* Cambridge University Press, 2008.

[19] Ryan O'Donnell, Yi Wu, and Yuan Zhou. Optimal lower bounds for locality-sensitive hashing (except when q is tiny). *TOCT*, 6(1):5:1–5:13, 2014.

[20] Mihai Pătraşcu and Mikkel Thorup. On the k-independence required by linear probing and minwise independence. In Samson Abramsky, Cyril Gavoille, Claude Kirchner, Friedhelm Meyer auf der Heide, and Paul G. Spirakis, editors, *Automata, Languages and Programming*, pages 715–726, Berlin, Heidelberg, 2010. Springer Berlin Heidelberg.

[21] Jeanette P. Schmidt, Alan Siegel, and Aravind Srinivasan. Chernoff-hoeffding bounds for applications with limited independence. *SIAM J. Discrete Math.*, 8(2):223–250, 1995.

[22] Gregory Shakhnarovich, Trevor Darrell, and Piotr Indyk. *Nearest-Neighbor Methods in Learning and Vision: Theory and Practice (Neural Information Processing).* The MIT Press, 2006.

[23] Anshumali Shrivastava. Optimal densification for fast and accurate minwise hashing. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, volume 70 of *Proceedings of Machine Learning Research*, pages 3154–3163. PMLR, 2017.

[24] Anshumali Shrivastava and Ping Li. Densifying one permutation hashing via rotation for fast near neighbor search. In *Proceedings of the 31th International Conference on Machine Learning, ICML 2014, Beijing, China, 21-26 June 2014*, volume 32 of *JMLR Workshop and Conference Proceedings*, pages 557–565. JMLR.org, 2014.

[25] Anshumali Shrivastava and Ping Li. Improved densification of one permutation hashing. In Nevin L. Zhang and Jin Tian, editors, *Proceedings of the Thirtieth Conference on Uncertainty in Artificial Intelligence, UAI 2014, Quebec City, Quebec, Canada, July 23-27, 2014*, pages 732–741. AUAI Press, 2014.

[26] Mikkel Thorup. Bottom-k and priority sampling, set similarity and subset sums with minimal independence. In *Proceedings of the Forty-Fifth Annual ACM Symposium on Theory of Computing*, STOC '13, page 371–380, New York, NY, USA, 2013. Association for Computing Machinery.