

Near-Duplicate Text Alignment under Weighted Jaccard Similarity

Yuheng Zhang¹ Miao Qiao² Zhencan Peng¹ Dong Deng^{1*}

¹Rutgers University ²University of Auckland

{yuheng.zhang,zhencan.peng,dong.deng}@cs.rutgers.edu,miao.qiao@auckland.edu

Abstract

Near-duplicate text alignment is the task of identifying all subsequences (i.e., substrings) in a collection of texts that are similar to a given query. Traditional approaches rely on seeding–extension–filtering heuristics, which lack accuracy guarantees and require many hard-to-tune parameters. More recent methods leverage min-hash techniques. They propose to group all the subsequences in each text by their min-hash and index the groups. When a query arrives, they can use the index to find all the min-hash sketches that are similar to the query’s sketch and then return the corresponding subsequences as the results efficiently. Thus these methods guarantee to identify all subsequences whose estimated Jaccard similarity with the query exceed a user-provided threshold. However, these methods only support unweighted Jaccard similarity, which cannot capture token importance or frequency, limiting their effectiveness in real-world scenarios where tokens carry weights, such as TF-IDF.

In this paper, we address this limitation by supporting weighted Jaccard similarity using consistent weighted sampling. We design an algorithm `MonoActive` to group all subsequences in a text by their consistent weighted sampling. We analyze the complexity of our algorithm. For raw count term frequency (where a token’s weight is proportional to its frequency in the text), we prove `MonoActive` generates $O(n + n \log f)$ groups (each group occupies $O(1)$ space) in expectation for a text with n tokens, where f is the maximum token frequency in the text. We further prove that our algorithm is optimal, meaning that any algorithm must generate $\Omega(n + n \log f)$ groups in expectation. Extensive experiments show that `MonoActive` outperforms the state-of-the-art by up to 4.7× in speed and reduces index size by up to 30%, with superior scalability.

PVLDB Reference Format:

Yuheng Zhang¹ Miao Qiao² Zhencan Peng¹ Dong Deng¹.
Near-Duplicate Text Alignment under Weighted Jaccard Similarity. PVLDB,
19(9): XXX-XXX, 2026.
doi:XX.XX/XXX.XX

PVLDB Artifact Availability:

The source code, data, and/or other artifacts have been made available at
<https://github.com/Yveh/OptAlign>.

*Dong Deng is the corresponding author.

This work is licensed under the Creative Commons BY-NC-ND 4.0 International License. Visit <https://creativecommons.org/licenses/by-nc-nd/4.0/> to view a copy of this license. For any use beyond those covered by this license, obtain permission by emailing info@vldb.org. Copyright is held by the owner/author(s). Publication rights licensed to the VLDB Endowment.

Proceedings of the VLDB Endowment, Vol. 19, No. 9 ISSN 2150-8097.
doi:XX.XX/XXX.XX

1 Introduction

This paper studies the near-duplicate text alignment problem [15, 33, 34, 37, 38, 52]. Given a collection of data texts, the task takes a short query text and returns all subsequences (i.e., substrings) of the data texts that are similar to the query. Near-duplicate text alignment has become increasingly important in the era of large language models (LLMs), with key applications in test set leakage (also known as data contamination) detection [29], training data deduplication [27], and memorization analysis [9, 50]. Beyond LLMs, this problem also plays a crucial role in domains such as bioinformatics [2], log analysis [13], and plagiarism detection [35, 39].

Due to the high computational cost of near-duplicate text alignment, previous methods often adopt the seeding–extension–filtering heuristic [3, 5, 7, 22, 23, 26, 31, 35, 42, 43, 51]. However, these methods lack accuracy guarantees and often involve many hard-to-tune hyperparameters [17]. To address these limitations, recent studies have proposed to use min-hash techniques [6, 28] for near-duplicate text alignment. These approaches guarantee to retrieve all subsequences whose estimated Jaccard similarities with the query exceed a user-defined threshold [15, 33, 34, 52]. However, a key limitation is that they only support the unweighted Jaccard similarity, which treats all tokens equally—regardless of their frequency or importance (note depending on the tokenizer, a token can be a word [53], a q -gram [18], a byte-pair-encoding [19], etc.).

To highlight the issue of unweighted Jaccard similarity, consider $Q = \text{AAAAAATTTTTC}$, $T = \text{AAAAAATTTTG}$, and $S = \text{AATTGCC}$. Intuitively, Q is much more similar to T than to S . Yet, under unweighted Jaccard similarity (with each token as a 2-gram), both Q and T , as well as Q and S yield the same similarity score: $4/7$. In contrast, weighted Jaccard similarity correctly reflects the difference: when each token in a text is weighted by its frequency in the text, Q and T have a similarity of 0.8, while Q and S only score 0.2. Moreover, unweighted Jaccard similarity fails to distinguish between stop words and content words. For example, consider $Q = \text{“I read about Einstein in a book”}$, $T = \text{“I studied Einstein through a book”}$, $S = \text{“I roamed about in a castle”}$. Semantically, Q is clearly closer to T than to S . However, the unweighted Jaccard similarity is the same for both: $4/9$ (with each token as a word). If a meaningful token weighting scheme such as TF-IDF [48] is applied, then stop words, e.g., *I*, *in*, *a*, *about*, *through*, could have near-zero weight and the other words may have near-one weight, and thereby drawing clear distinctions. Under this scheme, the weighted Jaccard similarity between Q and T is around 0.5, while that between Q and S drops to nearly zero.

Existing Methods. Existing methods propose to group, as the indexing process, the subsequences in each data text based on their min-hash values. This way, when a query arrives, they use the index

to identify all the min-hash sketches that are similar to the query’s min-hash sketch and then return the corresponding subsequences as the results efficiently [15, 33, 34]. A key observation made by these methods is that the nearby subsequences of a text tend to share the same min-hash (this is because appending a token to a subsequence most likely would not change the min-hash of the subsequence). Given a text with n tokens, they group the $O(n^2)$ subsequences in the text into $O(n)$ groups while represent each group with a tuple of $O(1)$ space.

Limitations of Existing Methods. In many real-world scenarios, tokens are associated with weights, such as TF-IDF (term frequency-inverse document frequency) weights [32]. Estimating weighted Jaccard similarity requires consistent weighted sampling (CWS) [24], as standard min-hash[6] applies only to the unweighted case. Existing methods do not support CWS or weighted Jaccard similarity. The only exception is AllAlign [15] which supports multi-set Jaccard, a special case of weighted Jaccard where each token’s weight is its frequency in the text. AllAlign is a greedy algorithm which groups the subsequences by their “multi-set min-hash”. However, it lacks a formal complexity analysis, including an upper bound on the number of groups generated and its time complexity. Due to its recursive nature, analyzing its complexity is particularly difficult.

Our Approach. In this paper, we propose an efficient algorithm MonoActive to group the subsequences of a text by their CWS or multi-set min-hash. We rigorously analyze the complexity for MonoActive. Specifically, for multi-set Jaccard similarity, we prove that MonoActive produces $O(n + n \log f)$ groups for a text with n tokens in expectation, where f is the maximum token frequency in the text. Each group occupies only $O(1)$ space. Furthermore, we prove that MonoActive is *worst-case optimal* by presenting a lower bound analysis on the number of groups. That is, there exists text instances for which any algorithm under the hash-based framework must generate at least $\Omega(n + n \log f)$ groups in expectation. We further develop an optimization that improves the time complexity of our algorithm from $O(nf \log n)$ to $O(n \log n + n \log n \log f)$ and space complexity from $O(nf)$ to $O(n + n \log f)$. Finally, we show our algorithm can be generalized to group the subsequences in a text based on their consistent weighted samplings, as long as the weight of a token in a text is monotonically increasing with its frequency in the text and is independent of other properties of the text. For example, for logarithmic term frequency (where the weight of a token t with frequency f_t is proportional to $\log(f_t + 1)$), MonoActive generates $O(n + n \log \log f)$ groups in expectation.

In summary, we make the following contributions in this paper.

- We develop MonoActive, the first algorithm for near-duplicate text alignment under weighted Jaccard similarity.
- We rigorously analyze the complexity of MonoActive and design optimizations to reduce its time and space complexities.
- For the special case of multi-set Jaccard similarity, we prove MonoActive is optimal, while the greedy algorithm AllAlign lacks theoretical guarantees.
- For multi-set Jaccard similarity, experimental results show that MonoActive outperforms AllAlign by up to 26× in index construction time, reduces index size by up to 30%, and improves query latency by up to 3×. The performance gain increases as the text length n grows, exhibiting superior scalability.

The rest of the paper is organized as follows. We introduce preliminary knowledge in Section 2 and the framework in Section 3. Section 4 presents our grouping algorithm and Section 5 extends it to weighted Jaccard similarity. Section 6 show experiment results, Section 7 reviews related work, and Section 8 concludes the paper.

2 Preliminaries

We first consider a special case where the tokens are weighted by their frequencies in the text (i.e., term frequency). In this case, the weighted Jaccard similarity degrades to multi-set Jaccard similarity.

2.1 Multi-Set Jaccard Similarity

We first define notations that will be used in the paper. A text T is a sequence of tokens, where $T[i]$ is the i -th token in the text. We define $f(t, T)$ as the frequency of a token t in T , i.e., the number of occurrences of t in T . A text T can be uniquely mapped to a set which exclusively contains all the distinct tokens of T as elements. For ease of presentation, we refer to T both as a text and a set. We use $|T|$ to denote the length of the text T . Furthermore, we use $T[i, j]$ to represent the subsequence of T from the i -th token to the j -th token (inclusive), where $1 \leq i \leq j \leq |T|$. Note all the definitions naturally extend to subsequences. Thus $T[i, j]$ is both a subsequence and a token set. Given two texts T and S , their *multi-set Jaccard similarity* is

$$J_{T,S} = \frac{\sum_{t \in T \cup S} \min(f(t, T), f(t, S))}{\sum_{t \in T \cup S} \max(f(t, T), f(t, S))}.$$

EXAMPLE 1. Consider two texts $T = ABBC$ and $S = BCD$ where each letter denotes a token. The union $T \cup S = \{A, B, C, D\}$. The token frequencies are $f(A, T) = 1, f(B, T) = 2, f(C, T) = 1, f(D, T) = 0, f(A, S) = 0, f(B, S) = 1, f(C, S) = 1, \text{ and } f(D, S) = 1$. Therefore, their multi-set Jaccard similarity is $J_{T,S} = \frac{2}{5}$.

2.2 Min-Hash for Multi-Set Jaccard Similarity

The multi-set Jaccard similarity of two texts can be efficiently and accurately estimated by their *min-hash* sketches. Specifically, let $h(t, x)$ be a random universal hash function that takes a token t and a positive integer x as input and outputs a non-negative integer hash value. The (multi-set) min-hash of a text T is

$$h(T) = \min\{h(t, x) \mid t \in T, 1 \leq x \leq f(t, T)\}. \quad (1)$$

EXAMPLE 2. Consider the text T and the hash function h from the running example in the caption of Fig. 1, we have $h(T) = h(B, 4) = 1$.

Given a random universal hash function h and two texts T and S , the probability that they share the same min-hash is equivalent to their multi-set Jaccard similarity. Formally, we have

$$\Pr(h(T) = h(S)) = J_{T,S}.$$

This is because there are $\sum_{t \in T \cup S} \max(f(t, T), f(t, S))$ unique hash values and $\sum_{t \in T \cup S} \min(f(t, T), f(t, S))$ common hash values in the two texts T and S . T and S share the same min-hash if and only if the smallest hash value of all unique hash values is among their common hash values. The probability of the latter is exactly $J_{T,S}$.

Min-Hash Sketch. With k independent random universal hash functions h_1, \dots, h_k , the min-hash sketch of a text T consists of k

min-hash values $h_1(T), \dots, h_k(T)$. The multi-set Jaccard similarity of two texts T and S can be unbiasedly estimated by

$$\hat{J}_{T,S} = \frac{1}{k} \sum_{i=1}^k \mathbf{1}\{h_i(T) = h_i(S)\} \quad (2)$$

where $\mathbf{1}$ is an indicator function [6].

Implementation Details. The hash function $h(t, x)$ can be drawn from a family \mathcal{H} of universal hash functions $h(t, x) = (a_1 t + a_2 x + b) \bmod p$ where p is a large prime and a_1, a_2, b are randomly chosen integers module p with $a_1 \neq 0, a_2 \neq 0$. We assume there is no hash collision under the universal hash function throughout the paper.

2.3 Near-Duplicate Text Alignment

Near-duplicate text alignment is formally defined as below.

DEFINITION 1. Given a collection of texts D , a query Q , and a similarity threshold $\theta \in [0, 1]$, the near-duplicate text alignment problem returns all the subsequences $T[i, j]$, where $T \in D$ and $1 \leq i \leq j \leq |T|$, such that $\hat{J}_{Q,T[i,j]} \geq \theta$.

EXAMPLE 3. Consider a collection of two texts $D = \{T, S\}$, where $T = \text{ABBCDE}$ and $S = \text{BCCDEF}$, and a query $Q = \text{ACE}$. Let the similarity threshold be $\theta = 0.5$. Suppose the similarity estimation is accurate. The near-duplicate text alignment problem returns three subsequences $T[1, 6]$, $T[4, 6]$, and $S[3, 5]$. Note that $S[2, 5] = \text{CCDE}$ is not a result as $\hat{J}_{Q,S[2,5]} = \frac{2}{5} < \theta$.

3 The Framework

To improve the query performance of near-duplicate text alignment, it is natural to index the min-hash sketches of all subsequences within each text in D . This way, when a query arrives, one shall use the index to identify all the min-hash sketches that are similar to the query's sketch and then return the corresponding subsequences as the results efficiently.

The main challenge of indexing the min-hash sketches is the cost – the number of min-hash sketches in a text grows quadratically with the length of the text (as there are $O(n^2)$ subsequences in a text with n tokens). To reduce the index size, we have a key observation: the min-hash of adjacent subsequences of a text T , such as $T[i, j]$ and $T[i, j+1]$, tend to be the same. Thus, the min-hash of adjacent subsequences can be grouped and compactly represented. For this purpose, we define the compact window.

3.1 Compact Window and Partition

Below, we formally define the compact window [15], a structure used to represent the groups of subsequences where the subsequences in each group share the same min-hash.

DEFINITION 2 (COMPACT WINDOW [15]). Given a text T and a random universal hash function h , a compact window in T is a tuple $\langle T, h, v, a, b, c, d \rangle$ satisfying that

- $1 \leq a \leq b \leq c \leq d \leq |T|$ and
- $h(T[i, j]) = v$ for every integers $i \in [a, b]$ and $j \in [c, d]$.

EXAMPLE 4. Consider the running example in Figure 1. Figure 1(a) shows the min-hash of all the subsequences of the text T . Specifically, the integer in the cell (i, j) is the min-hash of the subsequence

$T[i, j]$. For example, the min-hash of $T[1, 10]$ is shown in the cell $(1, 10)$, which is 1. As one can verify, both $\langle T, h, 1, 1, 2, 8, 10 \rangle$ and $\langle T, h, 2, 4, 5, 5, 10 \rangle$ (the highlighted rectangle) are compact windows.

Hereinafter, for two integers p, q , we abuse $[p, q]$ to denote the integer set $\{p, p+1, \dots, q\}$. Let $[a, b] \times [c, d]$ be the set of all integer pairs (i, j) where $i \in [a, b]$ and $j \in [c, d]$. A compact window $\langle T, h, v, a, b, c, d \rangle$ represents all the subsequences $T[i, j]$ where $(i, j) \in [a, b] \times [c, d]$. We aim to generate a set of compact windows such that each subsequence in T is represented by one and only one compact window in the set (i.e., the set of compact windows is a lossless compression of the min-hash of all subsequences in a text). Formally, we define the concept of *partition* as below.

DEFINITION 3 (PARTITION). Given a text T and a hash function h , a partition $\mathcal{P}(T, h)$ is a set of compact windows $\{W_1, W_2, \dots, W_l\}$, where $l = |\mathcal{P}(T, h)|$ and $W_x = \langle T, h, v_x, a_x, b_x, c_x, d_x \rangle$, that satisfies

- **Disjointness:** For any two compact windows W_x and W_y in $\mathcal{P}(T, h)$,

$$([a_x, b_x] \times [c_x, d_x]) \cap ([a_y, b_y] \times [c_y, d_y]) = \emptyset.$$

- **Coverage:** The union of all compact windows in $\mathcal{P}(T, h)$ covers all the subsequences in T , i.e.,

$$\{(i, j) \mid 1 \leq i \leq j \leq |T|\} \subseteq \bigcup_{x=1}^l ([a_x, b_x] \times [c_x, d_x]).$$

EXAMPLE 5. Consider the running example in Figure 1. Figures 1(b) and 1(c) show two example partitions with 13 and 17 compact windows respectively, as outlined by the colors. Each colored rectangle $[a, b] \times [c, d]$ in the figures corresponds to a compact window $\langle T, h, v, a, b, c, d \rangle$. All the subsequences in this rectangle share the same min-hash v .

3.2 Indexing and Query Processing

Algorithm 1 shows the pseudo-code of indexing in our framework. It takes a collection D of data texts and an integer k as input and produces k inverted indexes of compact windows. For this purpose, it first randomly selects k independent hash functions h_1, h_2, \dots, h_k . Then, for each text T in D , and each hash function h_i , it generates a partition $\mathcal{P}(T, h_i)$. The partition generation algorithm will be described in the next section. For each compact window $\langle T, h_i, v, a, b, c, d \rangle$ in the partition, it is appended to the inverted list $I_i[v]$. Finally, k inverted indexes I_1, \dots, I_k are returned.

Algorithm 2 presents the pseudo-code of query processing in our framework. The input consists of a query text Q , a multi-set Jaccard similarity threshold θ , the integer k , the same k random hash functions h_1, h_2, \dots, h_k used during indexing, and the k inverted indexes I_1, \dots, I_k . It first calculates the k min-hash v_1, \dots, v_k of Q and then retrieves the k corresponding inverted lists $I_1[v_1], \dots, I_k[v_k]$ from the k inverted indexes. The near-duplicate subsequences can be identified by a simple plane sweep algorithm over the compact windows in the k inverted lists. Specifically, the plane sweep algorithm produces all the cells (x, y) where there are at least $\lceil k\theta \rceil$ compact windows $\langle T, h, v, a, b, c, d \rangle$ in the k inverted lists satisfying $(x, y) \in [a, b] \times [c, d]$. This is because the estimated multi-set Jaccard similarity $\hat{J}_{Q,T[x,y]} \geq \theta$. The simple plane sweep algorithm is detailed in [33], while an optimized version is presented in [34]. Due to space constraints, we omit the details in this paper.

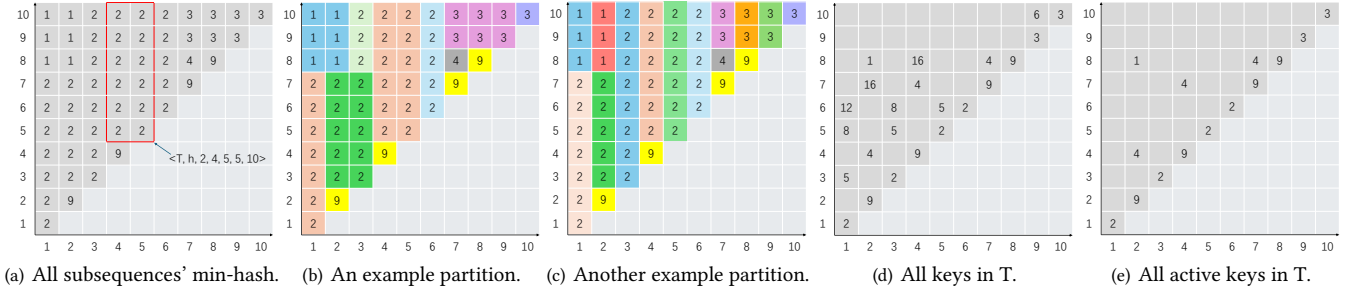


Figure 1: A running example used throughout the paper with a text $T = ABABAABBCC$ and a hash function h defined as follows:
 $h(A, 1) = 2, h(A, 2) = 5, h(A, 3) = 8, h(A, 4) = 12, h(B, 1) = 9, h(B, 2) = 4, h(B, 3) = 16, h(B, 4) = 1, h(C, 1) = 3$, and $h(C, 2) = 6$.

Algorithm 1: Indexing(D, k)

Input: D : a collection of data texts; k : an integer.

Output: I_1, I_2, \dots, I_k : k inverted indexes.

```

1 randomly select  $k$  hash functions  $h_1, h_2, \dots, h_k$  from  $\mathcal{H}$ ;
2 foreach  $T \in D$  do
3   foreach  $i \in [1, k]$  do
4      $\mathcal{P} \leftarrow \text{PARTITIONGENERATION}(T, h_i)$ ;
5     foreach  $\langle T, h_i, v, a, b, c, d \rangle \in \mathcal{P}$  do
6        $\text{append } \langle T, h_i, v, a, b, c, d \rangle$  to  $I_i[v]$ ;
7 return  $I_1, I_2, \dots, I_k$ ;

```

Algorithm 2: QueryProcessing($Q, \theta, k, h_1, \dots, h_k, I_1, \dots, I_k$)

Input: Q : a query text; θ : a similarity threshold; k : an integer;
 h_1, \dots, h_k : hash functions; I_1, \dots, I_k : inverted indexes.

Output: All subsequences $T[x, y]$ in D where $\hat{J}_{Q, T[x, y]} \geq \theta$.

```

1 foreach  $i \in [1, k]$  do  $v_i \leftarrow h_i(Q)$ ;
2 return  $\text{PLANE SWEEP}(k, \theta, I_1[v_1], \dots, I_k[v_k])$ ;

```

4 Partition Generation

Given a text T and a hash function h , there exist many possible partitions. In our framework, both the indexing and query costs scale with the partition size. This raises a natural question: how can we generate a small partition, and what is the smallest possible partition? In this section, we study the partition generation problem.

DEFINITION 4 (PARTITION GENERATION). *Given a text T and a random universal hash function h , the partition generation problem is to generate a partition $\mathcal{P}(T, h)$.*

4.1 Monotonic Partitioning

In this section, we present our partition generation algorithm. To this end, we first define the hash value set of a subsequence, which contains all the hash values of the subsequence.

DEFINITION 5 (HASH VALUE SET). *Given a random universal hash function h , the hash value set of a subsequence $T[i, j]$ is*

$$H(T[i, j], h) = \{h(t, x) \mid t \in T[i, j], 1 \leq x \leq f(t, T[i, j])\}.$$

Based on Equation 1, the *min-hash* of a subsequence is the smallest hash value in the hash value set of the subsequence, i.e.,

$$h(T[i, j]) = \min H(T[i, j], h). \quad (3)$$

We omit the hash function h in the hash value set for brevity.

EXAMPLE 6. *Consider the T and h from the running example. Figures 2(a) and 2(b) show the hash value sets $H(T[3, 6])$ and $H(T)$. In the figures, each integer in the highlighted cell (t, x) represents the*

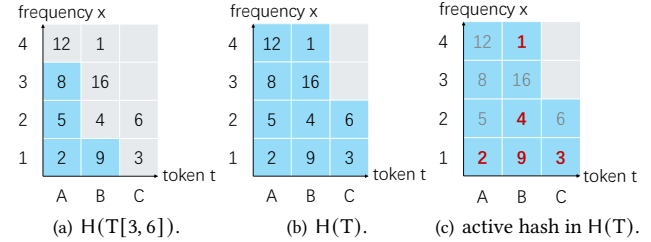


Figure 2: Examples of hash value sets and active hash values.

hash value $h(t, x)$ in the corresponding hash value set. Specifically, we have $H(T[3, 6]) = \{2, 5, 8, 9\}$. Thus $h(T[3, 6]) = \min H(T[3, 6]) = 2$.

The min-hash of any subsequence of a text T must belong to the hash value set $H(T)$ of T , as formalized below.

LEMMA 1. $h(T[i, j]) \in H(T)$ for any $1 \leq i \leq j \leq |T|$.

PROOF. Based on Definition 5, $H(T[i, j]) \subseteq H(T)$. Based on Equation 3, $h(T[i, j]) = \min H(T[i, j])$. Thus $h(T[i, j]) \in H(T)$. \square

Based on Lemma 1, the min-hash of any subsequence of T belongs to $H(T)$. Thus, for each hash value $v = h(t, x) \in H(T)$, we propose to generate a few disjoint compact windows to represent all the subsequences whose min-hash is v . The compact windows generated along the way must form a partition.

Specifically, we observe that if the min-hash of a subsequence $T[i, j]$ comes from a token t and a frequency x such that $h(t, x) = \min H(T[i, j])$, then $x \leq f(t, T[i, j])$, and thus there must exist a subsequence $T[p, q]$ of $T[i, j]$ (i.e., $i \leq p \leq q \leq j$) such that $T[p, q]$ starts and ends with token t and has exactly x occurrences of t ; otherwise $h(t, x) \notin H(T[i, j])$ and cannot be the min-hash of $T[i, j]$ based on Lemma 1. Next, we formalize our observation.

DEFINITION 6 (KEY). *Given a text T , a key s is a pair of integers (p, q) with $1 \leq p \leq q \leq |T|$ such that $T[p] = T[q]$. We refer to (p, q) as the coordinates of s , define $s.x = p$ and $s.y = q$. The hash value of s is $h(T[q], f(T[q], T[p, q]))$, abbreviated as $h(T, p, q)$.*

DEFINITION 7 (KEY SET). *The key set of a subsequence $T[i, j]$, denoted as $K(T[i, j])$ is $\{(p, q) \mid [p, q] \subseteq [i, j] \text{ and } (p, q) \text{ is a key}\}$.*

EXAMPLE 7. *Consider the running example in Figure 1. Figure 1(d) shows all the keys in $K(T)$ and their hash values, where the integer in cell (p, q) represents the hash value of key $(p, q) \in K(T)$ (in contrast, Figure 1(a) shows the min-hash of the subsequence $T[p, q]$ in cell (p, q)). In total, there are 23 keys in $K(T)$. The x -value and y -value of the key $(1, 3)$ are respectively 1 and 3. The hash value of $(1, 3)$ is $h(T, 1, 3) = h(A, f(A, T[1, 3])) = h(A, 2) = 5$.*

The min-hash of a subsequence is exactly the smallest hash value of all its keys, as formalized below.

LEMMA 2. $h(T[i, j]) = \min\{h(p, q, T) \mid (p, q) \in K(T[i, j])\}$.

PROOF. Based on Definitions 5, 6, and 7, we have $H(T[i, j]) = \{h(p, q, T) \mid (p, q) \in K(T[i, j])\}$. Based on Equation 1, we have $h(T[i, j]) = \min\{h(p, q, T) \mid (p, q) \in K(T[i, j])\}$. \square

EXAMPLE 8. In Figure 1(d), the key set $K(T[i, j])$ consists of all keys on the bottom-right of the cell (i, j) , i.e., (p, q) with $p \geq i$ and $q \leq j$. For example, $K(T[1, 3]) = \{(1, 1), (1, 3), (2, 2), (3, 3)\}$. By Lemma 2, the min-hash of $T[1, 3]$ is $h(T[1, 3]) = \min\{5, 2, 9, 2\} = 2$.

We say a subsequence $T[i, j]$ contains a key (p, q) if and only if $(p, q) \in K(T[i, j])$. Lemma 3 follows directly from Lemma 2.

LEMMA 3. A subsequence has min-hash v if and only if it contains a key with hash value v and no key with a smaller hash value.

Lemma 3 motivates us to visit all keys in $K(T)$ in ascending order of their hash values (breaking ties arbitrarily). When visiting a key (p, q) with hash value v , we identify all subsequences that contain (p, q) but no visited key. The min-hash of all these subsequences must be v based on Lemma 3. We then group them into disjoint compact windows. Next, we characterize these subsequences.

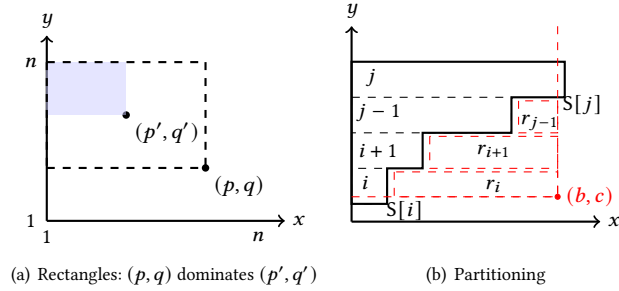


Figure 3: Monotonic partitioning.

DEFINITION 8 (RECTANGLE OF A KEY). Given a key (p, q) in a text T , where $n = |T|$, its rectangle $\text{rec}(p, q) = [1, p] \times [q, n] = \{(i, j) \mid i \in [1, p], j \in [q, n]\}$. For a set L of keys, $\text{rec}(L)$ is the union of $\text{rec}(p, q)$ with $(p, q) \in L$.

Figure 3(a) shows the rectangles of two keys (p, q) (outlined by dashed lines) and (p', q') (the shaded rectangle). By Definitions 7 and 8, a subsequence $T[i, j]$ contains a key (p, q) if and only if $(i, j) \in \text{rec}(p, q)$, i.e., the rectangle $\text{rec}(p, q)$ precisely characterizes all subsequences that contain the key (p, q) . To exclude subsequences that contain visited keys, we maintain a “skyline” to track the union of the rectangles of visited keys, as formalized below.

DEFINITION 9 (DOMINANCE). Consider two keys (p, q) and (p', q') . (p, q) dominates (p', q') if and only if $[p, q] \subset [p', q']$.

Informally, a key dominates all keys located on its top-left side, as illustrated in Figure 3(a), regardless of their hash values. Note that a key does not dominate itself based on the definition. In Example 7, the key $(1, 1)$ dominates another key $(1, 3)$ as $[1, 1] \subset [1, 3]$.

DEFINITION 10 (SKYLINE). The skyline of a set L of keys, denoted as $S(L)$, is the subset of keys in L that are not dominated by any key in L : $S(L) = \{(p', q') \in L \mid \nexists (p, q) \in L, (p, q) \text{ dominates } (p', q')\}$.

LEMMA 4. For a key set L and its skyline $S(L)$, $\text{rec}(L) = \text{rec}(S(L))$.

PROOF. By definitions, a key (p, q) dominates another key (p', q') if and only if $\text{rec}(p', q') \subset \text{rec}(p, q)$. Furthermore, by the definition of skyline, for any key s of L not in the skyline $S(L)$, there must be another key s' in the skyline that dominates s . Thus $\text{rec}(s) \subset \text{rec}(s')$. Therefore $\text{rec}(L) = \text{rec}(S(L))$. \square

When visiting a key $(p, q) \in K(T)$ with hash value v , let L be the set of all visited keys. Based on the discussion above, all subsequences in $\text{rec}(p, q) \setminus \text{rec}(S(L))$ must have min-hash v . As shown in Figure 3 (the red dot (b, c) is (p, q) and the indented black line is the skyline), these subsequences collectively form a “staircase shape”. We generate one compact window for “each step of the staircase” (red dashed rectangles).

In summary, our partitioning algorithm works as follows. Given a text T , we traverse all keys in $K(T)$ in ascending order of their hash values, maintaining a dynamic skyline of the visited keys. When visiting a key (p, q) with hash value v , let L be the set of visited keys. For each key in $S(L)$ that is dominated by (p, q) , we generate a compact window. We then update the skyline by adding (p, q) if it is not dominated by any key in the skyline, and removing all keys in the skyline that are dominated by (p, q) . As a result, the skyline is updated to $S(L \cup \{(p, q)\})$. The following sections elaborate on these details.

EXAMPLE 9. Consider the running example. There are 23 keys in T in total. The first 8 of them in the order of hash values are $(2, 8)$, $(1, 1)$, $(3, 3)$, $(5, 5)$, $(6, 6)$, $(9, 9)$, $(10, 10)$, and $(2, 4)$. As illustrated in Figure 4, we first visit $(2, 8)$. The skyline is empty. A compact window $(T, h, 1, 1, 2, 8, 10)$ is generated and $(2, 8)$ is added to the skyline. Then we visit $(1, 1)$. A compact window is generated and $(1, 1)$ is added to the skyline. Next, we visit $(3, 3)$. Notice that two compact windows are generated. Since the key $(2, 8)$ in the skyline is dominated by $(3, 3)$, it is removed from the skyline, while $(3, 3)$ is added. Later we visit $(5, 5)$ and one compact window is generated. The process is repeated. When we visit the key $(2, 4)$, we find $(2, 4)$ is dominated by the key $(3, 3)$ on the skyline. Thus no compact window will be generated and the skyline remains the same. In the end, 13 compact windows are generated as shown in Figure 1(b).

4.2 The Monotonic Partitioning Algorithm

Before presenting the pseudo-code of our algorithm, we first discuss how to efficiently find all keys in a skyline that are dominated by a key. This can be achieved by binary search, as formalized next.

DEFINITION 11 (COORDINATE ORDER). A skyline S is in coordinate order if its keys are ordered lexicographically by their coordinates. We denote $S[i]$ as the i -th key in S under this order.

The following lemmas show a skyline is a totally ordered set.

LEMMA 5. Let S be a skyline of any set of keys. For any two keys (p, q) and (p', q') in S , if $p \leq p'$, then $q \leq q'$.

PROOF. Suppose that $p \leq p'$ and $q > q'$. Then (p, q) dominates (p', q') , since $p \leq p'$ and $q > q'$ implies $[p', q'] \subset [p, q]$. This contradicts the definition of the skyline, where no key is dominated by another. Therefore, we must have $q \leq q'$. \square

LEMMA 6. Let S be the skyline of a set of keys in coordinate order. Then $S[1].x \leq \dots \leq S[l].x$ and $S[1].y \leq \dots \leq S[l].y$ where $l = |S|$.

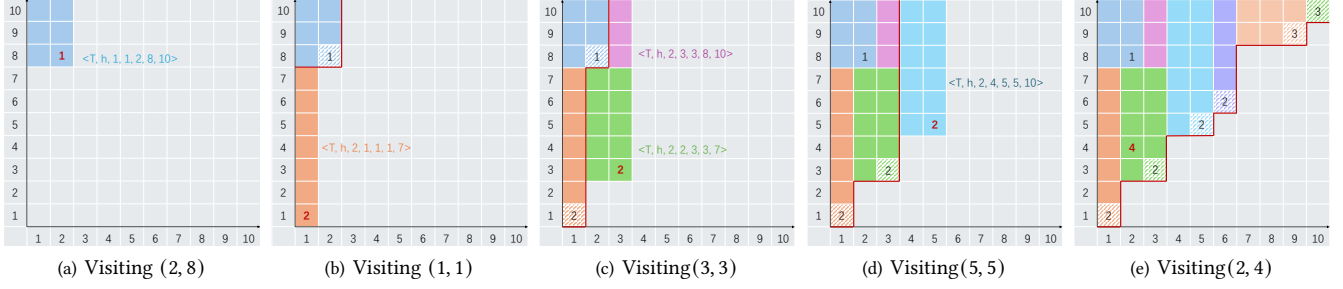


Figure 4: An example of monotonic partitioning.

PROOF. Consider any two consecutive keys $S[i]$ and $S[i+1]$ for $1 \leq i < l$. Since the keys in S are sorted lexicographically by their coordinates, we have $S[i].x \leq S[i+1].x$. By Lemma 5, this implies $S[i].y \leq S[i+1].y$. This completes the proof. \square

Lemma 7 shows that by a binary search, one can determine if a key is dominated by any key in a skyline under coordinate order.

LEMMA 7. *Given a skyline S in coordinate order and a key (b, c) , let j be the largest index such that $S[j].y \leq c$. There exists a key in S that dominates (b, c) if and only if j exists and $S[j]$ dominates (b, c) .*

PROOF. Suppose there exists a key $S[i]$ in S that dominates (b, c) , i.e., $b \leq S[i].x$ and $S[i].y \leq c$. Then j exists by the definition of j . Since j is the largest index with $S[j].y \leq c$, we have $S[i].y \leq S[j].y$, which by Lemma 6, implies $S[i].x \leq S[j].x$. Thus, $b \leq S[i].x \leq S[j].x$ and $S[j].y \leq c$. As $S[i]$ dominates (b, c) , $S[j]$ dominates (b, c) . On the other hand, if j exists and $S[j]$ dominates (b, c) , then clearly S contains a key (namely $S[j]$) that dominates (b, c) . \square

Algorithm 4 shows the pseudo-code our monotonic partitioning algorithm. It takes a text T and a hash function h as input and produces a partition $\mathcal{P}(T, h)$. It first generates all the keys in the text using the procedure GENERATEKEYS (Line 1). GENERATEKEYS (Algorithm 3) first builds an inverted index to keep track of the positions of each distinct token in T (Lines 1-3). For each inverted list (t, A) (Line 4), it enumerates every pair of positions $A[i]$ and $A[j]$ where $i \leq j$ (Lines 5-6). Then $(A[i], A[j])$ must be a key. The key, along with its hash value $h(t, j - i + 1)$, is added to the key set K (Line 7). Finally, the procedure returns the set K of all keys in T , along with their hash values (Line 8).

Once the key set K is generated, Algorithm 4 initializes the skyline with two guard keys $(0, 0)$ and $(n+1, n+1)$ in Line 2. The guard keys help simplify the algorithm by avoiding a few corner cases. The keys in S will always be kept in coordinate order using a self-balanced binary search tree. After that, we visit the keys in K in ascending order of their hash values; let (b, c) be the key currently being visited and v be its hash value (Line 3), denote by L the set of keys that have already been visited. S maintains the skyline on L upon the insertion of (b, c) . Specifically, if (b, c) has been dominated by any key in S , then $\text{rec}(b, c) \setminus \text{rec}(L)$ is empty and no subsequence in $\text{rec}(b, c)$ would have min-hash v . It means all subsequences in $\text{rec}(b, c)$ have been represented by compact windows generated earlier. In this case, we generate no compact window: skip (b, c) and jump to the next key (Line 5).

Based on Lemma 7, we can efficiently determine whether (b, c) is dominated by any key in S via binary search. Note that if the result is $j' = 0$ (Line 4), then $S[j']$ refers to the guard key $(0, 0)$, which does not dominate (b, c) . This case corresponds to the scenario in Lemma 7 where no such index j exists.

Algorithm 3: GenerateKeys(T, h)

Input: T : a text of length n ; h : a universal hash function.
Output: K : the set of keys in T .
1 $\text{map} \leftarrow \emptyset$, a map that maps a token to an array of integers;
2 **foreach** $1 \leq i \leq n$ **do**
3 append i to the array of token $T[i]$, i.e., $\text{map}[T[i]]$;
4 **foreach** token t in map (denote by A the array of $\text{map}[t]$) **do**
5 **foreach** $i \in [1, |A|]$ **do**
6 **foreach** $j \in [i, |A|]$ **do**
7 add $(A[i], A[j], h(t, j - i + 1))$ to K ;
8 **return** K ;

Algorithm 4: MonotonicPartitioning(T, h)

Input: T : a text of length n ; h : a universal hash function.
Output: A partition $\mathcal{P}(T, h)$.
1 $K \leftarrow \text{GENERATEKEYS}(h, T)$;
2 $S \leftarrow \{(0, 0), (n+1, n+1)\}$, a skyline kept in coordinate order;
3 **foreach** $(b, c, v) \in K$ in ascending order of the hash value v **do**
4 $j' \leftarrow$ binary search the largest index such that $S[j'].y \leq c$;
5 **if** $S[j']$ dominates (b, c) **then continue**;
6 $i \leftarrow$ binary search the largest index such that $S[i].y < c$;
7 $j \leftarrow$ binary search the smallest index such that $S[j].x > b$;
8 $c' \leftarrow c$;
9 **foreach** $i \leq k < j$ **do**
10 $a \leftarrow S[k].x$; $d \leftarrow S[k+1].y - 1$;
11 // this is to avoid the case when $S[i+1].y = c$
12 **if** $a \leq b$ and $c' \leq d$ **then**
13 add $\langle T, h, v, a, b, c', d \rangle$ to $\mathcal{P}(T, h)$;
14 $c' \leftarrow S[k+1].y$;
15 remove $S[i+1], S[i+2], \dots, S[j-1]$ from S ;
16 add (b, c) to S ;
17 **return** $\mathcal{P}(T, h)$;

Algorithm 5: GenerateActiveKeys(T, h)

// Replace Lines 4-7, Algorithm 3 with:
1 **foreach** entry $(t, A) \in \text{map}$ **do**
2 $\text{minkey} \leftarrow +\infty$;
3 **foreach** $1 \leq x \leq |A|$ **do**
4 **if** $\text{minkey} > h(t, x)$ **then**
5 $\text{minkey} \leftarrow h(t, x)$;
6 **foreach** $1 \leq i \leq |A| - x + 1$ **do**
7 add $(A[i], A[i+x-1], h(t, x))$ to K ;

When (b, c) is not dominated by any key in S , we first find all the keys in S that are dominated by (b, c) . Lemma 14 (C1) shall prove that with two binary searches on S (Lines 6-7), $S[i+1]$ to $S[j-1]$ are exactly the keys in S dominated by (b, c) . Lines 8-13 generate up to $j-i$ compact windows (red dashed rectangles in Figure 3(b)). Lemma 14 (C2) shall prove that these compact windows are disjoint and cover all the subsequences $T[i, j]$ where $(i, j) \in \text{rec}(b, c) \setminus \text{rec}(L)$, whose min-hash must be v (by Lemma 3 and Lemma 4). We update S by removing these dominated keys (Line 14) and inserting the new key (b, c) (Line 15).

In practice, Line 4 and Line 6 can be combined in one binary search. We use two binary searches for ease of presentation later.

EXAMPLE 10. Consider the running example. As shown in Figure 4(c), when visiting $(b = 3, c = 3)$, we have $v = 2$ and $S = \{(0, 0), (1, 1), (2, 8), (11, 11)\}$. The binary search would find $j' = 2$ as $S[2].y = 1 \leq c$ and $S[3].y = 8 > c$. Clearly, $S[j'] = (1, 1)$ does not dominate $(3, 3)$. Then the next two binary searches would find $i = 2$ and $j = 4$. Next, when $k = 2$, a compact window $\langle T, h, v = 2, a = 2, b = 3, c' = 3, d = 7 \rangle$ is generated and c' becomes 8. When $k = 3$, another compact window $\langle T, h, 2, 3, 3, 8, 10 \rangle$ is produced. Finally, the key $S[3] = (2, 8)$ is removed from the skyline, while $(3, 3)$ is added.

Let f_T be the maximum token frequency in T , i.e., $f_T = \max_{t \in T} f(t, T)$.

THEOREM 1. Given a text T and a hash function h , Algorithm 4 generates a partition $\mathcal{P}(T, h)$.

THEOREM 2. In expectation, the partition $\mathcal{P}(T, h)$ generated by Algorithm 4 has $O(n + n \log f_T)$ compact windows, where $n = |T|$.

THEOREM 3. When given the length and the maximum token frequency of a text, Algorithm 4, the Monotonic Partitioning algorithm, is optimal in terms of the expected partition size in the worst case.

We prove Theorem 1 in Appendix A (due to space limitation) and Theorem 2 in Section 4.3. Note that the “expectation” in Theorem 2 is due to the nature of the hash function as opposed to any assumption on T . We also provide a lower bound in Theorem 3, with the proof provided in Appendix B for interested readers.

4.3 Active Keys and Complexity Analysis

The total number of keys generated by Algorithm 3 is $O(\sum_{t \in T} f^2(t, T))$, which can go up to $O(n^2)$ – when all the tokens are duplicated (i.e., $f_T = n$). This leads to a running time of $O(n^2 \log n)$ and space $O(n^2)$ just for producing, sorting and storing the keys.

We observe that not all keys in $K(T)$ need to be enumerated. Specifically, in Algorithm 4, when visiting a key $(b, c) \in K(T)$, the key is skipped if it is dominated by a key in the skyline of visited keys (Line 5). Although it is non-trivial to determine all the skipped keys in advance, we identify a subset of them – referred to as *non-active keys* – which, as formalized below, are guaranteed to be dominated and can thus be safely skipped during key generation.

DEFINITION 12 (ACTIVE HASH VALUE). Given a text T and a hash function h , a hash value $h(t, x)$ for a token $t \in T$ is said to be active if and only if $h(t, x) < h(t, x')$ for all positive integers $x' < x$.

That is to say, a hash value $h(t, x)$ is active if and only if it is the smallest among $h(t, 1), h(t, 2), \dots, h(t, x-1), h(t, x)$. For example, Figure 2(c) shows all the active hash values in the hash value set

$H(T)$ in Figure 2(b). The bold, red integers are active hash values, while the rest (i.e., gray ones) are non-active hash values.

DEFINITION 13 (ACTIVE KEY). A key (p, q) is active if and only if its hash value $h(p, q, T)$ is active for token $T[p]$ (note that $T[p] = T[q]$ for a key). The set of active keys in a text T is denoted as $X(T)$.

For example, Figure 1(e) shows the active key set $X(T)$ of T under the hash function h from the running example. Although T contains 23 keys in total, only 14 of them are active keys.

Note that a key is added to the skyline in Algorithm 4 if and only if it is not skipped, i.e., the condition in Line 5 evaluates to false.

LEMMA 8. No non-active key is added to the skyline in Algorithm 4.

PROOF. Let $(p, q) \in K(T)$ be a non-active key. Then its hash value $h(t, x)$ must be non-active, where $t = T[p] = T[q]$ and $x = f(t, T[p, q])$. Let $p = q_1 < q_2 < \dots < q_x = q$ be the positions such that $T[q_i] = t$ for all $1 \leq i \leq x$. Since $h(t, x)$ is non-active, by Definition 12, there exists some $1 \leq i < x$ such that $h(t, i) < h(t, x)$. Then the corresponding key (p, q_i) has a smaller hash value than (p, q) and satisfies $[p, q_i] \subset [p, q]$, hence it dominates (p, q) . Moreover, since $h(t, i) < h(t, x)$, (p, q_i) is visited before (p, q) .

Now, consider two cases. (1) If (p, q_i) is in the skyline at the time (p, q) is visited, then (p, q) will be skipped by Lemma 7. (2) If (p, q_i) is not in the skyline, it must be pruned by a key dominates (p, q_i) . By Lemma 9 (dominance is transitive), there exists a key in the skyline that dominates (p, q_i) and thus also dominates (p, q) . In this case, (p, q) will also be skipped. In either case, (p, q) will not be added to the skyline. Therefore, only active keys can enter the skyline in Line 15, and all non-active keys are skipped. \square

LEMMA 9. Dominance is transitive, i.e., if (p, q) dominates (p', q') and (p', q') dominates (p'', q'') , then (p, q) dominates (p'', q'') .

PROOF. This is because $[p, q] \subset [p', q'] \subset [p'', q'']$. \square

LEMMA 10. The total number of compact windows generated by Algorithm 4 is no more than $2|X(T)|$.

PROOF. Consider Lines 9-13. The algorithm produces at most $j-i$ compact windows. Note that inserting the key (b, c) results in the removal of $j-i-1$ other keys from the skyline, and these removed keys will not be reinserted. Thus, among the $j-i$ compact windows generated, one corresponds to the insertion of (b, c) and one to the removal of each key. Thus one compact window is produced upon the insertion of a key and one compact window is produced upon the removal of a key. Moreover, by Lemma 8, only active keys can be added to the skyline. Hence, the total number of compact windows generated by Algorithm 4 is at most $2|X(T)|$. \square

LEMMA 11. In expectation, $|X(T)| = O(n + n \log f_T)$, where $n = |T|$.

PROOF. For any $t \in T$ and $i \in [f(t, T)]$, the probability that $h(t, i)$ is active, i.e., it is smaller than all the hash values $h(t, j)$ for all $j < i$ is $\frac{1}{i}$. Therefore, the total number of active hash values among $h(t, 1), \dots, h(t, f(t, T))$ in expectation is $\sum_{i \in [f(t, T)]} \frac{1}{i} = O(1 + \ln(f(t, T)))$, seeing that even when $f(t, T) = 1$, the summation is still 1. Moreover, there are $O(f(t, T))$ keys with hash value $h(t, i)$ for any $i \in [f(t, T)]$. Therefore, in expectation, the total number of active keys in T , i.e., $|X(T)|$, is $O(\sum_{t \in T} (f(t, T)(1 + \ln(f(t, T)))) = O(n + \sum_{t \in T} (f(t, T) \ln(f(t, T))) = O(n + n \log f_T)$. \square

Combining Lemmas 10 and 11, in expectation, the total number of compact windows in the partition $\mathcal{P}(T, h)$ generated by Algorithm 4 is $O(|X(T)|) = O(n + n \log f_T)$, which proves Theorem 2.

To incorporate this optimization in our algorithm, we replace procedure GENERATEKEYS with a similar procedure GENERATEACTIVEKEYS (Algorithm 5). It generates active keys by enumerating, for each token t , every possible frequency x of t : the algorithm maintains the smallest hash value encountered in *minkey* and only when the current hash value $h(t, x) < \text{minkey}$, it produces all the keys whose hash values are $h(t, x)$ and updates *minkey* as $h(t, x)$.

An additional optimization is to sort the active hash values prior to generating the corresponding active keys, rather than generating and subsequently sorting all active keys. This approach is justified by two key observations: (1) active keys associated with the same hash value can be ordered arbitrarily, and (2) the number of active hash values does not exceed the total number of active keys.

THEOREM 4. *Algorithm 4 with GENERATEACTIVEKEYS in Line 1 has time complexity $O(|X(T)| \log n)$ and space complexity $O(|X(T)|)$. In expectation, the time complexity is $O(n \log n + n \log n \log f_T)$, and space complexity is $O(n \log f_T + n)$.*

PROOF. Because for any $i \in [n]$, (i, i) is an active key, thus $|X(T)| \geq n$. The size of the skyline S is $O(n)$ as each $b \in [n]$ will have at most one key in S . The binary search and each update to the skyline take $O(\log n)$. GENERATEACTIVEKEYS only produces active keys $X(T)$. Each active key costs up to three binary searches of the skyline and two updates and corresponds to up to two compact windows. In addition, GENERATEACTIVEKEYS takes $O(|X(T)| + n) = O(|X(T)|)$ time; sorting the active hash values takes $O(n \log n)$ time as there are at most n active hash values. Thus the time complexity is $O(|X(T)| \log n)$. The space complexity is $O(|X(T)|)$. By Lemma 10, in expectation, $|X(T)| = O(n + n \log f_T)$. Therefore, the time complexity is $O(n \log n + n \log n \log f_T)$; the space complexity is $O(n \log f_T + n)$ in expectation. \square

Lemma 15 in Appendix A shall prove the correctness of this optimization. Note that there are $O(n f_T)$ keys in $K(T)$. Thus, the time and space complexities of vanilla Algorithm 4 (without the active key optimization) are respectively $O(n f_T \log n)$ and $O(n f_T)$.

5 Generalization to Weighted Jaccard Similarity

Consider a text T where each distinct token t is associated with a weight. The weight is determined by a weight function $w(t, T)$ where $w(t, T) > 0$ for $t \in T$ and $w(t, T) = 0$ for $t \notin T$. The weighted Jaccard similarity of two texts T and S is defined as

$$J_{T,S}^w = \frac{\sum_{t \in T \cup S} \min(w(t, T), w(t, S))}{\sum_{t \in T \cup S} \max(w(t, T), w(t, S))}.$$

To estimate weighted jaccard similarity, one can use **improved consistent weighted sampling** [25]. In a nutshell, it designs a hash family \mathcal{H} . Given a text T , for each distinct token t in T , a hash function $h \in \mathcal{H}$ takes the token t and its weight $w(t, T)$ as input and outputs a hash value, denoted as $h(t, w(t, T))$. With hash function h , we define the *weighted min-hash* of a text T as the smallest hash value among all distinct tokens in T , denoted as $h(T, w)$. Formally,

$$h(T, w) = \min\{h(t, w(t, T)) \mid t \in T\}. \quad (4)$$

Algorithm 6: Improved Consistent Weighted Sampling [25]

```

1 Class ICWS:
2   def init(self,  $\Sigma$ : the vocabulary of tokens):
3     foreach token  $t \in \Sigma$  do
4       self. $r_t \sim \text{Gamma}(2, 1)$ 
5       self. $c_t \sim \text{Gamma}(2, 1)$ 
6       self. $\beta_t \sim \text{Uniform}(0, 1)$ 
7   def  $h(t$ : a token,  $weight$ : a positive real value):
8      $y = \exp(r_t (\lfloor \frac{\ln weight}{r_t} + \beta_t \rfloor - \beta_t))$ ;
9      $a = \frac{1}{y} \cdot \frac{c_t}{\exp(r_t)}$ ;
10    return HashValue( $t, y, a$ );

```

```

11 Class HashValue:
12   def init(self,  $t, y, a$ ): self. $t = t$ , self. $y = y$ , self. $a = a$ ;
13   def operator= $(v_1, v_2)$ : return  $v_1.t = v_2.t$  and  $v_1.y = v_2.y$ ;
14   def operator< $(v_1, v_2)$ : return  $v_1.a < v_2.a$ ;

```

The hash family \mathcal{H} designed in improved consistent weighted sampling scheme guarantees that for any two texts T and S and weight function w , $\Pr[h(T, w) = h(S, w)] = J_{T,S}^w$.

Thus the weighted Jaccard similarity of two texts T and S can be accurately and unbiasedly estimated by k independent hash functions h_1, \dots, h_k randomly drawn from the hash family \mathcal{H} as

$$\hat{J}_{T,S}^w = \frac{1}{k} \sum_{i=1}^k \mathbf{1}\{h_i(T, w) = h_i(S, w)\}. \quad (5)$$

Implementation Details. Algorithm 6 implements the improved consistent weighted sampling scheme [25]. Drawing a random hash function $h \in \mathcal{H}$ is the same as instantiating a new object o of ICWS class (Lines 1-10). The constructor (Lines 2-6) samples three independent random variables for every token t in the vocabulary Σ from the ICWS specified distributions. The hash computation (Lines 7-10) $o.h(t, w(t, T))$ of o returns the hash value based on token t and the weight of t in T . Note that an object is initialized with a random seed; the set of possible objects form the hash family.

A hash value (t, y, a) is an instance of the HashValue, where the comparators $=$ and $<$ are defined. For simplicity, we assume there are no hash collisions: for any two hash values (t_1, y_1, a_1) and (t_2, y_2, a_2) , if $t_1 \neq t_2$, then $a_1 \neq a_2$. Comparator $<$ provides hash values a total order, i.e., for any two hash values $v_1 = (t_1, y_1, a_1)$ and $v_2 = (t_2, y_2, a_2)$, we say $v_1 < v_2$ if and only if $a_1 < a_2$. Besides, for a token t , r_t , c_t and β_t depend solely on t . Thus the weight determines both y and a , i.e., there is a one-to-one mapping between y and a in a hash value (t, y, a) . Thus, $v_1 = v_2$ iff. $t_1 = t_2$ and $y_1 = y_2$.

Consistency and Uniformity. The weighted min-hash under improved consistent weight sampling has the nice properties below:

- **Uniformity:** Denote by (t, y, a) the weighted min-hash of a text T , then (t, y) is distributed uniformly over $\cup_{t' \in \Sigma} \{t'\} \times [0, w(t', T)]$.
- **Consistency:** Let $v = (t, y, a)$ be the weighted min-hash of a text T . Given a text S with $w(t', S) \leq w(t', T)$ for all $t' \in \Sigma$. If $y \leq w(t, S)$, then v must also be the weighted min-hash of S .

Assumption on the Weight Function (AoW): *Given a text T , we assume that for any token t in T , its weight is a monotonically increasing function of its frequency $f(t, T)$ and is independent of*

tf weight functions	idf weight functions
binary: $\mathbf{1}\{f(t, T) > 0\}$	unary: 1
raw count: $f(t, T)$	standard: $\log \frac{N}{N_t}$
logarithmic: $\log(f(t, T) + 1)$	smooth: $\log(\frac{N + N_t}{N_t}) + 1$
squared: $f(t, T)^2$	probabilistic: $\log \frac{N - N_t}{N_t}$

Table 1: A snippet of TF and IDF weights we support. Here $N = |D|$ is the number of texts in the corpus D and $N_t = |\{S \in D \mid t \in S\}|$ is the number of texts in the corpus containing t . N is a global constant while N_t is a constant local to the token t .

other properties of T . With AoW, we can simplify the notation of the weight of a token t in T as $w(t, x)$ where $x = f(t, T)$. We can also simplify the hash function as $h_w(t, x) \doteq h(t, w(t, x))$ where x is a positive integer. The AoW, i.e., for any $1 \leq x \leq x'$, $w(t, x) \leq w(t, x')$, is reasonable especially when it comes to TF-IDF weights.

Term Frequency-Inverse Document Frequency (TF-IDF) [48]. TF-IDF is a widely used weighting scheme that captures both the importance of a token within a text and its rarity across a corpus D . The TF-IDF weight function is defined as $w(t, T) = \text{tf} \cdot \text{idf}$ where tf is the TF weight and idf is the IDF weight. A few example TF and IDF weights we support (but not limited to) are listed in Table 1.

Hash Value Set. We omit the weight function w in the notation of a hash value $h_w(t, x)$ when the context is clear and abbreviate the hash value as $h(t, x)$. We can define the hash value set of a subsequence in the same way as Definition 5.

DEFINITION 14. Given a random hash function h from \mathcal{H} and a weight function w , the hash value set of a subsequence $T[i, j]$ is

$$H(T[i, j], h) = \{h(t, x) \mid t \in T[i, j], 1 \leq x \leq f(t, T[i, j])\}.$$

Next we show that, under the assumption of AoW, the hash values under consistent weighted sampling have the property below.

LEMMA 12. Given a random hash function h from \mathcal{H} and a weight function w under the assumption of AoW, we have $h(t, 1) \geq h(t, 2) \geq \dots \geq h(t, x)$ for any token t in any text T where $x = f(t, T)$.

PROOF. Under the assumption of AoW, we have $w(t, 1) \leq w(t, 2) \leq \dots \leq w(t, x)$. Now consider texts T_1, T_2, \dots, T_x where T_i consists of exactly i copies of token t . Then T_i has only one hash value $h(t, i)$, which must be its weighted min-hash by Equation 4. Consider any $j > i$. Let $h(t, i) = (t, y, a)$ and $h(t, j) = (t, y', a')$, which are the weighted min-hash of T_i and T_j . By uniformity, y distributes uniformly over $[0, w(t, i)]$ and so does y' over $[0, w(t, j)]$. By consistency, because $w(t, i) \leq w(t, j)$, if $y' \leq w(t, i)$, we have $h(t, i) = h(t, j)$. Otherwise, i.e., $y' > w(t, i)$, because $y \leq w(t, i)$, we have $y < y'$. Under the same t , i.e., r_t, c_t, β_t are fixed, we have $a \propto \frac{1}{y}$ (Line 9, Algorithm 6). Thus we have $a > a'$, which indicates $h(t, i) > h(t, j)$. In both cases, we have $h(t, i) \geq h(t, j)$. Thus $h(t, 1) \geq h(t, 2) \geq \dots \geq h(t, x)$ for any t and T . \square

Based on Lemma 12 and Equation 4, we have

$$h(T[i, j]) = \min H(T[i, j], h). \quad (6)$$

With the above definitions, all the definitions, lemmas, theorems, algorithms, and conclusions naturally extend to the weighted min-hash. For example, a hash value $h(t, x)$ is active if and only if it is strictly greater than all of $h(t, 1), \dots, h(t, x - 1)$.

LEMMA 13. In expectation, $|X(T)|$ is $O(n)$ for binary TF, $O(n + n \log f_T)$ for raw count TF, $O(n + \log \log f_T)$ for logarithmic TF, and $O(n + n \log f_T)$ for squared TF, each combined with any IDF in Table 1.

PROOF. Given a text with n tokens and a weight function w . Consider a token $t \in T$. Let $x = f(t, T)$. Clearly, $h(t, 1)$ must be an active hash value. For any $i \in [2, x]$, based on the proof of Lemma 12, $h(t, i)$ is active if and only if $h(t, i) < h(t, i - 1)$. Let $h(t, i) = (t, y_i, a_i)$ and $h(t, i - 1) = (t, y_{i-1}, a_{i-1})$. $h(t, i) < h(t, i - 1)$ if and only if $y_i \in (w(t, i - 1), w(t, i)]$. By uniformity, y_i distributed uniformly over $[0, w(t, i)]$. Thus the probability that $h(t, i)$ is active is $\frac{w(t, i) - w(t, i - 1)}{w(t, i)}$. There are $x - i + 1 = O(x)$ keys corresponding to the hash value $h(t, i)$ (imagine a sliding window of size i over a string of length x). Thus in expectation, the number of active keys in T is $\mathbb{E}[X(T)] = O(\sum_{t \in T} f(t, T) \sum_{i=1}^{f(t, T)} \frac{w(t, i) - w(t, i - 1)}{w(t, i)})$. Note that $w(t, 0) = 0$. Clearly, for $x \geq 1$, when $w(t, x) = x \cdot \text{IDF}$, we have $\mathbb{E}[X(T)] = O(n + n \log f_T)$. When $w(t, x) = 1 \cdot \text{IDF}$, $\mathbb{E}[X(T)] = O(n)$. For $w(t, x) = \log(x + 1) \cdot \text{IDF}$, $\mathbb{E}[X(T)] = O(n + n \log \log f_T)$ based on Taylor expansion. For $w(t, x) = x^2 \cdot \text{IDF}$, $\mathbb{E}[X(T)] = O(n + n \log f_T)$. \square

A caveat is that, instead of breaking ties arbitrarily for keys with the same hash value when visiting the key set of a text, we need to order the keys firstly by their hash values and secondly by their frequencies. That is to say for two keys (p, q) and (p', q') with the same hash value, if $f(T[p], T[p, q]) > f(T[p'], T[p', q'])$, we should visit (p, q) before (p', q') . If $f(T[p], T[p, q]) = f(T[p'], T[p', q'])$, we can break tie arbitrarily.

6 Experiment

Environment. We implemented all the methods, including the baselines, using C++. The programs were compiled with GCC 11.4.0 and optimized using the `-O3` flag. All experiments were conducted on a machine equipped with an Intel Xeon Gold 6212 CPU @ 2.40GHz and 1 TB of memory, running Ubuntu 18.04.5.

Datasets. We conducted experiments on two datasets, PAN [36] and OWT [20]. PAN (<https://pan.webis.de/data.html>) is a benchmark for external plagiarism detection [36]. Each text in PAN is a book. Tokens are generated by splitting words based on whitespace. In total, there are 11,093 texts and 642,380,109 tokens in the dataset. OWT consists of 8 million web texts highly ranked on Reddit. Each token is a byte-pair-encoding generated by the GPT-2 tokenizer [40]. The vocabulary size is 50,257.

Parameters and Settings. There are three parameters, the text length n , the maximum token frequency f , and the min-hash sketch size k (see Section 2.2). We report the compact window generation time (i.e., partition time) and the number of compact windows generated (i.e., partition size) per text. To avoid text-dependent bias, for each configuration of parameters, we randomly choose 10 text under the configuration and report the average partition time/size.

To fix the text length n , texts with more than n tokens were truncated to their first n tokens. Texts with fewer than n tokens were concatenated until they reached or exceeded n tokens, after which they were truncated to the first n tokens as a single text.

To evaluate the impact of the maximum token frequency f , we fixed the text length n and used the first 100,000 texts in OWT and

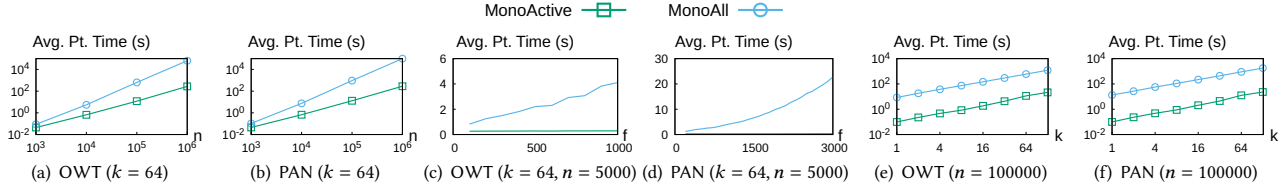


Figure 5: Evaluating the active hash optimization. k : sketch size, n : text length, f : maximum token frequency.

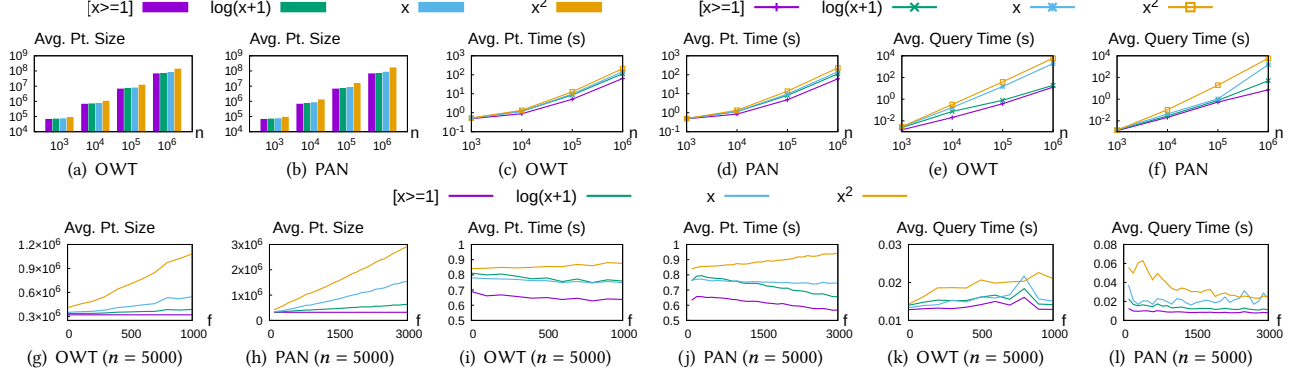


Figure 6: Evaluating weighted Jaccard similarity under various weight functions. $k = 64$ in all experiments.

the first 10,000 texts in PAN. These texts were grouped based on their maximum token frequencies using intervals of 100. For example, texts with maximum token frequencies in the range $[1, 100]$ formed one group, while those in $[101, 200]$ formed another.

Compared Methods. We consider 3 methods in our experiment.

- MonoAll: Our vanilla monotonic partitioning, Algorithm 4.
- MonoActive: Our monotonic partitioning with active hash optimization (Section 4.3). It calls Algorithm 5 to generate keys.
- AllAlign: A greedy partitioning algorithm for multi-set Jaccard similarity [15]. AllAlign generates compact windows in recursion. In each iteration, it takes a rectangle of the shape $[a, b] \times [a, c]$ as input and partitions all the subsequences $T[i, j]$ in this rectangle into a few compact windows and one or more smaller rectangles of that shape. The smaller rectangles are recursively partitioned until no rectangles left. At the beginning, the input rectangle is $[1, n] \times [1, n]$. However, AllAlign is greedy and its time and space complexities are unknown.

6.1 Evaluating Active Hash Optimization

This section evaluates our two methods MonoActive and MonoAll on two datasets OWT and PAN. Since the optimization in MonoActive does not change the generated compact windows, the partition size of MonoActive and MonoAll are the same, we only consider the partition time. Recall that in Theorem 4, the time complexities of MonoActive and MonoAll are respectively $O(n \log n \log f)$ and $O(nf \log n)$ when the maximum frequency $f > 1$.

Figures 5(a)-5(b) show that the partition time of two methods increases with the text length. MonoActive consistently outperforms MonoAll because MonoActive avoids unnecessary computations for non-active keys. The ratio between the partition time of MonoAll and that of MonoActive echoes our complexity analysis.

Figures 5(c)-5(d) show that with $k = 64$ and $n = 5000$, as the maximum frequency f increases, the runtime of MonoActive remained nearly constant whereas that of MonoAll grew almost linearly with f . For example, on PAN, when $f \in [901, 1000]$, $f \in [1901, 2000]$, and $f \in [2901, 3000]$, MonoActive took 0.27, 0.29, 0.33 seconds for partition generation on average, while MonoAll took 3.98, 10.56,

22.59 seconds. When f increased around threefold, the runtime of MonoActive and MonoAll respectively increased $1.2\times$ and $5.7\times$. This observation is consistent with our complexity analysis.

Figures 5(e)-5(f) show that both methods exhibit linear runtime growth with k . This is because the partition generation is independent across the k random universal hash functions. Thus, hereinafter, we skip evaluating the impact of the parameter k .

6.2 Evaluating Weighted Jaccard Similarity

This section evaluates near-duplicate text alignment under weighted Jaccard similarity. Our algorithm MonoActive is the only method that supports this problem. We evaluated four token weight functions $w(t, x)$: (1) binary TF weight $[x \geq 1]$; (2) logarithmic TF weight $\log(x+1)$; (3) raw count TF weight x ; and (4) squared TF weight x^2 ; all combined with unary IDF weight.

Figures 6(a) and 6(b) report the partition sizes (i.e., the number of compact windows generated) as the text length n varies from 10^3 to 10^6 . As observed, the squared weight consistently yields the largest partitions, followed by the raw count weight, whereas the binary weight produces the smallest partitions. For all weight functions, the partition size increases approximately linearly with n . For instance, when n increases from 1,000 to 1 million on PAN, the partition sizes under the four weight functions increase from 64K, 67K, 71K, and 89K to 64M, 68M, 82M, and 158M, respectively. This is because the partition sizes of the binary, logarithmic, raw count, and squared weights are respectively $O(n)$, $O(n + n \log \log f)$, $O(n + n \log f)$, and $O(n + n \log f)$, according to Lemma 13. Although squared and raw count weights have the same complexity, the constant factor of squared weight is larger than that of raw count weight.

Figures 6(c) and 6(d) present the average partitioning time (i.e., the compact window generation time) under the four weight functions as the text length n increases from 10^3 to 10^6 . Among them, the binary weight consistently achieved the fastest partitioning time, followed by the logarithmic weight, while the squared weight resulted in the slowest runtime. All weight functions exhibit quasi-linear growth with respect to n . For example, when n increased from 10^3 to 10^6 on PAN, the partitioning times under the four

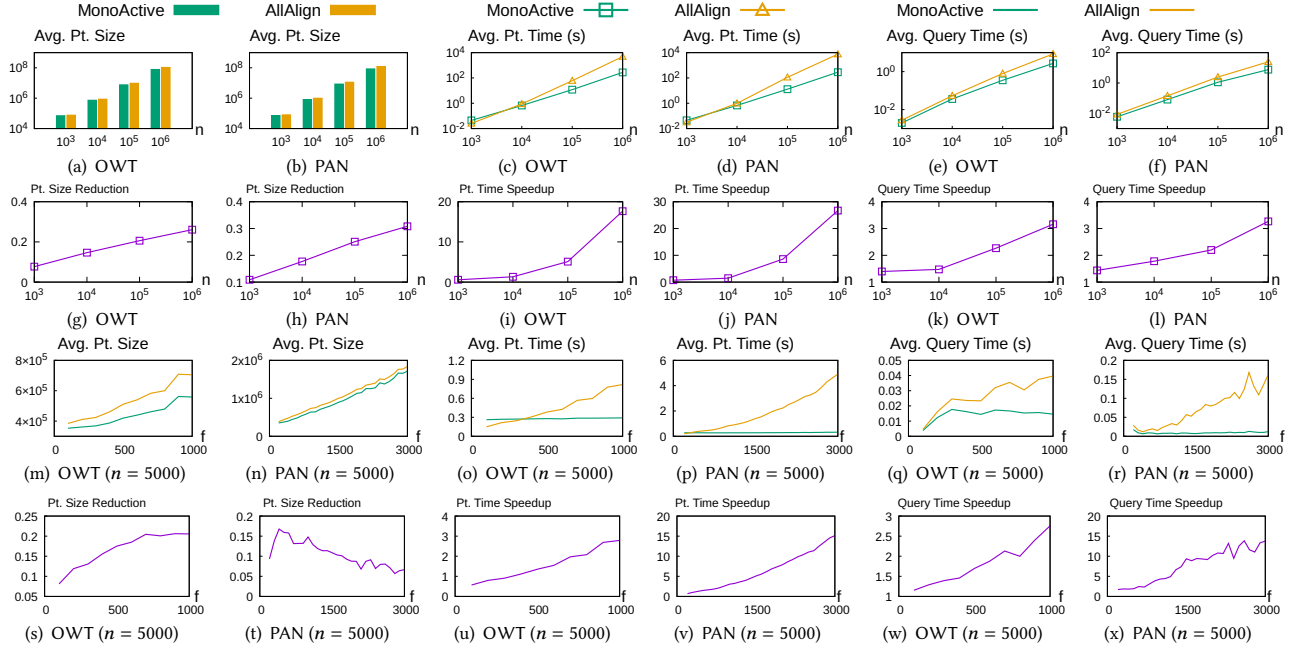


Figure 7: Comparing with the state-of-the-art for multi-set Jaccard similarity. $k = 64$ in all experiments.

weights increased from 0.47, 0.50, 0.48, and 0.49 seconds to 60.96, 110.99, 143.43, and 229.02 seconds, respectively. This aligns with our complexity analysis, which shows that the runtime of MonoActive scales with $|X(T)| \log n$ for all weight functions.

Figures 6(e) and 6(f) show the average query latency when varying the text length n from 10^3 to 10^6 . The queries are sampled from the same dataset of the data texts but are disjoint from them. Each query is a text of length 100, and the same set of queries is used across all values of n . Similarly, the squared weight had the largest query latency, followed by raw count weight, while the binary weight had the smallest latency. This is because query processing cost grows with the number of compact windows sharing the same hash values between the query and the data texts (a.k.a., collided compact windows). Weight functions that generate more compact windows are more likely to incur a larger number of such collisions, thereby increasing query latency.

Figures 6(g)–6(h) show the average partition sizes generated under the four weight functions, with $n = 5,000$ and $k = 64$ fixed while varying the maximum token frequency f . As f increased, the partition size under the binary weight remained constant. The partition size under logarithmic increased modestly, whereas those under raw count and squared weights grew sublinearly with f . For example, when f increased from the range $[901, 1000]$ to $[2901, 3000]$, the partition sizes under the binary, logarithmic, raw count, and squared weights increased from 320K, 412K, 607K, and 1M to 320K, 625K, 1.5M, and 2.8M, respectively. These results are consistent with our theoretical analysis. The numbers of compact windows generated by binary, logarithmic, raw count, and squared weights scale with n , $n \log \log f$, $n \log f$, and $n \log f$, respectively.

Figures 6(i)–6(j) report the average partitioning time of the four weight functions under varying f . The partitioning time either remained nearly constant or decreased slightly across all four weight functions. This is because a significant part of MonoActive’s runtime is spent on sorting the active hash values, the number of which is proportional to the number of distinct tokens in the text. As f increases while n remains fixed, the number of distinct tokens decreases, thereby reducing the sorting cost. Finally, Figures 6(k)–6(l)

show the average query latency of the four weight functions under varying f . The two datasets showed different behaviors. This is because the query time depends heavily on the query text (which determines the number of collided compact windows).

6.3 Comparing with the State-of-the-Art

This section compares our method MonoActive with AllAlign, the state-of-the-art method for multi-set Jaccard similarity.

Figures 7(a) and 7(b) show, when varying the text length n from 10^3 to 10^6 , the average number of compact windows generated by MonoActive and AllAlign (i.e., partition sizes, which are proportional to their index sizes) on two datasets OWT and PAN, while Figures 7(g) and 7(h) show the reduction ratio (i.e., $1 - \frac{\text{MonoActive}}{\text{AllAlign}}$). As shown, MonoActive consistently generated less compact windows than AllAlign, up to 30.82%. More importantly, the gap between the number of compact windows for MonoActive and AllAlign steadily widens as n increases. For example, when n increased from 1000 to 1 million, the gap grew from 10.86% to 30.82% on PAN and from 7.69% to 26.08% on OWT. This is because AllAlign’s recursive approach divides the rectangles into multiple rectangular regions, imposing early boundaries that may split what is a single compact window in MonoActive into multiple. The number of compact windows generated by MonoActive grew linearly with n , which is consistent with our complexity analysis.

Figures 7(c) and 7(d) show, when varying the text length n from 10^3 to 10^6 , the average compact window generation time (i.e., partition time, which is proportional to the index time) of AllAlign and MonoActive, respectively, while Figures 7(i) and 7(j) show the speedup of MonoActive over AllAlign. MonoActive was significantly faster than AllAlign. For example, on the PAN dataset, when $n = 10^6$, AllAlign took 7581s on average to generate the compact windows, whereas MonoActive only spent 284s, achieving a speedup of 26.7 \times . This is because AllAlign produces more compact windows than MonoActive and AllAlign is recursive whereas MonoActive is non-recursive. Moreover, the speedup of MonoActive over AllAlign increased as n grew. For example, on

the PAN dataset, when n increased from 10^4 to 10^6 , the speedup improved from 1.46 to 26.7. A similar trend is observed on OWT.

Figures 7(e) and 7(f) show the average query latency of the two methods when varying the text length n from 10^3 to 10^6 , while Figures 7(k) and 7(l) show the speedup of MonoActive over AllAlign. The queries are generated in the same way as described earlier. The query latency of MonoActive was up to $3.15\times$ faster than that of AllAlign, though the number of compact windows was only up to 30.82% less. For example, on the PAN dataset, when n increased from 10^3 to 10^6 , the speedup improved from 1.4 to 3.15. A similar trend is observed on OWT. This is because, although the number of collided compact windows grows linearly with the number of generated compact windows in expectation, the cost of query processing scales quadratically with the number of collided windows in the worst case. Therefore, reducing the number of compact windows can lead to a significant speedup in query time.

We also evaluate the impact of the maximum token frequency f . Figures 7(m)–7(r) present the number of compact windows generated, the compact window generation time, and the query latency of the two methods, under varying f with fixed text length $n = 5000$. Figures 7(s)–7(x) further report the relative improvements. As f increases, the generation time of MonoActive remains largely stable, whereas that of AllAlign increases significantly. For example, on the PAN dataset, when f increases from 600 to 3000, the speedup improved from 1.41 to 16.26. In addition, the gap in the number of compact windows produced by the two methods first grows and then decreases with larger f . This trend arises because higher token frequencies lead to deeper recursion and more early boundaries in AllAlign, which fragment compact windows and increase computational cost. However, as f approaches n , the text instances approach the worst-case scenario and both methods produce a large number of compact windows. The gap of the query latency increases with the increase of f . This is because the collided windows are dominated by tokens with the highest frequencies. In AllAlign, the high-frequency tokens trigger deeper recursion, resulting in more finely fragmented compact windows. These fragments increase the number of collided compact windows and the number of iterations in the inner loop. As a result, the query latency grows superlinearly with the number of collided compact windows.

Scalability. As shown in Figures 7(g)–7(l), the performance gain of MonoActive over AllAlign increased as the text length n grew in terms of partition size, partition generation time, and query latency. Thus MonoActive scales better than AllAlign. This is attributed to the complexity guarantees of our algorithm.

7 Related Work

Near-Duplicate Text Alignment. Most of existing methods for near-duplicate text alignment rely on rule-based heuristics and the “seeding-extension-filtering” framework [3, 5, 7, 22, 23, 26, 31, 43]. They first find “seed matches” between the data texts and the query. Various kinds of seeds have been proposed, such as fingerprints [35], super-shingles [8], $0 \bmod p$ [31], fixed-length windows [51], q-grams [42], and sentences [41]. Then they extend the seed matches as far as possible to form candidates. Finally, candidates failing predefined criteria (e.g., length or overlap thresholds) are filtered. However, these heuristics are highly sensitive to the hard-to-tune

hyper-parameters [17] such as the granularity of the seeds and various kinds of thresholds [1, 41]. They also lack accuracy guarantees.

A few recent works propose to use the min-hash techniques [6] for near-duplicate text alignment [15, 33, 34, 52]. They introduce the concept of “compact windows” to group nearby subsequences sharing the same min-hash. It has been shown that when duplicate tokens have the same hash value, the $O(kn^2)$ min-hashes in a text with n tokens can be compressed in compact windows using $O(kn)$ space and $O(kn \log n)$ time, where k is the sketch size [15]. Along this line, an algorithm is developed to group the $O(n^2)$ bottom- k sketches in a text with n tokens into compact windows using $O(nk^2)$ space and $O(n \log n + nk)$ time [52]. To further reduce the space cost, another algorithm is designed to group the $O(kn^2)$ one-permutation hashing [28] min-hash sketches into compact windows using $O(n+k)$ space and $O(n \log n + k)$ time [34]. These algorithms are used to evaluate the memorization behaviour in large language models (LLMs), revealing that up to 10% of texts generated by GPT-2 [40] had near-duplicates in its training data. It also shows that the min-hash sketches of all subsequences no shorter than t tokens in a text with n tokens can be grouped into $O(\frac{n}{t})$ compact windows on average [33]. However, none of these works can deal with weighted min-hash (i.e., consistent weighted sampling [25]).

Min-Hash Sketch. Min-hash was originally introduced in statistics for coordinated sampling [4] and later adapted for database applications by Flajolet and Martin [16]. Broder [7] employed the min-hash sketch to detect near-duplicate web pages. Several variants of the min-hash sketch have been proposed to improve the sketching time of the classic min-hash sketch, including the bottom- k sketch [49], one-permutation hashing (OPH) [28], and fast similarity sketch [12]. The number of min-hashes generated by one-permutation hashing is not fixed. A few works propose techniques to address this issue [45–47].

Weighted Jaccard Similarity Estimation. Many techniques have been proposed to estimate the weighted Jaccard similarity [14, 25, 30, 54]. Specifically, [21] extends the classic min-hash to estimate the multi-set Jaccard similarity. A method dealing with integer weights is mentioned in [10]. It was extended by [11] to support a more general weight function. Consistent weighted sampling (CWS) is first proposed in [30] to estimate weighted Jaccard similarity. Ioffe proposes improved consistent weighted sampling, which simplified CWS and guarantees worst-case constant time for each non-zero weight [25]. Shrivastava proposes to use rejected sampling to estimate the weighted Jaccard Similarity [44].

8 Conclusion

In conclusion, this paper extends near-duplicate text alignment to support weighted Jaccard similarity by leveraging consistent weighted sampling. We introduce MonoActive, an efficient and theoretically optimal algorithm for grouping subsequences based on their consistent weighted samplings. Our analysis establishes tight bounds on the number of groups generated, and our experiments demonstrate substantial improvements over state-of-the-art methods in both index time and index size. These results highlight the practicality and scalability of our approach for real-world text alignment tasks where token weights matter.

References

- [1] E. Agirre, C. Banea, D. M. Cer, M. T. Diab, A. Gonzalez-Agirre, R. Mihalcea, G. Rigau, and J. Wiebe. Semeval-2016 task 1: Semantic textual similarity, monolingual and cross-lingual evaluation. In *SEMEVAL*, pages 497–511. The Association for Computer Linguistics, 2016.
- [2] S. F. Altschul, W. Gish, W. Miller, E. W. Myers, and D. J. Lipman. Basic local alignment search tool. *Journal of Molecular Biology*, 215(3):403–410, 1990.
- [3] R. A. Baeza-Yates and B. A. Ribeiro-Neto. *Modern Information Retrieval*. ACM Press / Addison-Wesley, 1999.
- [4] K. R. W. Brewer, L. J. Early, and S. F. Joyce. Selecting several samples from a single population. *Australian Journal of Statistics*, 14(3):231–239, 1972.
- [5] S. Brin, J. Davis, and H. Garcia-Molina. Copy detection mechanisms for digital documents. In *SIGMOD*, pages 398–409. ACM Press, 1995.
- [6] A. Z. Broder. On the resemblance and containment of documents. In *SEQUENCES*, pages 21–29. IEEE, 1997.
- [7] A. Z. Broder, S. C. Glassman, M. S. Manasse, and G. Zweig. Syntactic clustering of the web. *Comput. Networks*, 29(8-13):1157–1166, 1997.
- [8] A. Z. Broder, S. C. Glassman, M. S. Manasse, and G. Zweig. Syntactic clustering of the web. *Computer networks and ISDN systems*, 29(8-13):1157–1166, 1997.
- [9] N. Carlini, D. Ippolito, M. Jagielski, K. Lee, F. Tramèr, and C. Zhang. Quantifying memorization across neural language models. *CoRR*, abs/2202.07646, 2022.
- [10] M. S. Charikar. Similarity estimation techniques from rounding algorithms. In *Proceedings of the thirty-fourth annual ACM symposium on Theory of computing*, pages 380–388, 2002.
- [11] O. Chum, J. Philbin, A. Zisserman, et al. Near duplicate image detection: Min-hash and tf-idf weighting. In *Bmvc*, volume 810, pages 812–815, 2008.
- [12] S. Dahlgaard, M. B. T. Knudsen, and M. Thorup. Fast similarity sketching. In *2017 IEEE 58th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 663–671. IEEE, 2017.
- [13] H. Ding, J. Zhai, D. Deng, and S. Ma. The case for learned provenance graph storage systems. In J. A. Calandrino and C. Troncoso, editors, *32nd USENIX Security Symposium, USENIX Security 2023, Anaheim, CA, USA, August 9-11, 2023*, pages 3277–3294. USENIX Association, 2023.
- [14] O. Ertl. Bagminhash-minwise hashing algorithm for weighted sets. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 1368–1377, 2018.
- [15] W. Feng and D. Deng. Allign: Aligning all-pair near-duplicate passages in long texts. In *SIGMOD*, pages 541–553. ACM, 2021.
- [16] P. Flajolet and G. N. Martin. Probabilistic counting algorithms for data base applications. *Journal of computer and system sciences*, 31(2):182–209, 1985.
- [17] T. Foltýnek, N. Meuschke, and B. Gipp. Academic plagiarism detection: A systematic literature review. *ACM Comput. Surv.*, 52(6):112:1–112:42, 2020.
- [18] A. Franz and T. Brants. All our n-gram are belong to you. *Google Machine Translation Team*, 20, 2006.
- [19] P. Gage. A new algorithm for data compression. *C Users J*, 12(2):23–38, feb 1994.
- [20] A. Gokaslan and V. Cohen. Openwebtext corpus.
- [21] S. Gollapudi and R. Panigrahy. Exploiting asymmetry in hierarchical topic extraction. In *Proceedings of the 15th ACM international conference on Information and knowledge management*, pages 475–482, 2006.
- [22] O. A. Hamid, B. Behzadi, S. Christoph, and M. R. Henzinger. Detecting the origin of text segments efficiently. In *WWW*, pages 61–70. ACM, 2009.
- [23] T. C. Hoar and J. Zobel. Methods for identifying versioned and plagiarized documents. *J. Assoc. Inf. Sci. Technol.*, 54(3):203–215, 2003.
- [24] S. Ioffe. Improved consistent sampling, weighted minhash and L1 sketching. In *ICDM*, pages 246–255. IEEE Computer Society, 2010.
- [25] S. Ioffe. Improved consistent sampling, weighted minhash and L1 sketching. In *2010 IEEE international conference on data mining*, pages 246–255. IEEE, 2010.
- [26] J. W. Kim, K. S. Candan, and J. Tatemura. Efficient overlap and content reuse detection in blogs and online news articles. In *WWW*, pages 81–90. ACM, 2009.
- [27] K. Lee, D. Ippolito, A. Nystrom, C. Zhang, D. Eck, C. Callison-Burch, and N. Carlini. Deduplicating training data makes language models better. In *ACL*, pages 8424–8445, 2022.
- [28] P. Li, A. B. Owen, and C. Zhang. One permutation hashing. In *NIPS*, pages 3122–3130, 2012.
- [29] I. Magar and R. Schwartz. Data contamination: From memorization to exploitation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, *ACL 2022, Dublin, Ireland, May 22-27, 2022*, pages 157–165. Association for Computational Linguistics, 2022.
- [30] M. Manasse, F. McSherry, and K. Talwar. Consistent weighted sampling. *Unpublished technical report* <http://research.microsoft.com/en-us/people/manasse,2>, 2010.
- [31] U. Manber. Finding similar files in a large file system. In *USENIX Winter 1994 Technical Conference, San Francisco, California, USA, January 17-21, 1994, Conference Proceedings*, pages 1–10. USENIX Association, 1994.
- [32] G. Pasi and G. Bordogna. Introduction to the minitrack on intelligent information access and retrieval. In T. X. Bui, editor, *56th Hawaii International Conference on System Sciences, HICSS 2023, Maui, Hawaii, USA, January 3-6, 2023*, pages 4169–4170. ScholarSpace, 2023.
- [33] Z. Peng, Z. Wang, and D. Deng. Near-duplicate sequence search at scale for large language model memorization evaluation. *Proc. ACM Manag. Data*, 1(2):179:1–179:18, 2023.
- [34] Z. Peng, Y. Zhang, and D. Deng. Near-duplicate text alignment with one permutation hashing. *Proc. ACM Manag. Data*, 2(4):200:1–200:26, 2024.
- [35] M. Potthast, A. Barrón-Cedeño, A. Eiselt, B. Stein, and P. Rosso. Overview of the 2nd international competition on plagiarism detection. In *CLEF 2010 Labs and Workshops, Notebook Papers*, volume 1176 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2010.
- [36] M. Potthast, A. Eiselt, A. Barrón-Cedeño, B. Stein, and P. Rosso. Overview of the 3rd international competition on plagiarism detection. In *CLEF 2011 Labs and Workshop, Notebook Papers*, volume 1177 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2011.
- [37] M. Potthast, T. Gollub, M. Hagen, J. Kiesel, M. Michel, A. Oberländer, M. Tippmann, A. Barrón-Cedeño, P. Gupta, P. Rosso, and B. Stein. Overview of the 4th international competition on plagiarism detection. In *CLEF 2012 Evaluation Labs and Workshop, Notebook Papers*, volume 1178 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2012.
- [38] M. Potthast, M. Hagen, A. Beyer, M. Busse, M. Tippmann, P. Rosso, and B. Stein. Overview of the 6th international competition on plagiarism detection. In *Working Notes for CLEF 2014 Conference*, volume 1180 of *CEUR Workshop Proceedings*, pages 845–876. CEUR-WS.org, 2014.
- [39] M. Potthast, M. Hagen, T. Gollub, M. Tippmann, J. Kiesel, P. Rosso, E. Stammatos, and B. Stein. Overview of the 5th international competition on plagiarism detection. In *Working Notes for CLEF 2013 Conference*, volume 1179 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2013.
- [40] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [41] M. A. Sanchez-Perez, A. F. Gelbukh, and G. Sidorov. Adaptive algorithm for plagiarism detection: The best-performing approach at PAN 2014 text alignment competition. In *6th International Conference of the CLEF Association*, volume 9283 of *Lecture Notes in Computer Science*, pages 402–413. Springer, 2015.
- [42] S. Schleimer, D. S. Wilkerson, and A. Aiken. Winnowing: Local algorithms for document fingerprinting. In *SIGMOD*, pages 76–85. ACM, 2003.
- [43] J. Seo and W. B. Croft. Local text reuse detection. In *SIGIR*, pages 571–578. ACM, 2008.
- [44] A. Shrivastava. Simple and efficient weighted minwise hashing. *Advances in Neural Information Processing Systems*, 29, 2016.
- [45] A. Shrivastava. Optimal densification for fast and accurate minwise hashing. In *International Conference on Machine Learning*, pages 3154–3163. PMLR, 2017.
- [46] A. Shrivastava and P. Li. Densifying one permutation hashing via rotation for fast near neighbor search. In *International Conference on Machine Learning*, pages 557–565. PMLR, 2014.
- [47] A. Shrivastava and P. Li. Improved densification of one permutation hashing. *arXiv preprint arXiv:1406.4784*, 2014.
- [48] A. Singhal et al. Modern information retrieval: A brief overview. *IEEE Data Eng. Bull.*, 24(4):35–43, 2001.
- [49] M. Thorup. Bottom-k and priority sampling, set similarity and subset sums with minimal independence. In *STOC*, pages 371–380. ACM, 2013.
- [50] T. Vu, X. He, G. Haffari, and E. Shareghi. Koala: An index for quantifying overlaps with pre-training corpora. In Y. Feng and E. Lefever, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023 - System Demonstrations, Singapore, December 6-10, 2023*, pages 90–98. Association for Computational Linguistics, 2023.
- [51] P. Wang, C. Xiao, J. Qin, W. Wang, X. Zhang, and Y. Ishikawa. Local similarity search for unstructured text. In *SIGMOD*, pages 1991–2005. ACM, 2016.
- [52] Z. Wang, C. Zuo, and D. Deng. Txtalign: Efficient near-duplicate text alignment search via bottom-k sketches for plagiarism detection. In *SIGMOD*, pages 1146–1159. ACM, 2022.
- [53] J. J. Webster and C. Kit. Tokenization as the initial phase in nlp. In *COLING 1992 volume 4: The 14th international conference on computational linguistics*, 1992.
- [54] W. Wu, B. Li, L. Chen, C. Zhang, and S. Y. Philip. Improved consistent weighted sampling revisited. *IEEE Transactions on Knowledge and Data Engineering*, 31(12):2332–2345, 2018.

A Correctness of Monotonic Partitioning

This section serves as a proof of Theorem 1. Let $n = |T|$.

Augmented Skyline. Since every key $(p, q) \in K(T)$ has $(p, q) \in [1, n] \times [1, n]$, the two guard keys $(0, 0)$ and $(n + 1, n + 1)$ neither dominate nor are dominated by any key in $K(T)$. For any key set $L \subseteq K(T)$, set $\{(0, 0), (n + 1, n + 1)\} \cup S(L)$ is a skyline – no key in the set dominates another. We refer to this set as the augmented skyline of L . All the previous discussions about a skyline naturally apply to the corresponding augmented skyline.

By Definition 8, $\text{rec}(\{(0, 0), (n + 1, n + 1)\} \cup S(L)) = \text{rec}(S(L)) \cup \text{rec}(0, 0) \cup \text{rec}(n + 1, n + 1) = \text{rec}(S(L))$. This is because $\text{rec}(0, 0) = [1, 0] \times [0, n] = \phi$ and $\text{rec}(n + 1, n + 1) = [1, n + 1] \times [n + 1, n] = \phi$.

Lemma 7 tests the existence of any key in L that dominates (b, c) . If not, we need to update the augmented skyline to include (b, c) and, at the same time, generate compact windows according to Lemma 14 below. Please see Figure 3(b) as an illustration.

LEMMA 14. *Given a key set $L \subseteq K(T)$ and a key $(b, c) \notin L$. Let S be the augmented skyline of L in coordinate order. Let*

- *i be the largest index such that $S[i].y < c$ and*
- *j be the smallest index such that $S[j].x > b$.*

The two guard keys in S ensure that i and j exist. We claim C1 & C2.

C1 *If $i + 1 \leq j - 1$, then the keys in S dominated by (b, c) are exclusively $S[i + 1], S[i + 2], \dots, S[j - 1]$; otherwise, there is no key in S dominated by (b, c) .*

C2 *If (b, c) is not dominated by any key in S , then*

$$r_k \doteq \begin{cases} [S[k].x + 1, b] \times [c, S[k + 1].y - 1], & k = i, \\ [S[k].x + 1, b] \times [S[k].y, S[k + 1].y - 1], & k \in [i + 1, j - 1]. \end{cases}$$

are $j - i$ mutually disjoint rectangles that jointly cover the growing part of the skyline when including (b, c) , i.e.,

$$r_i \cup r_{i+1} \cup \dots \cup r_{j-1} = \text{rec}(b, c) \setminus \text{rec}(S).$$

PROOF. We prove C1 in two complementary cases.

Case $i + 1 > j - 1$. Assume there is a key $S[k]$ that is dominated by (b, c) . Then $S[k].x \leq b < S[j].x$ and $S[i].y < c \leq S[k].y$. By Lemma 6, $i < k$ and $k < j$. Hence $i \leq k - 1 \leq j - 2$, which contradicts $i + 1 > j - 1$. Thus, no key in S is dominated by (b, c) .

Case $i + 1 \leq j - 1$. For any $k \notin [i + 1, j - 1]$, either $k \leq i$ or $k \geq j$. By the definition of i , $S[k].y < c$ for $k \leq i$, and by the definition of j , $S[k].x > b$ for $k \geq j$. For any $k \leq i$ or $k \geq j$, (b, c) does not dominate $S[k]$. For any $k \in [i + 1, j - 1]$, by the definitions of i and j , we have $S[k].x \leq b$ and $S[k].y \geq c$. Since $(b, c) \notin L$, we must have $(b, c) \neq S[k]$. Thus $[b, c] \subset [S[k].x, S[k].y]$. Thus, (b, c) dominates $S[k]$ for $\forall k \in [i + 1, j - 1]$. This proves C1.

We prove C2 below by considering when (b, c) is not dominated by any key in S . We first show $S[i + 1]$ exists in the augmented skyline by showing $i < j$. If otherwise, i.e., $i \geq j$, by Lemma 6, $S[i].x \geq S[j].x > b$, in addition, $S[i].y < c$ by definition, thus $S[i]$ dominates (b, c) , contradiction. $S[j]$ exists and thus $S[i + 1]$ exists.

We then show $r_i, r_{i+1}, \dots, r_{j-1}$ are mutually disjoint: the y -ranges of these rectangles, i.e., $[c, S[i + 1].y - 1], [S[i + 1].y, S[i + 2].y], \dots, [S[j - 2].y, S[j - 1].y], [S[j - 1].y, S[j].y]$ are mutually disjoint.

Next, consider any $(p, q) \in \text{rec}(b, c) \setminus \text{rec}(S)$, we show there exists $k \in [i, j - 1]$ such that $(p, q) \in r_k$. Since $(p, q) \in \text{rec}(b, c)$ and thus $q \in [c, n]$, q will fall in exactly one of the 3 cases below.

Case 1 : $q \in [c, S[i + 1].y - 1]$,

Case 2 : $q \in [S[k].y, S[k + 1].y - 1]$ for some $k \in [i + 1, j - 1]$,

Case 3 : $q \in [S[j].y, n]$.

In Case 1, the definition of i ensures that $S[i].y < c$. Besides, for any $S[u]$ in the augmented skyline, $S[u].y - 1 \leq n$, and thus $[c, S[i + 1].y - 1] \subseteq [S[i].y, n]$. (p, q) is not dominated by any key in S including $S[i]$ and $q \in [c, S[i + 1].y - 1]$, $p \notin [1, S[i].x]$. As $(p, q) \in \text{rec}(b, c)$, $p \in [1, b]$, and thus $p \in [S[i].x + 1, b]$. Thus, $(p, q) \in r_i$.

In Case 2, consider $k \in [i + 1, j - 1]$ such that $q \in [S[k].y, S[k + 1].y - 1]$. As $(p, q) \notin \text{rec}(S[k]) = [1, S[k].x] \times [S[k].y, n]$ and $p \in [S[k].y, S[k + 1].y - 1] \subseteq [S[k].y, n]$, $p \notin [1, S[k].x]$. Together with $p \in [1, b]$, we have $p \in [S[k].x + 1, b]$. Thus $(p, q) \in r_k$.

Lastly, we show Case 3 is impossible. If $S[j] = (n + 1, n + 1)$, then $[S[j].y, n] = \phi$; otherwise, since (p, q) is not dominated by $S[j]$ and $q \in [S[j].y, n]$, we have $p \notin [1, S[j].x]$. As $p \in [1, b]$, $S[j].x > b \geq p > S[j].x$, contradiction. Thus, Case 3 is impossible.

In conclusion, $S[i + 1]$ exists and thus the ranges of the three cases are well defined. Case 3 is impossible, thus q can only fall in Case 1 or 2, and in either case, there is $k \in [i, j - 1]$ s.t. $(p, q) \in r_k$. \square

Now we prove Theorem 1. Firstly, it is trivial to verify GENERATEKEYS produces $K(T)$. Algorithm 4 examines all keys in $K(T)$ in the increasing order of their hash values. Let (b, c) with hash value v be the current visiting key and L be the set of visited keys. We prove by induction that at the end of each for-loop of Algorithm 4, S is the augmented skyline of $L \cup \{(b, c)\}$ and the set of compact windows $\mathcal{P}(b, c)$ generated during the loop constitutes a partition of $\text{rec}(b, c) \setminus \text{rec}(S)$. Thus the union of all generated compact windows forms a partition of $\text{rec}(S)$. Note that a guard key (p, q) has $\text{rec}(p, q) = \phi$, thus $\text{rec}(S) = \text{rec}(S(L))$.

We assume S is the augmented skyline of L at the beginning of each for loop. This is true at the beginning of the first loop, when $L = \phi$ and $S = \{(0, 0), (n + 1, n + 1)\}$. Clearly, S is the augmented skyline of L . Next, we consider two cases for each loop.

In the case there exists a key in S dominates (b, c) , by Lemma 7, Line 5 evaluates to true and the loop ends¹. S remains unchanged. The augmented skyline of $L \cup \{(b, c)\}$ is S as (b, c) is dominated. $\mathcal{P}(b, c) = \phi$ is a partition of $\text{rec}(b, c) \setminus \text{rec}(S) = \phi$.

In the case (b, c) is not dominated by any key in $S(L)$, by Lemma 7, Line 5 evaluates to false. Furthermore, by Lemma 14(C1), all keys dominated by (b, c) in S (if there is any) are precisely $S[i + 1], \dots, S[j - 1]$. By the definition of skyline, $S \cup \{(b, c)\} \setminus \{S[i + 1], \dots, S[j - 1]\}$ must be the augmented skyline of $L \cup \{(b, c)\}$. Next, we show that $\mathcal{P}(b, c)$ generated during the loop is a partition of $\text{rec}(b, c) \setminus \text{rec}(L)$.

Consider Lemma 14 (C2) and determine i, j and r_k (denoted as $[a_k, b_k] \times [c_k, d_k]$), $k \in [i, j - 1]$ accordingly. Note that $r_k \neq \phi$ only if $a_k \leq b_k$ and $c_k \leq d_k$. The tuples added to $\mathcal{P}(T, h)$ in Lines 9-13 are

$$\mathcal{P}(b, c) = \{W_k \doteq \langle T, h, v, a_k, b_k, c_k, d_k \rangle | k \in [i, j - 1], r_k \neq \emptyset\}.$$

First, any $W_k \in \mathcal{P}(b, c)$ is a compact window because by Lemma 3 and Lemma 4, all subsequences in $\text{rec}(b, c) \setminus \text{rec}(S(L))$ have min-hash v . Secondly, all the compact windows W_k in $\mathcal{P}(b, c)$ constitute a partition of $\text{rec}(b, c) \setminus \text{rec}(S)$ because by Lemma 14 (C2), the rectangles $\{r_k | k \in [i, j - 1]\}$, form a partition of $\text{rec}(b, c) \setminus \text{rec}(S)$.

¹When using Lemma 7 on $S(L)$, j does not exist only if Line 5 on augmented skyline S has $S[j']$ taking the guard key $(0, 0)$ – Line 5 also tests false; otherwise, the same key is returned, i.e., $S(L)[j] = S[j']$. Thus, Line 5 correctly facilitates Lemma 7.

Since along the loop, the increments $\text{rec}(b, c) \setminus \text{rec}(S)$ are disjoint, the union of all generated $\mathcal{P}(b, c)$ forms a partition of $\text{rec}(S(L))$.

The resulting $\mathcal{P}(T, h)$ satisfies the coverage condition in Definition 3: for every subsequence $T[x, y]$ of T , (x, y) is in $\text{rec}(S(L))$ and is thus covered by exactly one compact window. Denote by v the hash value of subsequence $T[x, y]$. By Lemma 3, $T[x, y]$ must contain a key, say (b, c) , with hash value v . Thus, $(x, y) \in \text{rec}(b, c) \subseteq \text{rec}(S(L))$. This completes the proof of Theorem 1.

LEMMA 15. *Algorithm 4 with GENERATEACTIVEKEYS produces a partition of T .*

PROOF. It is easy to verify GENERATEACTIVEKEYS produces all the active keys $X(T)$. Based on Lemma 8, for any non-active key (b, c) , it is skipped in our algorithm, which means $\mathcal{P}(b, c) = \emptyset$. Based on the proof of Theorem 1, Algorithm 4 with GENERATEACTIVEKEYS produces a partitioning of T . \square

B Lower Bound Analysis

This section analyzes the lower bound of the partition generation problem. A hard case is when all the tokens are the same, i.e., there is only one token t in the text T . Next we show every partition $\mathcal{P}(T, h)$ has $\Omega(n \log n)$ compact windows in expectation in the worst case. Formally, we define two subsequences are *mergeable* as below.

DEFINITION 15 (MERGEABLE). *Given text T and hash function h , two subsequences $T[i, j]$ and $T[p, q]$ are mergeable if they can be represented by the same compact window $\langle T, h, v, a, b, c, d \rangle$, i.e., $i, p \in [a, b]$, $b \leq c, j, q \in [c, d]$, and $h(T[x, y]) = v, \forall x \in [a, b], y \in [c, d]$.*

Clearly, for a set of subsequences from the same text that are mutually not mergeable, it needs at least one compact window in the partition for each subsequence in the set. Thus we analyze the lower bound of the partition generation problem by constructing a set of subsequences that are mutually not mergeable.

LEMMA 16. *For a text T of length n of one token t , i.e., all tokens in T are duplicate, any partition $\mathcal{P}(T, h)$ of T contains $\Omega(n + n \log n)$ compact windows in expectation. The expectation is introduced by the randomness of the hash function as opposed to any other assumption.*

PROOF. Let $h(t, x_1), h(t, x_2), \dots, h(t, x_m)$ be all the active hash of the token t where $1 \leq x_1 < x_2 < \dots < x_m \leq n = |T|$. That means, $h(t, x_j) > h(t, x_i)$ for every $1 \leq j < i \leq m$. Construct a set of subsequences $S = S_1 \cup S_2 \cup \dots \cup S_m$ where $S_i = \{T[p, q] \mid q - p + 1 = x_i\}$. That is to say, S_i is the set of all subsequences in T containing exactly x_i tokens. Consider two sets S_i and S_j where $1 \leq i \neq j \leq m$. First, they must be disjoint as the lengths of their subsequences are different. Furthermore, based on the definition of active hash, the min-hash of subsequences in S_i and S_j are respectively $h(t, x_i)$ and $h(t, x_j)$. Since $x_i \neq x_j$, the subsequences from S_i and S_j are not mergeable based on Definition 15.

Assume that there are two distinct subsequences $T[p_1, q_1]$ and $T[p_2, q_2]$ from the same set S_i that are mergeable. As they all of the same length and are distinct, either $p_1 < p_2$ or $p_1 > p_2$. Without loss of generality, assume $p_1 < p_2$. As the two subsequences are mergeable, there exists a compact window $\langle T, h, v, a, b, c, d \rangle$ representing both subsequences. Thus $p_1 < p_2 \leq b \leq c \leq q_1$.

The subsequence $T[b, c]$ then contains less than x_i tokens because $c - b + 1 \leq q_1 - p_2 + 1 = (q_1 - p_1 + 1) + p_1 - p_2 = x_i + p_1 - p_2 < x_i$. Thus $h(T[b, c]) > x_i$ since $h(t, x_i)$ is active hash, contradicting the definition of a compact window, i.e., $h(T[b, c]) = h(T[p_1, q_1]) = h(t, x_i)$.

Based on the discussion above, the subsequences in S are mutually not mergeable. Next, we calculate the expected size of S . There are $n - x_i + 1$ subsequences in S_i . Thus $|S| = \sum_{i=1}^m (n - x_i + 1)$. Let us consider an arbitrary index $j \in [1, n]$. We observe that $j = x_i$ for some $i \in [1, m]$ if and only if $h(t, j)$ is the smallest among $h(t, 1), h(t, 2), \dots, h(t, j)$. The probability is $\frac{1}{j}$. Thus we have

$$\begin{aligned} \mathbb{E}[|S|] &= \mathbb{E}\left[\sum_{i=1}^m (n - x_i + 1)\right] = \sum_{i=1}^n (n - i + 1) \frac{1}{i} \\ &= (n + 1) \sum_{i=1}^n \frac{1}{i} - n \geq (n + 1) \ln n - n = \Omega(n \log n). \end{aligned} \quad (7)$$

Note that when $n = 1$, $\mathbb{E}[|S|] = 1$, besides, $n \leq n \log n$ when $n \geq 2$. So $\mathbb{E}[|S|] = \Omega(n \log n + n)$. The subsequences in S are mutually not mergeable. One compact window is needed to represent each of the subsequence in S . Thus any partition $\mathcal{P}(T, h)$ contains $\Omega(n \log n + n)$ compact windows in expectation. \square

Lemma 16 shows a hard case in terms of partition size. The reason we keep factor n in the lower bound is for the special case when $n = 1$. This will be useful when we conduct frequency-aware analysis below.

When $f_T \ll n$. Next we consider the cases when we know the length and the maximum token frequency f_T of a text.

LEMMA 17. *Consider integer $n > 0$ and $f_T \in [1, n]$, there is a text T of length n and maximum frequency f_T such that any partition $\mathcal{P}(T, h)$ of T contains $\Omega(n + n \log f_T)$ compact windows in expectation.*

PROOF. Without loss of generality, assume that n is a multiple of f_T . Let m be $\frac{n}{f_T}$. Let $T = t_1 t_1 \dots t_1 t_2 t_2 \dots t_2 \dots t_m t_m \dots t_m$ where each t_i (where $1 \leq i \leq m$) appears f_T times next to each other in T .

We construct a set of subsequences $S = S^1 \cup S^2 \cup \dots \cup S^m$ where S^i is constructed in the same way as in the proof in Lemma 16 for t_i . The subsequences within each S^i are mutually not mergeable as shown earlier. Next we show any two subsequences $T[p_1, q_1]$ and $T[p_2, q_2]$ from S^i and S^j (where $i \neq j$) are not mergeable. This is obvious as $h(T[p_1, q_1]) = h(t_i, x)$ for some $x \in [1, f_i]$ and $h(T[p_2, q_2]) = h(t_j, y)$ for some $y \in [1, f_j]$. However, $h(t_i, x) \neq h(t_j, y)$ assuming the universal hash function h has no collision.

Based on the discussion above, it needs $\Omega(f_T(1 + \log f_T))$ compact windows to represent all the subsequences $T[p, q]$ where $T[p] = [q] = t_i$ for each $1 \leq i \leq m$. Thus in total, it needs $\Omega(\sum_{i \in [m]} f_T(1 + \log f_T)) = \Omega(n + n \log f_T)$ compact windows to represent all the subsequences in T . Thus every partition $\mathcal{P}(T, h)$ has $\Theta(n + n \log f_T)$ compact windows. \square

Now we prove Theorem 3. Theorem 1 shows that the expected number of compact windows of Algorithm 4 is $O(n + n \log f_T)$ while Lemma 17 shows a case of T , given n and f_T , such that no partition can have the number of compact windows smaller than $\Omega(n + n \log f_T)$ in expectation. Thus, given n and f_T , Algorithm 4 is optimal in the worst case.