

Approximate Hamming Distance in a Stream

Raphaël Clifford¹ and Tatiana Starikovskaya²

1 University of Bristol, Bristol, UK

raphael.clifford@bristol.ac.uk

2 University of Bristol, Bristol, UK

tat.starikovskaya@gmail.com

Abstract

We consider the problem of computing a $(1+\varepsilon)$ -approximation of the Hamming distance between a pattern of length n and successive substrings of a stream. We first look at the one-way randomised communication complexity of this problem. We show the following:

- If Alice and Bob both share the pattern and Alice has the first half of the stream and Bob the second half, then there is an $\mathcal{O}(\varepsilon^{-4} \log^2 n)$ bit randomised one-way communication protocol.
- If Alice has the pattern, Bob the first half of the stream and Charlie the second half, then there is an $\mathcal{O}(\varepsilon^{-2} \sqrt{n} \log n)$ bit randomised one-way communication protocol.

We then go on to develop small space streaming algorithms for $(1 + \varepsilon)$ -approximate Hamming distance which give worst case running time guarantees per arriving symbol.

- For binary input alphabets there is an $\mathcal{O}(\varepsilon^{-3} \sqrt{n} \log^2 n)$ space and $\mathcal{O}(\varepsilon^{-2} \log n)$ time streaming $(1 + \varepsilon)$ -approximate Hamming distance algorithm.
- For general input alphabets there is an $\mathcal{O}(\varepsilon^{-5} \sqrt{n} \log^4 n)$ space and $\mathcal{O}(\varepsilon^{-4} \log^3 n)$ time streaming $(1 + \varepsilon)$ -approximate Hamming distance algorithm.

1998 ACM Subject Classification F.2 Analysis of algorithms and problem complexity

Keywords and phrases Hamming distance, communication complexity, data stream model

Digital Object Identifier 10.4230/LIPIcs.ICALP.2016.20

1 Introduction

We study the complexity of one of the most basic problems in pattern matching, that of approximating the Hamming distance. Given a pattern P of length n the task is to output a $(1 + \varepsilon)$ -approximation of the Hamming distance between P and every n -length substring of a longer text. We provide the first efficient one-way randomised communication protocols as well as a new, fast and space efficient streaming algorithm for this problem.

The general task of efficiently computing the Hamming distances offline between a pattern and a text has been studied for many years. When the input is binary and the text has length proportional to that of the pattern, then all outputs can be computed exactly in $\mathcal{O}(n \log n)$ time by repeated application of the fast Fourier transform [14]. For larger alphabets, $\mathcal{O}(n \sqrt{n \log n})$ time solutions were first discovered in the 1980s [1, 22]. The fastest randomised algorithm for $(1 + \varepsilon)$ -approximate Hamming distance computation for large alphabets was due for many years to Karloff from 1993 [20] running in $\mathcal{O}(\varepsilon^{-2} n \log^2 n)$ time overall. In a breakthrough paper in 2015 a new algorithm was given improving the time complexity to $\mathcal{O}(\varepsilon^{-1} n \log^3 n \log \varepsilon^{-1})$ [21]. These fast methods all require linear space and up until this point no sublinear space solutions have been known.

The first basic question that arises is whether it is in fact possible to give a $(1 + \varepsilon)$ -approximation to the Hamming distance in a streaming setting while using only sublinear space. In order to explore this question we start our study by considering two natural



© Raphaël Clifford, Tatiana Starikovskaya;
licensed under Creative Commons License CC-BY

43rd International Colloquium on Automata, Languages, and Programming (ICALP 2016).

Editors: Ioannis Chatzigiannakis, Michael Mitzenmacher, Yuval Rabani, and Davide Sangiorgi;
Article No. 20; pp. 20:1–20:14



Leibniz International Proceedings in Informatics
Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany



communication complexity problems which may also be of independent interest. Any lower bound for these communication problems will give a lower bound for the space usage of a corresponding streaming algorithm. This follows from a standard reduction where a space efficient streaming algorithm is converted into a communication protocol by taking a snapshot of memory after some symbol of the input has been read in and then sending this snapshot to the other player. On the other hand, the communication upper bounds we provide will set targets for space bounds for algorithms in the streaming setting.

Any streaming pattern matching algorithm using a pattern of length n can be reduced to repeated application of a streaming algorithm that runs on texts of length $2n$. This is done by splitting the stream into substreams of length $2n$ which overlap by n symbols. As a result we consider communication problems with these parameter settings for pattern and text length.

► **Problem 1.** *Consider a text T of length $2n$ and a pattern P of length n . Let Alice hold the information about the first half of the text and the whole of the pattern, and let Bob hold the information about the second half of the text and the whole of the pattern. Bob must output $(1 + \epsilon)$ -approximations of the Hamming distance for each alignment of P and T .*

A lower bound for the communication complexity of this problem follows from a combination of the lower bound for the communication complexity of a windowed counting problem introduced by Datar et al. in 2002 [13] and the one-way communication complexity lower bound for approximating the Hamming distance between two n -bit strings from [18].

For the first part consider the following communication problem. Assume that there is a bit vector B of length $2n$. Let Alice hold the information about the first half of B , and let Bob hold the information about the second half of B . Bob must output $(1 + \epsilon)$ -approximations of the number of set bits in each window of length n . Datar et al. showed that Alice will have to send to Bob $\Omega(\epsilon^{-1} \log^2 \epsilon^{-1} n)$ bits of information. There is a straightforward reduction from this basic counting problem to Problem 1 which then gives us the same lower bound. We set $T = B$ and $P = 00 \dots 0$ and then a $(1 + \epsilon)$ -approximation of the Hamming distance at an alignment i of P and T gives a $(1 + \epsilon)$ -approximation of the number of set bits in the window $T[i, i + n - 1]$. For the second part we use Theorem 4.1 from [18] which states that the one-way communication complexity of $(1 + \epsilon)$ -approximate Hamming distance for two strings of length n is $\Omega(\epsilon^{-2} \log n)$ for constant error probability. Combining these two lower bounds together we get a lower bound of $\Omega(\epsilon^{-2} \log n + \epsilon^{-1} \log^2 \epsilon^{-1} n)$ for the communication complexity of Problem 1.

Our first result is an efficient one-way communication protocol for Problem 1 whose complexity is only slightly higher than this lower bound. In our protocol Alice uses the fact that Bob knows the pattern as well to give an efficient encoding for parts of her half of the text which are at small Hamming distance from the pattern.

► **Theorem 1.** *Problem 1 has one-way randomised communication complexity $\mathcal{O}(\epsilon^{-4} \log^2 n)$.*

As a model for streaming pattern matching, this communication upper bound requires that a copy of the pattern is available at all times. Our main interest is however in algorithms whose total space complexity is sublinear in the pattern size. In order to model this situation more accurately we now consider a stronger three party communication problem.

► **Problem 2.** *Assume that there is a text T of length $2n$ and a pattern P of length n . Let Alice hold the information about the pattern, let Bob hold the information about the first half of the text, and let Charlie hold the information about the second half of the text. Alice will send one message to Bob who will then send one message to Charlie. Charlie must output $(1 + \epsilon)$ -approximations of the Hamming distance for each alignment of P and T .*

Somewhat surprisingly, we are still able to obtain a sublinear protocol for this new problem although the bound is higher than for the simpler Problem 1. The main technical elements of this communication protocol combine the newly introduced idea of approximate periods with succinct run-length encoded representations of the input.

► **Theorem 2.** *Problem 2 has one-way randomised communication complexity $\mathcal{O}(\frac{\sqrt{n}}{\varepsilon^2} \log n)$.*

Having investigated the communication complexity of $(1 + \varepsilon)$ -approximate Hamming distance we can now define the streaming $(1 + \varepsilon)$ -approximate Hamming distance problem.

► **Problem 3.** *Consider a pattern P of length n and a stream arriving one symbol at a time. We must output a $(1 + \varepsilon)$ -approximation to the Hamming distance between P and the latest n -length suffix of the stream as soon as a new symbol arrives. In this setting we cannot, for example, store a copy of the pattern or stream without accounting for it in our space usage.*

The upper bounds for the communication complexity of Problem 2 suggest space upper bounds we shall aim for in order to develop an optimal algorithm for the $(1 + \varepsilon)$ -approximate Hamming distance in the streaming setting. We make the first step towards this direction and show two randomised sublinear-space algorithms for the problem. We start by giving a solution for the case when both the pattern and the text are binary strings.

► **Theorem 3.** *When both P and T are binary, there is an algorithm for Problem 3 which uses $\mathcal{O}(\varepsilon^{-3} \sqrt{n} \log^2 n)$ bits of space and runs in $\mathcal{O}(\varepsilon^{-2} \log n)$ time per arriving symbol.*

The same bounds hold for alphabets of constant size σ as we can map each occurrence of the i^{th} symbol of the alphabet in the pattern or in the text to a binary string $0^{i-1}10^{\sigma-i}$, which will result in doubling the Hamming distance between the pattern and the text at each particular alignment.

For polynomial size alphabets our bounds are higher by a factor $\varepsilon^{-2} \log^2 n$ and our approach is based on the mapping idea of Karloff [20]. In that paper he showed that there exists $\Theta(\varepsilon^{-2} \log^2 n)$ mappings map_j of the alphabet onto $\{0, 1\}$ such that an $(1 + \varepsilon/3)$ -approximation of the Hamming distance between P and T at a particular alignment can be given by a normalised average of the Hamming distances between $map_j(P)$ and $map_j(T)$ at this alignment. Moreover, Karloff showed that the mappings can be generated in $\mathcal{O}(\varepsilon^{-2} \log^3 n)$ space and $\mathcal{O}(\log n)$ time per symbol. For each pattern-text pair mapped on to a binary alphabet we then run the algorithm of Theorem 3 to find $(1 + \varepsilon/3)$ -approximations and finally obtain:

► **Theorem 4.** *There is an algorithm for Problem 3 which uses $\mathcal{O}(\varepsilon^{-5} \sqrt{n} \log^4 n)$ bits of space and runs in $\mathcal{O}(\varepsilon^{-4} \log^3 n)$ time per arriving symbol.*

Our solution has guaranteed worst case complexity per arriving symbol and uses roughly the square root of the space required by the known offline $(1 + \varepsilon)$ -approximate algorithms. A key technical innovation for our space reduction is the notion we introduce of a super-sketch. This is a compact and efficiently updateable representation of consecutive text substrings which we require to be able to achieve sublinear space. For simplicity we will make the natural assumption throughout that $\varepsilon < 1/2$.

1.1 Related work and lower bounds

The one-way communication complexity of a number of variants of Hamming distance computation has been studied over the years. These includes $(1 + \varepsilon)$ -approximation [18],

the so called gap Hamming problem [9] and a bounded version known as k -mismatch [15]. However all this previous work has assumed that both Alice and Bob have strings of the same length and so need only give a single output. There has also been great interest in efficient streaming algorithms over the last 20 years, following the seminal work of [2]. In relation specifically to pattern matching problems, where space is not limited but where an output must be computed after every new symbol of the text arrives, the Hamming distance between the pattern and the latest suffix of the stream can be computed online in $O(\sqrt{n \log n})$ worst-case time per arriving symbol or $O(\sqrt{k} \log k + \log n)$ time for the k -mismatch version [11]. Both these methods however require $\Theta(n)$ space. Using the same approach, a number of other approximate pattern matching algorithms have also been transformed into efficient linear space online algorithms including [4, 3, 5, 8, 7, 6, 23]. In 2013 a small space streaming pattern matching algorithm was shown for parameterised matching [17] and in 2016 for the k -mismatch problem [10]. The latter k -mismatch paper is of particular relevance to our work. In [10] as a part of a space-efficient streaming algorithm for the k -mismatch problem, the authors presented a $(1 + \varepsilon)$ -approximate algorithm with space $\mathcal{O}(\varepsilon^{-2} k^2 \log^7 n)$ and running time $\mathcal{O}(\varepsilon^{-2} \log^5 n)$ per arriving symbol that returns a $(1 + \varepsilon)$ -approximation for all alignments of the pattern and text where the Hamming distance is at most k . The algorithm we give in this current paper can be seen a generalisation of this work, both removing the requirement for a prespecified threshold k and also using less space when $k \gtrsim n^{1/4}$.

2 Overview

In this section we give an overview of the main ideas needed for our results. We will make extensive use of sketching. Alon, Matias and Szegedy were first to show that sketching can be used to approximate frequency statistics of a stream with a particular emphasis on F_2 [2]. Later their sketching technique was generalized to allow approximation of $\|x_1 - x_2\|_p$ for two vectors x_1 and x_2 and any $p \in (0; 2]$ by Indyk et al. [16, 12]. We will use the sketches of Indyk et al. to show the communication complexity upper bounds. These sketches are based on p -stable distributions and have the advantage that they can be used even for large-size alphabets. For our streaming algorithm where we assume that the input alphabet is binary we will use simpler sketches based on the original technique of Alon et al.

2.1 Communication complexity

To show communication complexity bounds we will be using sketches based on p -stable distributions (see [16] and Sections 4.1 and 5.1 of [12]). Setting σ to be the alphabet size, a sketch of a string x is defined as a vector $sk(x)$ of length $\Theta(\varepsilon^{-2})$ such that

$$sk(x)[i] = \sum_j Y_{i,j} \cdot x[j]$$

where each $Y_{i,j}$ is drawn independently from a random stable distribution with parameter $p \leq \varepsilon / \log \sigma$. For two vectors x_1 and x_2 it can be shown that with constant probability the median of values $|sk(x_1)[i] - sk(x_2)[i]|$, appropriately scaled, is a $(1 + \varepsilon)$ -approximation of the Hamming distance. Importantly, variables $Y_{i,j}$ can be generated when we need them with the help of Nisan's pseudo-random generator, which requires only $\mathcal{O}(\log^2 n)$ random bits.

2.1.1 Problem 1 – both Alice and Bob know the pattern

The main idea of our communication complexity upper bound for Problem 1 is that if the Hamming distance between the text and the pattern at a particular alignment is (relatively) small, then Alice and Bob can use the pattern to describe the part of the text aligned with the pattern.

At each alignment the pattern can be divided into two parts – a prefix, aligned with Alice’s half of the text, and a suffix, aligned with Bob’s half of the text. Alice needs to transmit information that will help Bob approximate the Hamming distance between these different prefixes of the pattern and her half of the text. She does so by selecting a logarithmic number of prefixes of the pattern with Hamming distances $\Theta(\varepsilon^{-j})$ from the text. She then divides the part of the text aligned with each of these prefixes into blocks such that the mismatches are evenly spread across the blocks, and sends each block’s starting position and sketch to Bob.

When Bob wants to compute the Hamming distance between a prefix P' of the pattern and the text and he knows that this Hamming distance is at least $\Theta(\varepsilon^{-(j-1)})$, he uses the prefix P_j with Hamming distance $\Theta(\varepsilon^{-j})$ and the sketches of associated text blocks. The part of Alice’s text aligned with P' can be composed of several full blocks and at most one block suffix. Hamming distances between P' and the full blocks can be approximated with the help of the sketches. To approximate the Hamming distance between P' and the suffix of the block, Bob will substitute the suffix with the aligned part of P_j . As the number of mismatches between the suffix and P_j is small compared to $\Theta(\varepsilon^{-(j-1)})$, it will give a good approximation of the Hamming distance.

2.1.2 Problem 2 – only Alice knows the pattern

We start by reviewing some notation introduced in [10].

► **Definition 5.** The x -period of a string S of length n is the smallest integer $\ell > 1$ such that the Hamming distance between $S[1, n - \ell]$ and $S[\ell, n]$ is at most x .

► **Definition 6.** We define the ℓ -RLE encoding of S to be the ordered set of the run-length encodings of strings $S_i = S[i]S[\ell + i]S[2\ell + i] \dots S[\lfloor (n - i)/\ell \rfloor \cdot \ell + i]$, where $i \in [1, \ell]$. The size of the ℓ -RLE encoding is the total number of runs in the encodings of strings S_i .

► **Example 7.** Let $S = aabaabaabaabaabaac$. The 3-RLE encoding of S is: the run-length encoding $(a, 7)$ of $S_1 = aaaaaaa$, the run-length encoding $(a, 7)$ of $S_2 = aaaaaaa$, and the run-length encoding $(b, 6)(c, 1)$ of $S_3 = bbbbbc$. The size of the encoding is $1 + 1 + 2 = 4$.

Note that ℓ -RLE encoding of S is deterministic and lossless. In [10] it was also shown that if ℓ is the x -period of a string S for some integer x , then the size of the ℓ -RLE encoding is $\mathcal{O}(\ell + x)$. Intuitively, this is because each new run in the encoding of S_i corresponds to a mismatch between $S[1, n - \ell]$ and $S[\ell, n]$, and therefore there can be at most $\ell + x$ runs.

We now explain the idea of the communication protocol for Problem 2. Let the block size $B = \sqrt{n}$ and the threshold $\tau = 2\varepsilon^{-1}\sqrt{n}$. Bob will compute a sketch for each B^{th} suffix of his half of the text and send it to Charlie. Consider a particular alignment of the pattern and of the text.

Case 1: Hamming distance is large. The pattern can be divided into three parts: a prefix of length at most $B - 1$, a middle part aligned with one of the n/B sketched suffixes of Bob’s half of the text, and a suffix aligned with Charlie’s half of the text. If the Hamming

distance at the alignment is larger than τ , then the prefix can be discarded as it will change the Hamming distance by at most $B = (\varepsilon/2) \cdot \tau$. The Hamming distance between the rest of the pattern and the text can be approximated easily. Charlie has received the sketch of the middle part of the pattern as well as the sketch of the suffix of Alice's half of the text which aligns with it. Charlie can combine the sketch from Alice's part of the text with the information he has about his half of the text and then compare this sketch to the pattern sketch as required.

Case 2: Hamming distance is small. The main challenge is therefore alignments where the Hamming distance is smaller than τ . If the $(2 + \varepsilon)\tau$ -period of the pattern is larger than B , then there are at most n/B such alignments. In this case, Bob can simply send the Hamming distances for all these alignments to Bob. If the period is at most B , then Bob will find the first alignment with small Hamming distance and will use the ℓ -RLE encoding of the pattern and the full list of mismatches to describe the text. Using this description Charlie will be able to fully recover the corresponding suffix of the text and to compute the Hamming distances for all remaining alignments. The only technicality is that Bob does not know Charlie's half of the text and thus will not be able to compute the Hamming distances between the whole pattern and the text. We elaborate on this in Section 3.2.

2.2 A small space streaming algorithm

In our small space streaming algorithm we will use simpler sketches which provide a $(1 + \varepsilon)$ -approximation to the Hamming distance between two binary strings of the same length B . The method is now folklore but is essentially an application of the technique of the Johnson-Lindenstrauss lemma [19]. To do this we create a random $\varepsilon^{-2} \times B$ matrix M whose entries are from $\{-1, 1\}$. The sketch of a string x of length B is then defined to be equal to Mx , and it is known that the appropriately scaled square of the L_2 norm of the difference of the sketches of two strings gives a $(1 + \varepsilon)$ -approximation of the Hamming distance between them. We will also be using M to define sketches of strings of length $\ell < B$. In this case, we simply append the strings with $(B - \ell)$ zeros and use the method describe above. The original analysis applies here as well. Finally, we will use M to define "super-sketches" of strings of length $n - B$. Assume that a string of length $n - B$ is divided into $n/B - 1$ non-overlapping blocks of size B . A super-sketch is then defined to be a linear combination of the sketches of the blocks. We elaborate more on sketches and super-sketches in Section 4.

Now we give a high-level overview of our algorithm. The algorithm starts by preprocessing the pattern P . It computes and stores a super-sketch of each $(n - B)$ -length substring of P . The algorithm then processes the text in non-overlapping blocks of length B , computing a sketch for each block. The blocks' sketches can be maintained efficiently as we need to maintain only one sketch at a time. The algorithm also maintains a super-sketch of the last $n/B - 1$ blocks. To compute an approximation of the Hamming distance at a particular alignment, the algorithm divides the pattern into three parts: a prefix of length $(B - i)$, a middle part of length $(n - B)$, and a suffix of length i , where the middle part is aligned with a block border (see Figure 1).

The algorithm then starts by computing the $(1 + \varepsilon)$ -approximation of the Hamming distance between the middle part and the text with the help of the super-sketches. If the Hamming distance is large, the algorithm can simply discard the prefix and suffix parts. Otherwise, the algorithm also needs to approximate the Hamming distance between the prefix or the suffix of the pattern and the text. To approximate the Hamming distance between the prefix of the pattern and the text the idea is to use the information Alice transfers to

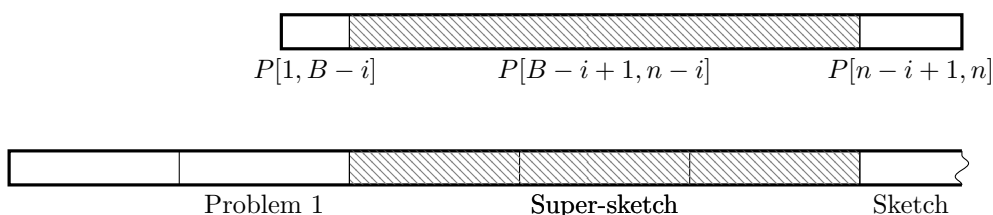


Figure 1 To estimate the Hamming distances at a particular position the algorithm uses a data structure containing the information Alice transfers to Bob in our solution for Problem 1 for the prefix $P[1, B - i]$, a super-sketch for the middle part $P[B - i + 1, n - i]$, and a sketch for the suffix $P[n - i + 1, n]$.

Bob in our solution for Problem 1. For the suffix, the algorithm will use the sketch of the part of the block between its start and the current alignment.

3 Communication complexity

In this section we show upper bounds for communication complexities of Problems 1 and 2.

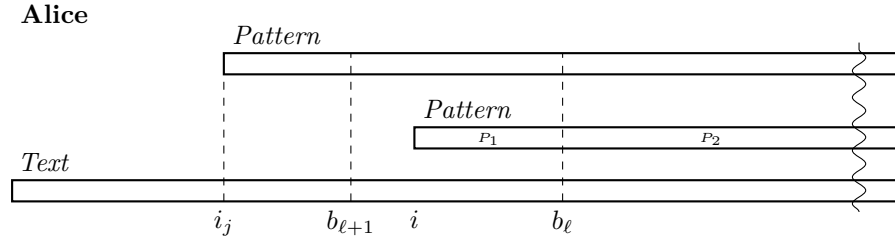
3.1 Problem 1

We start by showing an upper bound for the communication complexity of Problem 1. Remember that in this problem we have two players, Alice and Bob. Alice knows the first half of the text T and the pattern P , and Bob knows the second half of the text T and the pattern P . We will show that the communication complexity of this problem is $\mathcal{O}(\varepsilon^{-4} \log^2 n)$.

Let us first explain what Alice sends to Bob. For simplicity, we denote $k = 6/\varepsilon$. First, Alice selects $q = \lfloor \log_k n \rfloor$ positions $n \geq i_1 \geq i_2 \geq \dots \geq i_q \geq 1$ such that the Hamming distance between $T[i_j, n]$ and the prefix $P[1, n - i_j + 1]$ is at most k^{j+1} . She does this in turn starting from $j = 1$ and selecting the leftmost possible position for each j . Alice then sends to Bob $\mathcal{O}(k^2 \cdot \varepsilon^{-2} \log n) = \mathcal{O}(\varepsilon^{-4} \log n)$ bits of information for each j . She starts by dividing $T[i_j, n]$ into k^2 blocks such that the Hamming distance between each block and the corresponding substring of the pattern is at most k^{j-1} . If $n \geq b_1 > b_2 > \dots > b_{k^2} = i_j$ are the borders of the blocks, she sends Bob $b_1, b_2, \dots, b_{k^2} = i_j$ and the $(1 + \varepsilon/6)$ -approximate sketches of $T[b_\ell, n]$ for all $\ell \in [1, k^2]$. In total, Alice sends to Bob $\mathcal{O}(\varepsilon^{-4} \log^2 n)$ bits of information.

To see how Bob can use this information, consider a particular position i . At this position $P[1, n - i + 1]$ is aligned with Alice's half of the text, whereas $P[n - i + 2, n]$ is aligned with Bob's half of the text. As Bob knows the pattern, he can compute the exact Hamming distance between $P[n - i + 2, n]$ and his half of the text with no additional information. We now go on to explain how he can estimate the Hamming distance h between $P[1, n - i + 1]$ and Alice's half of the text.

Bob starts by locating the position i_j that is closest to i from the left, and the block $T[b_{\ell+1}, b_\ell]$ of $T[i_j, n]$ containing i (see Figure 2). The border b_ℓ divides the pattern into two parts, P_1 and P_2 . Let h_1 be the Hamming distance between P_1 and the text, and h_2 be the Hamming distance between P_2 and the text, $h_1 + h_2 = h$. To find a $(1 + \varepsilon)$ -approximation h'_2 of h_2 , Bob uses the sketch of $T[b_\ell, n]$. He cannot use sketches to estimate h_1 as P_1 is not aligned with the block $T[b_{\ell+1}, b_\ell]$, but he knows that there are only a few mismatches between $T[b_\ell, b_{\ell+1}]$ and the pattern aligned at the position i_j . So he estimates h_1 by computing



■ **Figure 2** Figure shows Alice's half of the text and the rightmost position $i_j < i$. Dashed lines show block borders for $T[i_j, n]$. Borders $b_{\ell+1}$ and b_ℓ are the closest to i from the left and from the right respectively. The border b_ℓ divides the pattern into two parts P_1 and P_2 . To estimate the Hamming distance h_1 between P_1 and T , Bob uses the pattern aligned at i_j . To estimate the Hamming distance h_2 between P_2 and T , he uses the sketch of $T[b_\ell, n]$.

the Hamming distance h'_1 between P_1 and the pattern aligned at the position i_j . The next lemma shows that Bob can output $h' = (h'_1 + h'_2)/(1 - \varepsilon/3)$ as a $(1 + \varepsilon)$ -approximation of h .

► **Lemma 8.** h' is a $(1 + \varepsilon)$ -approximation of h .

Proof. Remember that h'_1 is the Hamming distance between P_1 and the pattern aligned at the position i_j , and h_1 is the Hamming distance between P_1 and the text. The Hamming distance between the pattern aligned at the position i_j and $T[b_{\ell+1}, b_\ell]$ is at most k^{j-1} . Therefore,

$$h_1 - k^{j-1} \leq h'_1 \leq h_1 + k^{j-1}$$

On the other hand, h'_2 is a $(1 + \varepsilon/6)$ -approximation of h_2 . Hence,

$$h_1 + h_2 - k^{j-1} \leq h'_1 + h'_2 \leq h_1 + k^{j-1} + (1 + \varepsilon/6) \cdot h_2.$$

We now substitute $h = h_1 + h_2$ and estimate $h_2 \leq h$ to obtain

$$h - k^{j-1} \leq h'_1 + h'_2 \leq (1 + \varepsilon/6) \cdot h + k^{j-1}.$$

Finally, by our choice of i_j we have $h \geq k^{j+1}$, and therefore

$$(1 - \varepsilon/3) \cdot h \leq (1 - \varepsilon/6) \cdot h \leq h'_1 + h'_2 \leq (1 + \varepsilon/3) \cdot h.$$

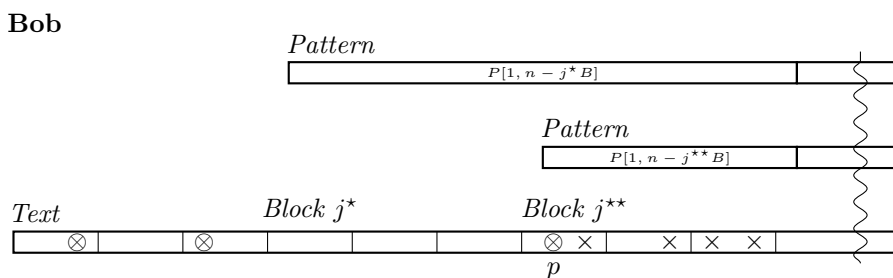
Dividing all parts of this inequality by $(1 - \varepsilon/3)$, we obtain

$$h \leq h' = (h'_1 + h'_2)/(1 - \varepsilon/3) \leq \frac{1 + \varepsilon/3}{1 - \varepsilon/3} h \leq (1 + \varepsilon) \cdot h. \quad \blacktriangleleft$$

3.2 Problem 2

In this section we show an upper bound for the communication complexity of Problem 2. Remember that in this problem we have three players, Bob, Charlie, and Alice. Bob knows the first half of the text T , Charlie knows the second half of the text T , and Alice knows the pattern P . We will show that the communication complexity of this problem is $\mathcal{O}(\varepsilon^{-2} \sqrt{n} \log n)$.

Let the block size $B = \sqrt{n}$ and the threshold $\tau = 2\varepsilon^{-1} \sqrt{n}$. We start by explaining what the players send to each other. Alice sends to Bob the following information:



■ **Figure 3** The figure shows Bob's half of the text. Crosses show alignments where the Hamming distance is at most τ . $P[1, n - j^*B]$ is the longest prefix with $\tau(2 + \varepsilon)$ -period smaller than B . Block $j^{**} \geq j^*$ is the first block containing a cross. Bob sends to Charlie sketches of text suffixes starting at blocks' borders, locations of all encircled crosses, and the last block.

1. $(1 + \varepsilon/2)$ -approximate sketches of suffixes $P[i, n]$ for all $i \in [1, B]$ (Charlie will use them to estimate large Hamming distances);
2. $(1 + \varepsilon/2)$ -approximate sketches of prefixes $P[1, n - jB]$ for all $j \in [1, n/B]$ (Bob will use them to find alignments with small Hamming distances);
3. The ℓ -RLE encoding of the longest prefix $P[1, n - j^*B]$ with $(2 + \varepsilon)\tau$ -period ℓ smaller than B (Bob will use it to describe the text).

Overall Alice sends $\mathcal{O}((n/B + B) \cdot \varepsilon^{-2} \log n + ((2 + \varepsilon)\tau + B) \cdot \log n) = \mathcal{O}(\varepsilon^{-2} \sqrt{n} \log n)$ bits of information.

Bob starts by forwarding the information he received from Alice to Charlie. Bob also sends him $(1 + \varepsilon/2)$ -approximate sketches of all suffixes $T[jB, n]$. Next, for each $j < j^*$ Bob uses the sketch of $P[1, n - jB]$ to find $(1 + \varepsilon/2)$ -approximations of Hamming distances in a block j . (Remember that Bob knows $T[1, n]$ and can compute a sketch for any its substring.) If the approximate value of the Hamming distance for some alignment is smaller than $(1 + \varepsilon/2)\tau$, he sends it to Charlie. Note that there is at most one such alignment in a block. Indeed, if we have two such alignments in the block, then the Hamming distance between the patterns at these alignments is at most $(2 + \varepsilon)\tau$, which is a contradiction with the $(2 + \varepsilon)\tau$ -period being larger than B . Moreover, Bob will not miss any alignment with the Hamming distance smaller than τ .

After that, Bob decodes $P[1, n - j^*B]$ from its ℓ -RLE encoding and computes the Hamming distances between $P[1, n - jB]$ and the text for all alignments in blocks $j \geq j^*$ precisely. He finds the first block $j^{**} \geq j^*$ where there is an alignment of $P[1, n - j^{**}B]$ with the Hamming distance at most τ . Bob sends Charlie the starting position p of this alignment and the positions of the mismatches. Finally, he sends Charlie all bits of the last block of his half of the text. Overall, Bob sends to Charlie $\mathcal{O}(\varepsilon^{-2} \sqrt{n} \log n + (\varepsilon^{-2} \log n + \log n) \cdot (n/B) + \tau) = \mathcal{O}(\varepsilon^{-2} \sqrt{n} \log n)$ bits of information.

We now explain how Charlie computes the Hamming distances. If the Hamming distance at a particular alignment i in a block $j < j^{**}$ is smaller than τ , then Charlie already knows its approximate value. If it is bigger than τ , then Charlie computes its approximation using the sketch of the longest suffix $P[jB - i, n]$ of P aligned with a block border, the sketch of $T[(j + 1)B, n]$, and $T[n + 1, 2n]$. Let h be the Hamming distance between the text and the pattern at the alignment i and let h' be the Hamming distance between $T[(j + 1)B, i + n - 1]$ and $P[jB - i, n]$.

► **Lemma 9.** $h \leq h' + B \leq (1 + \varepsilon) \cdot h$.

20:10 Approximate Hamming Distance in a Stream

Proof. The left inequality is trivial. To prove the right one, remember that $\tau \leq h$, which implies $B = (\varepsilon/2)\tau \leq (\varepsilon/2) \cdot h$. ◀

We now go on to the remaining blocks. The Hamming distances at alignments $i < p$ in the block j^{**} are bigger than τ and Charlie can find their approximation in the way described above. Charlie then decodes $P[1, n - j^*B]$ and recovers $T[p, n]$ by fixing the at most \sqrt{n} mismatches between $P[1, n - j^*B]$ and $T[p, p + n - j^{**}B + 1]$ and appending the last $p - (j^{**} - 1)B$ symbols of T (Remember that Charlie knows all symbols of the last block of $T[1, n]$). Using $T[p, n]$, $T[n + 1, 2n]$, and the sketch of P , he can approximate the Hamming distances for all alignments to the right of p .

4 Streaming algorithm

We now show a streaming algorithm for Problem 3. In this problem we are asked to output a $(1 + \varepsilon)$ -approximation of the Hamming distance between the pattern and the text at each alignment, and we do not assume that we store a copy of the pattern or of the text. For $\varepsilon < 1/2$, the algorithm uses $\mathcal{O}(\varepsilon^{-3} \sqrt{n} \log^{1.5} n)$ bits of space and its running time is $\mathcal{O}(\varepsilon^{-2} \log n)$ per arriving symbol. For simplicity, we will set $k = 1/\varepsilon > 2$ for the rest of this section.

Let $B = k\sqrt{n}$. The algorithm starts by selecting a $9k^2 \times B$ matrix M and a vector $(\sigma_1, \sigma_2, \dots, \sigma_{n/B-1})$ of i.u.d. ± 1 random variables. The algorithm then preprocesses the pattern P . It remembers the first B symbols of P , as well as a super-sketch of each $(n - B)$ -length substring of P . To compute the super-sketches the algorithm divides a substring into $(n/B - 1)$ blocks of length B , computes their sketches using M as described in Section 2, and then sums the sketches multiplying them by σ_i . The algorithm also computes sketches of the last B suffixes of P . The sketch of a suffix $P[n - i + 1, n]$ is defined to be equal to $M \cdot S_i$, where $S_i = P[n - i + 1, n] 0^{B-i}$. Finally, for each $i \in [1, B]$ and for each $j \in [0, \log_{1+\varepsilon} n]$ it stores the maximal length of pattern's prefix such that the Hamming distance between this prefix aligned at position i and the pattern is at most $(1 + \varepsilon)^j$, which takes $\mathcal{O}(\varepsilon^{-1} B \log^2 n)$ bits since $\log_{1+\varepsilon} n = \mathcal{O}(\varepsilon^{-1} \log n)$.

4.1 Text processing

The algorithm processes the text in non-overlapping blocks of length B . For each of the last n/B blocks the algorithm maintains its sketch and a data structure containing the information Alice transfers to Bob in our solution for Problem 1.

Let us start by explaining how the algorithm maintains the sketches. At the starting index of each block it initialises the block's sketch with a zero vector of length $9k^2$. When the j^{th} symbol of the block arrives, the algorithm adds the product of the j^{th} column of M and the symbol to the sketch in $\mathcal{O}(9k^2)$ time. While reading the block the algorithm also computes the super-sketch of the $(n - B)$ -length substring consisting of the $n/B - 1$ most recent blocks. Recall that the super-sketch is defined to be equal to the sum of the blocks' sketches multiplied by the variables σ_i . The total time needed for computing the sum is $\mathcal{O}(9k^2 n/B)$. The algorithm de-amortises this time over the block executing $\Omega(9k^2 n/B^2)$ steps per arriving symbol.

For each block the algorithm maintains a data structure containing the information Alice transmits to Bob in our solution for Problem 1. The algorithm starts computing the data structure when it has received the entire block. It then computes the Hamming distance between prefixes $P[1], P[1, 2], \dots, P[1, B]$ as being aligned at the right border of the block

and the block by running the fast Fourier transform algorithm on $P[1, B]$ and the block appended with B zeros, which takes $\mathcal{O}(B)$ space and $\mathcal{O}(B \log B)$ time in total [14]. The algorithm then finds i_1, i_2, \dots, i_q , where $q = \lceil \log_k B \rceil$ as defined in Problem 1 and for each i_j it computes the borders and the sketches of the blocks, where the sketches are defined with the help of the matrix M . Remember that the algorithm stores the block and the first B symbols of the pattern, so this could be done in a naive way, using symbol-by-symbol comparison. Finally, the algorithm builds binary search trees on i_1, i_2, \dots, i_q and the block borders for each i_j to allow fast access to the information. The total construction time of the data structure is $\mathcal{O}((B + k^2) \cdot \log n)$. Note that the data structure will only be used $n/B - 1 \geq 2$ blocks later, so we can de-amortise the construction time over the succeeding block executing $\Omega((1 + k^2/B) \cdot \log n)$ steps of the construction process per symbol. The data structure occupies $\mathcal{O}(k^4 \log^2 n)$ bits of space.

4.2 Hamming distance

To compute the Hamming distance at an alignment i , the algorithm divides the pattern into three parts: a prefix of length $(B - i)$, a middle part of length $(n - B)$, and a suffix of length i , where the middle part is aligned with a block border. The algorithm then starts by computing the square N of the norm of the difference between the super-sketches of the middle part and the corresponding text substring. Both super-sketches are already known as the middle part is an $(n - B)$ -length substring of the pattern and we store its super-sketch explicitly, while the super-sketch of the text substring was computed at the end of the preceding block. As both sketches have length $9k^2$, it takes $\mathcal{O}(9k^2)$ time. Next, the algorithm computes the Hamming distance H_s between the sketch of the suffix of the pattern and the part of the text block seen so far. This again takes $\mathcal{O}(9k^2)$ time. Finally, the algorithm computes an approximation H_p of the Hamming distance between the prefix and the text as described in Problem 1. With the help of the binary search trees, $i_j, b_{\ell+1}$ and b_ℓ can be found in $\mathcal{O}(\log \log n + \log \log k^2)$ time. Recall that b_ℓ divides the prefix into two parts. The Hamming distance between the second part of the prefix and the text can be approximated in $\mathcal{O}(9k^2)$ time with the help of the sketches as in Problem 1, but it is not possible to use symbol-by-symbol comparison for the first part as this would take too much time. Instead, the algorithm does binary search on the prefixes' lengths it calculated during the preprocessing step which allows him to find $(1 + \varepsilon)$ -approximation of the Hamming distance in $\mathcal{O}(\log \log_{1+\varepsilon} n)$ time. It then outputs $H_p + H_m + H_s$, where $H_m = \varepsilon^2 N / 9(1 - \varepsilon/3)$.

4.3 Analysis

The running time of the algorithm is $\mathcal{O}(\varepsilon^{-2} \log n)$ per arriving symbol. The space used is $\mathcal{O}(\varepsilon^{-3} \sqrt{n} \log^2 n)$ bits. We now need to show that $H_p + H_m + H_s$ is a $(1 + \varepsilon)$ -approximation of the Hamming distance with constant probability. It suffices to show that H_m is a $(1 + \varepsilon)$ -approximation of the Hamming distance between the middle part of the pattern and the text. Consider two binary strings t and p of length $(n - B)$. Let sk_t and sk_p be their super-sketches of length $9k^2$ calculated with the help of M and σ_i and let $N = \|sk_t - sk_p\|_2^2$ and $\tilde{H} = \varepsilon^2 N$, where $\varepsilon = \varepsilon/3$. We will show that \tilde{H} is a good approximation of the Hamming distance between t and p . Recall that t and p are binary, and therefore the Hamming distance between them is equal to $\|t - p\|_2^2$.

► **Lemma 10.** *With constant probability $(1 - \tilde{\varepsilon}) \cdot \|t - p\|_2^2 \leq \tilde{H} \leq (1 + \tilde{\varepsilon}) \cdot \|t - p\|_2^2$.*

Proof. Let t_i and p_i , $i \in [1, n/B - 1]$, be the blocks of t and p of length B . We have

$$\mathbb{E}[\tilde{H}] = \tilde{\varepsilon}^2 \cdot \mathbb{E} \left[\left\| \sum_i \sigma_i M \cdot (t_i - p_i) \right\|_2^2 \right] = \tilde{\varepsilon}^2 \sum_j \mathbb{E} \left[\left(\sum_i \sigma_i M_j \cdot (t_i - p_i) \right)^2 \right]$$

where M_j is the j^{th} row of M . As all rows of M are identically distributed, we have $\mathbb{E} \left[\left(\sum_i \sigma_i M_j \cdot (t_i - p_i) \right)^2 \right] = \mathbb{E} \left[\left(\sum_i \sigma_i M_1 \cdot (t_i - p_i) \right)^2 \right]$ for all j , which is equal to $\|t - p\|_2^2$ as if at least one of the inequalities $i_1 = i_2$ or $j_1 = j_2$ does not hold, then the variables $\sigma_{i_1} M_1[j_1]$ and $\sigma_{i_2} M_1[j_2]$ are independent and the expectation of $\sigma_{i_1} \sigma_{i_2} M_1[j_1] M_1[j_2]$ is equal to zero, and otherwise it is equal to one. So finally we have $\mathbb{E}[\tilde{H}] = \|t - p\|_2^2$.

We now compute the variance of H . We again use the fact that the rows of M are independent and identically distributed.

$$\text{Var}[\tilde{H}] = \tilde{\varepsilon}^2 \cdot \text{Var} \left[\left(\sum_i \sigma_i M_1 \cdot (t_i - p_i) \right)^2 \right] \leq \tilde{\varepsilon}^2 \cdot \mathbb{E} \left[\left(\sum_i \sigma_i M_1 \cdot (t_i - p_i) \right)^4 \right].$$

By Khintchine's inequality there exists a universal constant $c > 0$ such that

$$\text{Var}[\tilde{H}] \leq c \tilde{\varepsilon}^2 \cdot \mathbb{E} \left[\left(\sum_i \sigma_i M_1 \cdot (t_i - p_i) \right)^2 \right]^2 \leq c \tilde{\varepsilon}^2 \cdot \|t - p\|_2^4.$$

The claim then follows by Chebyshev's inequality. \blacktriangleleft

Let now $H = \varepsilon^2 N / 9(1 - \varepsilon/3) = \tilde{H} / (1 - \varepsilon/3)$. The probability H is a $(1 + \varepsilon)$ -approximation of the Hamming distance between t and p is at least the probability \tilde{H} is in $[(1 - \varepsilon/3) \cdot \|t - p\|_2^2, (1 - \varepsilon/3)(1 + \varepsilon) \cdot \|t - p\|_2^2]$, which in turn can be estimated from below as

$$\Pr \left[\tilde{H} \in [(1 - \varepsilon/3) \cdot \|t - p\|_2^2, (1 + \varepsilon/3) \cdot \|t - p\|_2^2] \right] \geq 1 - 1/c \text{ (Lemma 10.)}$$

To justify the last transition note that $(1 - \varepsilon/3)(1 + \varepsilon) \geq (1 + \varepsilon/3)$ for all $\varepsilon < 1$. From above it follows that with constant probabilities H_p , H_m , and H_s are $(1 + \varepsilon)$ -approximations of the Hamming distances for the prefix, the middle part, and the suffix of the pattern respectively. We note that the probabilities can be made arbitrarily small by Chebyshev's inequality if we run a constant number of independent instances of the algorithm in parallel and output the sum of the medians of the values H_p , H_m , H_s . Correctness of the algorithm follows by the union bound.

Acknowledgements. We are grateful to T.S. Jayram and Paul Beame for helpful and inspiring conversations about the problems in this paper and to Ely Porat for introducing the original streaming pattern matching problem to us and for explaining how to solve Problem 3 in $\mathcal{O}(n^{2/3} \text{poly}(1/\varepsilon))$ space. We were also informed that Ely Porat had independently developed a solution that uses $\mathcal{O}(\sqrt{n}/\varepsilon^2)$ space and for each alignment with Hamming distance H outputs some integer in the interval $[(1 - \varepsilon) \cdot H - 1/2\sqrt{n}, (1 + \varepsilon) \cdot H + 1/2\sqrt{n}]$ [24].

References

- 1 Karl Abrahamson. Generalized string matching. *SIAM Journal on Computing*, 16(6):1039–1051, 1987.

- 2 Noga Alon, Yossi Matias, and Mario Szegedy. The space complexity of approximating the frequency moments. In *STOC'00: Proc. 28th Annual ACM Symp. Theory of Computing*, pages 20–29. ACM, 1996.
- 3 Amihood Amir, Yonatan Aumann, Gary Benson, Avivit Levy, Ohad Lipsky, Ely Porat, Steven Skiena, and Uzi Vishne. Pattern matching with address errors: Rearrangement distances. *Journal of Computer System Sciences*, 75(6):359–370, 2009.
- 4 Amihood Amir, Yonatan Aumann, Oren Kapah, Avivit Levy, and Ely Porat. Approximate string matching with address bit errors. In *CPM'08: Proc. 19th Annual Symp. on Combinatorial Pattern Matching*, pages 118–129, 2008.
- 5 Amihood Amir, Yonatan Aumann, Moshe Lewenstein, and Ely Porat. Function matching. *SIAM Journal on Computing*, 35(5):1007–1022, 2006.
- 6 Amihood Amir, Richard Cole, Ramesh Hariharan, Moshe Lewenstein, and Ely Porat. Overlap matching. *Information and Computation*, 181(1):57–74, 2003.
- 7 Amihood Amir, Estrella Eisenberg, and Ely Porat. Swap and mismatch edit distance. *Algorithmica*, 45(1):109–120, 2006.
- 8 Amihood Amir, Martin Farach, and S. Muthu Muthukrishnan. Alphabet dependence in parameterized matching. *Information Processing Letters*, 49(3):111–115, 1994.
- 9 Amit Chakrabarti and Oded Regev. An optimal lower bound on the communication complexity of gap-hamming-distance. *SIAM Journal on Computing*, 41(5):1299–1317, 2012.
- 10 Raphaël Clifford, Allyx Fontaine, Ely Porat, Benjamin Sach, and Tatiana A. Starikovskaya. The k -mismatch problem revisited. In *SODA'16: Proc. 27th ACM-SIAM Symp. on Discrete Algorithms*, pages 2039–2052, 2016.
- 11 Raphaël Clifford and Benjamin Sach. Pseudo-realtime pattern matching: Closing the gap. In *CPM'10: Proc. 21st Annual Symp. on Combinatorial Pattern Matching*, pages 101–111, 2010.
- 12 Graham Cormode, Mayur Datar, Piotr Indyk, and S. Muthukrishnan. Comparing data streams using Hamming norms (how to zero in). *IEEE Trans. on Knowl. and Data Eng.*, 15(3):529–540, 2003.
- 13 Mayur Datar, Aristides Gionis, Piotr Indyk, and Rajeev Motwani. Maintaining stream statistics over sliding windows. *SIAM Journal on Computing*, 31(6):1794–1813, 2002.
- 14 M. Fischer and M. Paterson. String matching and other products. In *Proc. 7th SIAM-AMS Complexity of Comp.*, pages 113–125, 1974.
- 15 Wei Huang, Yaoyun Shi, Shengyu Zhang, and Yufan Zhu. The communication complexity of the Hamming distance problem. *Information Processing Letters*, 99(4):149–153, 2006.
- 16 Piotr Indyk. Stable distributions, pseudorandom generators, embeddings, and data stream computation. *Journal of the ACM*, 53(3):307–323, 2006.
- 17 Markus Jalsenius, Benny Porat, and Benjamin Sach. Parameterized matching in the streaming model. In *STACS'13: Proc. 30th Annual Symp. on Theoretical Aspects of Computer Science*, pages 400–411, 2013. [arXiv:1109.5269](https://arxiv.org/abs/1109.5269).
- 18 Thathachar S. Jayram and David P. Woodruff. Optimal bounds for Johnson-Lindenstrauss transforms and streaming problems with subconstant error. *ACM Transactions on Algorithms (TALG)*, 9(3):26, 2013.
- 19 William Johnson and Joram Lindenstrauss. Extensions of Lipschitz mappings into a Hilbert space. In *Proc. of the Conference in Modern Analysis and Probability*, volume 26 of *Contemporary Mathematics*, pages 189–206. American Mathematical Society, 1984.
- 20 H. Karloff. Fast algorithms for approximately counting mismatches. *Information Processing Letters*, 48(2):53–60, 1993.
- 21 Tsvi Kopelowitz and Ely Porat. Breaking the variance: Approximating the Hamming distance in $1/\epsilon$ time per alignment. In *FOCS'15: Proc. 56th Annual Symp. Foundations of Computer Science*, pages 601–613, 2015.

20:14 Approximate Hamming Distance in a Stream

- 22 S. Rao Kosaraju. Efficient string matching. Manuscript, 1987.
- 23 Gad M. Landau and Uzi Vishkin. Fast string matching with k differences. *Journal of Computer System Sciences*, 37(1):63–78, 1988.
- 24 Ely Porat. Personal communication, 2016.