

PATTERN MATCHING UNDER POLYNOMIAL TRANSFORMATION*

AYELET BUTMAN[†], PETER CLIFFORD[‡], RAPHAËL CLIFFORD[§], MARKUS JALSENIOUS[§],
NOA LEWENSTEIN[¶], BENNY PORAT^{||}, ELY PORAT^{||}, AND BENJAMIN SACH^{**}

Abstract. We consider a class of pattern matching problems where a normalizing polynomial transformation can be applied at every alignment of the pattern and text. Normalized pattern matching plays a key role in fields as diverse as image processing and musical information processing, where application specific transformations are often applied to the input. By considering a wide range of such transformations, we provide fast algorithms and the first lower bounds for both new and old problems. Given a pattern of length m and a longer text of length n , where both are assumed to contain integer values only, we first show $O(n \log m)$ time algorithms for pattern matching under linear transformations even when wildcard symbols can occur in the input. We then show how to extend the technique to polynomial transformations of arbitrary degree. Next we consider the problem of finding the minimum Hamming distance under polynomial transformation. We show that, for any $\varepsilon > 0$, there cannot exist an $O(nm^{1-\varepsilon})$ time algorithm for additive and linear transformations conditional on the hardness of the classic 3SUM problem. Finally, we consider a version of the Hamming distance problem under additive transformations with a bound k on the maximum distance that needs to be reported. We give a deterministic $O(nk \log k)$ time solution, which we then improve by careful use of randomization to $O(n\sqrt{k} \log k \log n)$ time for sufficiently small k . Our randomized solution outputs the correct answer at every position with high probability.

Key words. string algorithms, pattern matching, normalization, 3SUM-hardness

AMS subject classifications. 68W32, 68W20

DOI. 10.1137/110853327

1. Introduction. We consider pattern matching problems, where the task is to find the distance between a pattern and every substring of the text of suitable length. In the class of problems we consider, the values in the pattern can first be transformed so as to minimize this distance. Further, the selection of which transformation to apply, which is possibly distinct for each alignment, forms part of the problem that is to be solved. This class of problems generalizes the well-known problem of exact matching with wildcards [8, 11] as well as the set of problems known previously as transposition invariant matching [24], both of which come from the pattern match-

*Received by the editors October 28, 2011; accepted for publication (in revised form) January 15, 2013; published electronically April 3, 2013. The results in section 2, except for those relating to higher degree polynomials, appeared in preliminary form in “Self-normalised Distance with Don’t Cares,” CPM ’07. Section 4 contains revised and rewritten versions of results from “Jump-Matching with Errors,” SPIRE ’07. The research of the fourth and eighth authors was supported by the EPSRC.

<http://www.siam.org/journals/sicomp/42-2/85332.html>

[†]Department of Computer Science, Holon Institute of Technology, Holon, Israel (ayeletb@hit.ac.il).

[‡]Department of Statistics, University of Oxford, Oxford OX1 3TG, UK (peter.clifford@stats.ox.ac.uk).

[§]Department of Computer Science, University of Bristol, Bristol BS8 1UB, UK (clifford@cs.bris.ac.uk, markus.jalsenius@gmail.com).

[¶]Department of Computer Science, Netanya Academic College, Netanya 42365, Israel (noa.lewenstein@gmail.com).

^{||}Department of Computer Science, Bar-Ilan University, Ramat Gan 52900, Israel (bennyporat@gmail.com, porately@macs.biu.ac.il).

^{**}Department of Computer Science, University of Warwick, Coventry CV4 7AL, UK (sach@dcs.warwick.ac.uk).

ing literature. However, as we will see, it is considerably broader than both, with applications in both image processing and musical information retrieval.

By way of a first motivation for our work, consider a fundamental problem in image processing which is to measure the similarity between a small image segment or template and regions of comparable size within a larger scene. It is well known that the cross-correlation between the two can be computed efficiently at every position in the larger image using the fast Fourier transform (FFT). In practice, images may differ in a number of ways including being rotated, or scaled, or affected by noise. We consider here the case where the intensity or brightness of an image occurrence is unknown and where parts of either image contain *don't care* or *wildcard* pixels, i.e., pixels that are considered to be irrelevant as far as image similarity is concerned. As an example, a rectangular image segment may contain a facial image, and the objective is to identify the face in a larger scene. However, some faces in the larger scene are in shadow and others are in light. Furthermore, background pixels around the faces may be considered to be irrelevant for facial recognition and these should not affect the search algorithm.

In order to overcome the first difficulty of varying intensity within an image, a standard approach is to compute the *normalized* distance when comparing a template to part of a larger image. Thus both template and image are transformed or rescaled in order to make any matches found more meaningful and to allow comparisons between matches at different positions. Within the image processing literature the accepted method of normalization is to scale the mean and variance of the template and image segments. We take a slightly different, although related, approach to normalization which will allow us to show a number of natural generalizations.

We start by defining measures of distance between a pattern P and text T , where P is a string of length m and T is a string of length $n \geq m$, both over the integers. The squared L_2 or Euclidean distance between the pattern and the text at position i is

$$\sum_{j=0}^{m-1} (P[j] - T[i+j])^2.$$

In this case, for each $i \in \{0, \dots, n-m\}$, the pattern can be normalized, or fitted as closely as possible to the text, by transforming the input to minimize the distance.

In the case of degree one polynomial transformations, the normalized L_2 distance between the pattern and the text at position i can now be written as

$$\min_{\alpha, \beta} \sum_{j=0}^{m-1} (\alpha + \beta P[j] - T[i+j])^2,$$

where the minimization is over rational values of α and β . The minimization is per alignment of P and T , and hence the values of α and β may (and probably will) differ between the positions i .

We also consider the case when the input alphabet is augmented with the special wildcard symbol, denoted \star . A position where either the pattern or text has a wildcard will not contribute to the distance. That is, the minimization is carried out using the sum of the remaining terms. Details are given in the problem definitions in the next section.

1.1. Problems and our results. The words *shift* and *scale* are used to refer to additive and multiplicative transformations of the pattern, respectively. The input to

all our problems is a text T of length n and a pattern P of length m , and the output is a problem specific distance $d(i)$ between P and T at every position $i \in \{0, \dots, n-m\}$ of the text. To avoid overloading variable names, we give the distance $d(i)$ a unique name for each problem.

PROBLEM 1.1 (SHIFT- L_2^*). *Normalized L_2 distance under shifts. Wildcards are allowed. We define*

$$d_2^+(i) \stackrel{\text{def}}{=} \min_{\alpha} \sum_{j=0}^{m-1} (\alpha + P[j] - T[i+j])^2.$$

When either $P[j] = \star$ or $T[i+j] = \star$, the contribution of the pair to the sum $d_2^+(i)$ is taken to be zero. The minimization is carried out using the sum of the remaining terms.

Next we define the normalized L_2 distance under shifts and scaling, corresponding to a degree one polynomial transformation of the values of the pattern.

PROBLEM 1.2 (SHIFTS- L_2^*). *Normalized L_2 distance under shifts and scaling. Wildcards are allowed. We define*

$$d_2^1(i) \stackrel{\text{def}}{=} \min_{\alpha, \beta} \sum_{j=0}^{m-1} (\alpha + \beta P[j] - T[i+j])^2.$$

When either $P[j] = \star$ or $T[i+j] = \star$, the contribution of the pair to the sum $d_2^1(i)$ is taken to be zero. The minimization is carried out using the sum of the remaining terms.

We show that both SHIFT- L_2^* and SHIFTS- L_2^* can be solved in $O(n \log m)$ time by the use of FFTs of integer vectors. Our results are stated in Theorems 2.1 and 2.2. We assume the RAM model of computation throughout in order to be consistent with previous work on matching with wildcards. Further, our techniques also provide $O(n \log m)$ time solutions (Theorems 2.3 and 2.4) to the problems of exact shift matching with wildcards (SHIFT-EXACT *) and exact shift-scale matching with wildcards (SHIFTS-EXACT *), formally defined as follows.

PROBLEM 1.3 (SHIFT-EXACT *). *Normalized exact matching under shifts. Wildcards are allowed. We define*

$$d_E^+(i) \stackrel{\text{def}}{=} \begin{cases} 1, & \exists \alpha \text{ such that } \alpha + P[j] = T[i+j] \text{ for all } j \in \{0, \dots, m-1\}; \\ 0 & \text{otherwise.} \end{cases}$$

Every position j where either $P[j] = \star$ or $T[i+j] = \star$ is ignored.

PROBLEM 1.4 (SHIFTS-EXACT *). *Normalized exact matching under shifts and scaling. Wildcards are allowed. The problem is defined similarly to SHIFT-EXACT * , but now we check whether there exist α and β such that $\alpha + \beta P[j] = T[i+j]$ for all positions j (except positions where $P[j] = \star$ or $T[i+j] = \star$).*

We will also discuss extensions to pattern transformations under polynomials of higher degree in section 2. In terms of normalized L_2 distance, we give the following definition.

PROBLEM 1.5 (POLY- r - L_2^*). *Normalized L_2 distance under degree- r polynomial transformation. Wildcards are allowed. Let $f(x) = \alpha_0 + \alpha_1 x + \alpha_2 x^2 + \dots + \alpha_r x^r$ be a polynomial of degree r with $r \geq 1$. We define*

$$d_2^r(i) \stackrel{\text{def}}{=} \min_{\alpha_0, \dots, \alpha_r} \sum_{j=0}^{m-1} (f(P[j]) - T[i+j])^2.$$

When either $P[j] = \star$ or $T[i+j] = \star$, the contribution of the pair to the sum $d_2^F(i)$ is taken to be zero. The minimization is carried out using the sum of the remaining terms.

Note that the problem $\text{SHIFTSCALE-}L_2^\star$ is the same problem as $\text{POLY-}r\text{-}L_2^\star$ with degree $r = 1$. We will show that $\text{POLY-}r\text{-}L_2^\star$ can be solved in $O(rn \log m + r^w n)$ time, where w is the exponent for matrix multiplication (e.g., $w \approx 2.38$ when using the Coppersmith–Winograd algorithm).

The second main topic of our work is on normalized pattern matching problems under the Hamming distance. The Hamming distance is perhaps the most commonly considered measure of distance between strings in the field of pattern matching. We therefore define related normalized versions of our pattern matching problems in a way similar to how we defined our previous problems.

PROBLEM 1.6 (SHIFT-HAM). *Normalized Hamming distance under shifts. Wildcards are not allowed. We define*

$$d_H^+(i) \stackrel{\text{def}}{=} \min_{\alpha} |\{j \mid \alpha + P[j] \neq T[i+j]\}|.$$

PROBLEM 1.7 (SHIFTSCALE-HAM). *Normalized Hamming distance under shifts and scaling. Wildcards are not allowed. We define*

$$d_H^1(i) \stackrel{\text{def}}{=} \min_{\alpha, \beta} |\{j \mid \alpha + \beta P[j] \neq T[i+j]\}|.$$

Previously, it was shown that SHIFT-HAM , sometimes also referred to as transposition invariant matching, can be solved in $O(nm \log m)$ time [24]. It has been tempting to believe that it might be possible to improve this time complexity, particularly as there exist algorithms for standard nonnormalized pattern matching under the Hamming distance which take $O(n\sqrt{m \log m})$ time [1, 19]. We show by reductions from the well-known 3SUM problem that for both shift and shift-scale matching under the Hamming distance there cannot exist an $O(nm^{1-\varepsilon})$ time algorithm for any $\varepsilon > 0$ (Theorems 3.4 and 3.6).

To circumvent this new conjectured lower bound, we consider as our last problem a shift version of the k -mismatch problem. In the k -mismatch problem, the Hamming distance is to be reported at every alignment as long as it is at most k . If it is greater than k , then the algorithm is only required to report that the Hamming distance is large. We define the problem as follows.

PROBLEM 1.8 (SHIFT- k -MISMATCH). *Normalized k -mismatch under shifts. Wildcards are not allowed. We define*

$$d_M^+(i) \stackrel{\text{def}}{=} \min(d_H^+(i), k+1).$$

We first give a simple deterministic $O(nk \log k)$ time solution (Theorem 4.4). We then consider a decision version of the problem where we output only the locations i for which $d_H^+(i) \leq k$. The decision version is defined as follows.

PROBLEM 1.9 (SHIFT- k -DECISION). *Normalized k -mismatch decision problem under shifts. Wildcards are not allowed. We define*

$$d_D^+(i) \stackrel{\text{def}}{=} \begin{cases} 0, & d_H^+(i) \leq k, \\ 1 & \text{otherwise.} \end{cases}$$

Using randomization we show how to solve this problem in $O(cn\sqrt{k \log k} \log n)$ time for the case that $k < \sqrt{m/6}$ (Theorem 4.13). Here c is a constant that can be

chosen arbitrarily to fine tune the error probability. Namely, our algorithm outputs the correct answer at every alignment with probability at least $1 - 1/n^c$. We therefore succeed in breaking our newly introduced running time barrier provided by the reduction from 3SUM for a limited range of values of k .

1.2. Related work. Combinatorial pattern matching has concerned itself mainly with strings of symbolic characters where the distance between individual characters is specified by some convention. For the k -mismatch problem, an $O(nk)$ time algorithm was given in 1986 that uses constant time lowest common ancestor queries on the suffix tree of the pattern and text in a technique that has subsequently come to be known as “kangaroo hopping” [20]. Almost 20 years afterward, the asymptotic running time was finally improved in [4] to $O(n\sqrt{k\log k})$ time by a method based on filtering, the suffix tree (with kangaroo hopping), and FFTs. In 2002, a deterministic $O(n\log m)$ time solution for exact matching with wildcards was given by Cole and Hariharan [11] and further simplified in [8]. In the same paper by Cole and Hariharan, an $O(n\log(\max(m, N)))$ time algorithm for the exact shift matching problem we consider in section 2 was presented. Here N is the largest value in the input. The approach we take to provide a simpler solution for this problem is similar in spirit to that of [8].

There has also been some work in recent years on fast algorithms for distance calculation and approximate matching between numerical strings. A number of different metrics has been considered, with, for example, $O(n\sqrt{m\log m})$ time solutions found for the L_1 distance [5, 6, 9] and less-than matching problems [3] and an $O(\delta n\log m)$ time algorithm for the δ -bounded version of the L_∞ norm first discussed in [10] and then improved in [9, 23].

The most closely related work to ours comes under the heading of transposition invariant matching [21]. The original motivation for this problem was within musical information retrieval, where a musical search is to be performed invariant of pitch level transposition. The transposition invariant distance between two equal length strings A and B is defined to be $\min_\alpha d(A + \alpha, B)$, where $A + \alpha$ is the string obtained from A by adding α to every value, and the distance d between strings can be variously defined. Algorithms for transposition invariant Hamming distance, longest common subsequence (LCS), and Levenshtein (edit) distance, among others, were given in [24], whose time complexities are close to the known upper bounds without transposition. We show, in section 3, lower bounds for the special case of transposition invariant Hamming distance, which we named SHIFT-HAM. Normalized pattern matching is also of central interest in the image processing literature, where normalization is typically performed by scaling the mean and standard deviation of the template and each suitably sized image segment to be 0 and 1, respectively. An asymptotically fast method for performing normalized cross-correlation for template matching, also using FFTs, was given in [22]. The methods we give in section 2 are similar to the approach in [22] only in the use of FFTs to provide fast solutions. Due to the differences in the definition of normalization between our work and [22], the solutions we give are otherwise quite distinct.

As a general class of problems, pattern matching under polynomial transformation is, to the best of our knowledge, new. However, if we allow the degree of the polynomial transformation to increase to m , then determining for which alignment the normalized distance equals zero is equivalent to the known problem of function matching. Function matching has a deterministic $O(n|\Sigma_P|\log m)$ time solution, where $|\Sigma_P|$ is the size of the pattern alphabet, and a faster randomized algorithm, which runs in $O(n\log n)$ time and has failure probability $1/n$ [2].

1.3. Basic notation. For a string X of length ℓ , we write $X[i]$ to denote the i th character of X such that $X = X[0] X[1] X[2] \cdots X[\ell - 1]$ (the first index is always zero). The s -length substring of X starting at position i is denoted $X[i \dots i + s - 1]$. For two strings X and Y , the notion $X \parallel Y$ is used to denote the string formed by concatenating X and Y in that order. All strings in this paper are over the integer alphabet. Therefore, $X[j]Y[j]$ denotes the product of the numerical characters $X[j]$ and $Y[j]$. If strings X and Y are of equal length, we use the notation $X \cdot Y$ for the string with characters $(X \cdot Y)[i] = X[i]Y[i]$. This elementwise arithmetic is used similarly for addition, subtraction, division, and power. For example, the i th symbol of X^2/Y is $X[i]^2/Y[i]$. For a real value k , the scalar multiplication kX is the string $(kX)[i] = kX[i]$.

The notation $\text{Ham}(X, Y)$ will be used to denote the Hamming distance between equal length strings X and Y :

$$\text{Ham}(X, Y) \stackrel{\text{def}}{=} |\{i \mid X[i] \neq Y[i]\}|.$$

Throughout this paper we use T to denote the text and P for the pattern. We use n to denote the length of T and m for the length of P .

Our algorithms in section 2 make extensive use of FFTs. An important property of the FFT is that the cross-correlation, defined as

$$(T \otimes P)[i] \stackrel{\text{def}}{=} \sum_{j=0}^{m-1} P[j]T[i+j],$$

can be calculated accurately and efficiently for all $i \in \{0, \dots, n - m\}$ in $O(n \log m)$ time (see, e.g., [13, Chapter 32]). The time complexity is reduced from $O(n \log n)$ to $O(n \log m)$ using a standard splitting trick which partitions the text into $2m$ length substrings which overlap each other by m characters. When it is clear from the context we use \sum as an abbreviation for $\sum_{j=0}^{m-1}$.

We use \star for the single character wildcard symbol. Under arithmetics on strings, as defined above, we may think of a wildcard as having the value zero. This value is, however, inconsequential for our purposes, as all expressions in this paper have the property that whenever a wildcard symbol is involved in some arithmetic, it is multiplied by a zero.

We write $[n]$ to denote the set of integers $\{0 \dots n - 1\}$. We also say that $g(n) \in \tilde{\Omega}(h(n))$ if and only if $g(n) \in \Omega(h(n)/\log^c n)$ for some constant c , i.e., $g(n) \in \Omega(h(n))$ up to log factors.

1.4. Organization. The remainder of the paper is organized as follows. In section 2 we discuss normalized pattern distance under L_2 distance (SHIFT- L_2^\star and SHIFTS-SCALE- L_2^\star) and the decision variants (SHIFT-EXACT * and SHIFTS-SCALE-EXACT *). We also show how to extend the methods to transformations of higher degree polynomials (POLY- r - L_2^\star). Then in section 3 we give running time lower bounds for SHIFT-HAM and SHIFTS-SCALE-HAM by reduction from the 3SUM problem. In section 4 we introduce our new deterministic and randomized algorithms for SHIFT- k -MISMATCH and SHIFT- k -DECISION. Finally, we conclude in section 5 and set out some open problems.

2. Normalized L_2 distance. We give $O(n \log m)$ time solutions for shift and shift-scale versions of the normalized L_2 distance problem with wildcards. We further show that this enables us to solve the exact shift matching and exact shift-scale

ALGORITHM 1. Solution to SHIFT- L_2^* .

-
1. Construct P' from P such that $P'[j] = 0$ if $P[j] = \star$, and $P'[j] = 1$ otherwise. Construct T' from T similarly.
 2. Compute the following six cross-correlations:

$$\begin{aligned} C_1 &= (T^2 \cdot T') \otimes P' & C_3 &= T' \otimes (P^2 \cdot P') & C_5 &= T' \otimes (P \cdot P') \\ C_2 &= (T \cdot T') \otimes (P \cdot P') & C_4 &= (T \cdot T') \otimes P' & C_6 &= T' \otimes P' \end{aligned}$$
 3. Return $A = C_1 - 2C_2 + C_3 - ((C_4 - C_5)^2 / C_6)$. We have $d_2^+(i) = A[i]$. For positions i where $C_6[i] = 0$ we have $d_2^+(i) = 0$.
-

matching problems in the same time complexity for inputs containing wildcard symbols. Lastly, we show how to extend our solutions to normalization under polynomials of arbitrary degree.

2.1. Normalized L_2 distance under shifts. In order to handle wildcards, we define two new strings P' and T' obtained from P and T , respectively, such that $P'[j] = 0$ if $P[j] = \star$, and $P'[j] = 1$ otherwise. Similarly, $T'[i] = 0$ if $T[i] = \star$, and $T'[i] = 1$ otherwise. We can now express the shift-normalized L_2 distance at position i as

$$d_2^+(i) = \min_{\alpha} \sum_{j=0}^{m-1} \left((\alpha + P[j] - T[i+j])^2 \cdot P'[j] \cdot T'[i+j] \right).$$

Algorithm 1 shows how to compute $d_2^+(i)$ for all positions i . Correctness and running time are given in the following theorem.

THEOREM 2.1. *The shift version of the normalized L_2 distance with wildcards problem (SHIFT- L_2^*) can be solved in $O(n \log m)$ time.*

Proof. Consider Algorithm 1. We first analyze the running time. Step 1 requires only single passes over the input. Similarly, $(P^2 \cdot P')$, $(P \cdot P')$, $(T \cdot T')$, and $(T^2 \cdot T')$ can all be calculated in linear time once T' and P' are known. Using the FFT, the six cross-correlations in step 2 can be calculated in $O(n \log m)$ time. The final vector of step 3 is obtained in linear time. Thus, $O(n \log m)$ is the overall time complexity of the algorithm.

To show correctness we consider the minimum value of

$$(2.1) \quad A[i] = \sum_{j=0}^{m-1} \left((\alpha + P[j] - T[i+j])^2 \cdot P'[j] \cdot T'[i+j] \right).$$

This can be obtained by differentiating with respect to α and obtaining the minimizing value. Solving

$$\frac{\partial A[i]}{\partial \alpha} = 2 \sum_{j=0}^{m-1} \left((\alpha + P[j] - T[i+j]) \cdot P'[j] \cdot T'[i+j] \right) = 0$$

gives us the value

$$\hat{\alpha} = \frac{\sum \left((T[i+j] - P[j]) \cdot P'[j] \cdot T'[i+j] \right)}{\sum P'[j] \cdot T'[i+j]} = \frac{((T \cdot T') \otimes P') - ((T' \otimes (P \cdot P')))}{T' \otimes P'},$$

ALGORITHM 2. Solution to $\text{SHIFTSCALE-}L_2^*$.

1. Construct P' from P such that $P'[j] = 0$ if $P[j] = \star$, and $P'[j] = 1$ otherwise. Construct T' from T similarly.
 2. Compute the following six cross-correlations:
 $C_1 = (T^2 \cdot T') \otimes P' \quad C_3 = T' \otimes (P^2 \cdot P') \quad C_5 = T' \otimes (P \cdot P')$
 $C_2 = (T \cdot T') \otimes (P \cdot P') \quad C_4 = (T \cdot T') \otimes P' \quad C_6 = T' \otimes P'$
 3. Compute
 $B_1 = C_3 \cdot C_4 - C_2 \cdot C_5, \quad B_2 = C_3 \cdot C_6 - C_5^2, \quad B_3 = C_2 \cdot C_6 - C_4 \cdot C_5$
and compute $\hat{\alpha} = B_1/B_2$ and $\hat{\beta} = B_3/B_2$. At positions i where $C_6[i] = 0$, set $\hat{\alpha}[i] = \hat{\beta}[i] = 0$. At positions i where $B_2[i] = 0$ and $C_6[i] \neq 0$, set $\hat{\alpha}[i] = C_4/C_6$ and $\hat{\beta}[i] = 0$.
 4. Return $B = (\hat{\alpha}^2 \cdot C_6) + 2(\hat{\alpha} \cdot \hat{\beta} \cdot C_5) - 2(\hat{\alpha} \cdot C_4) + (\hat{\beta}^2 \cdot C_3) - 2(\hat{\beta} \cdot C_2) + C_1$.
We have $d_2^1(i) = B[i]$.
-

where $\hat{\alpha}[i]$ is the minimizing value at position i . Substituting $\alpha = \hat{\alpha}$ into (2.1) and expanding and collecting terms, we obtain the final answer as

$$A = C_1 - 2C_2 + C_3 - \frac{(C_4 - C_5)^2}{C_6},$$

where C_1, \dots, C_6 are the correlations defined in Algorithm 1.

Lastly, we observe that when $C_6[i] = (T' \otimes P')[i] = 0$ there is a wildcard at every position in the alignment of P and T . Here the shift-normalized L_2 distance is defined to be 0. \square

2.2. Normalized L_2 distance under shift-scale. Similarly to the shift version of normalized L_2 distance in the previous section, we can now solve the shift-scale version. The solution is slightly more involved, but the running time remains the same. Algorithm 2 sets out the main steps to achieve this, and the result is summarized in the following theorem.

THEOREM 2.2. *The shift-scale version of the normalized L_2 distance with wildcards problem ($\text{SHIFTSCALE-}L_2^*$) can be solved in $O(n \log m)$ time.*

Proof. Consider Algorithm 2. Notice that the same six correlations as in Algorithm 1 have to be calculated. The additional strings in step 3 require linear time, as well as producing the output in step 4. Hence the overall running time is $O(n \log m)$.

Similarly to (2.1) we can express the shift-scale version of the normalized L_2 distance at position i as

$$(2.2) \quad B[i] = \sum_{j=0}^{m-1} \left((\alpha + \beta P[j] - T[i+j])^2 \cdot P'[j] \cdot T'[i+j] \right).$$

By minimizing this expression with respect to both α and β we get a system of two simultaneous linear equations,

$$\frac{\partial B[i]}{\partial \alpha} = 2 \sum_{j=0}^{m-1} \left((\alpha + \beta P[j] - T[i+j]) \cdot P'[j] \cdot T'[i+j] \right) = 0,$$

$$\frac{\partial B[i]}{\partial \beta} = 2 \sum_{j=0}^{m-1} \left((\alpha + \beta P[j] - T[i+j]) \cdot P[j] \cdot P'[j] \cdot T'[i+j] \right) = 0.$$

By solving this system and using the definitions of B_1, \dots, B_3 in Algorithm 2, we get the minimizing values

$$\hat{\alpha} = \frac{B_1}{B_2} \quad \text{and} \quad \hat{\beta} = \frac{B_2}{B_3}.$$

For some positions i , the solution to the system might not be unique. This happens at alignments i for which every position $i + j$ has a wildcard; hence $C_6[i] = 0$. Here we avoid illegal division by zero by simply setting both $\hat{\alpha}[i]$ and $\hat{\beta}[i]$ to zero (any value would do). A nonunique solution also occurs at alignments i where all $P[j]$ are identical over every nonwildcard position $i + j$. This is characterized by $B_2[i] = 0$. To see this, observe that $C_5[i]^2 \leq C_3[i]C_6[i]$ by the Cauchy–Schwarz inequality. Here we set (arbitrarily) $\hat{\beta}[i] = 0$ and therefore obtain the minimizing value $\hat{\alpha}[i] = C_4/C_6$.

At stage 4, $\hat{\alpha}$ and $\hat{\beta}$ contain the minimizing values for α and β at every position. We substitute these into (2.2) and expand. This gives us the expression for B . \square

2.3. Exact shift and shift-scale matching with wildcards. For the exact shift matching problem with wildcards, SHIFT-EXACT*, a match is said to occur at location i if, for some shift α and for every position j in the pattern, either $\alpha + P[j] = T[i + j]$ or at least one of $P[j]$ and $T[i + j]$ is the wildcard symbol. Cole and Hariharan [11] introduced a new coding for this problem that maps the string elements into 0 for wildcards and complex numbers of modulus 1 otherwise. The FFT is then used to find the (complex) cross-correlation between these coded strings, and finally a shift match is declared at location i if the i th element of the modulus of the cross-correlation is equal to $(P' \otimes T')[i]$.

Our Algorithm 1 provides a straightforward alternative method for shift matching with wildcards. It has the advantage of using only simple integer codings. Since Algorithm 1 finds the minimum L_2 distance at location i , over all possible shifts, it is necessary only to test whether this distance is zero. The running time for the test is then $O(n \log m)$ since it is determined by the running time of Algorithm 1.

THEOREM 2.3. *The problem of exact shift matching with wildcards (SHIFT-EXACT*) can be solved in $O(n \log m)$ time.*

The exact shift-scale matching problem with wildcards, SHIFTS-SCALE-EXACT*, can be solved similarly by applying Algorithm 2.

THEOREM 2.4. *The problem of exact shift-scale matching with wildcards (SHIFTS-SCALE-EXACT*) can be solved in $O(n \log m)$ time.*

2.4. Normalized L_2 distance under higher degree transformations. We can now consider the problem of computing the normalized L_2 distance under general polynomial transformations. The problem, which we termed POLY- r - L_2^* , was defined in Problem 1.5. Recall that we let

$$f(x) = \alpha_0 + \alpha_1 x + \alpha_2 x^2 + \dots + \alpha_r x^r$$

be a polynomial of degree $r \geq 1$. As in the shift and shift-scale versions of the normalized L_2 distance, we consider the minimum value of

$$(2.3) \quad D[i] = \sum_{j=0}^{m-1} \left((f(P[j]) - T[i + j])^2 \cdot P'[j] \cdot T'[i + j] \right).$$

By differentiating with respect to each α_k in turn, giving

$$\frac{\partial D[i]}{\partial \alpha_k} = 2 \sum_{j=0}^{m-1} \left((f(P[j]) - T[i+j]) \cdot P[j]^k \cdot P'[j] \cdot T'[i+j] \right) = 0,$$

we obtain a system of $r+1$ linear equations in $r+1$ unknowns for each alignment i of the pattern and text. We need to solve these equations and then substitute the minimizing α_k values back into (2.3) as we did in the proof of Theorem 2.2. This procedure is captured by the following theorem.

THEOREM 2.5. *The normalized L_2 distance problem with wildcards under polynomial transformations of degree r ($\text{POLY-}r\text{-}L_2^*$) can be solved in $O(rn \log m + r^{2.38}n)$ time.*

Proof. To compute the coefficients for the first linear equation for α_0 we need to perform $O(r)$ cross-correlations. However, for each subsequent equation for $\alpha_1 \dots \alpha_r$ we need only perform a constant number of new cross-correlations. Therefore the total number of cross-correlations is $O(r)$ to give the coefficients of all the equations, taking $O(rn \log m)$ time overall. The time to solve the systems of $O(r)$ equations in $O(r)$ unknowns is $O(r^w)$ per alignment i , where w is the exponent for matrix multiplication. This gives $O(nr^w)$ time or $O(nr^{2.38})$ using the algorithm of Coppersmith and Winograd [12].

Once the equations have been solved, and the minimizing values of α_k calculated, they are then substituted into the expression for D in (2.3). To calculate the final values $D[i]$ we require $O(r)$ cross-correlations to be computed as well as $O(r^2)$ products of vectors of length m . The overall time complexity is therefore $O(rn \log m + r^{2.38}n + r^2m)$. \square

This method is of particular relevance for low degree polynomials, or at least polynomials whose degree is less than the number of distinct values in the pattern. However, if the degree r is greater than the number of distinct values in the pattern, then there exists a suitable polynomial f for any mapping we should choose. This gives us a straightforward $O(nm)$ time solution by considering each position of the pattern in the text independently and ignoring any values aligned with wildcards in either the pattern or text. For each such position we need only set $f(P[j])$ to be the mean of the values in the text that align with a value equal to $P[j]$ in the pattern.

3. Lower bounds for Hamming distance. In this section we will show that no $O(nm^{1-\varepsilon})$ time algorithm can exist for either SHIFT-HAM or SHIFTSCALE-HAM conditional on the hardness of the classic 3SUM problem. One formulation of the 3SUM problem is given below.

DEFINITION 3.1 (3SUM). *Given a set of s positive integers, determine whether there are three elements a, b, c in the set such that $a + b = c$.*

The 3SUM problem can be solved in $O(s^2)$ time, and it is a long standing conjecture that this is essentially the best possible. The problem has been extensively discussed in the literature, where Gajentaan and Overmars [17] were the first to introduce the concept of 3SUM -hardness (see definition below) to show that a wide range of problems in computational geometry is at least as hard as the 3SUM problem. One example is the GEOMBASE problem, defined below, which we will use in one of our reductions in this section. See [18] for a survey of problems from computational geometry whose hardness relies on that of 3SUM .

DEFINITION 3.2 (GEOMBASE). *Given a set of s points with integer coordinates on three horizontal lines $y = 0$, $y = 1$, and $y = 2$, determine whether there exists a nonhorizontal line containing three of the points.*

Although an $\tilde{\Omega}(s^2)$ lower bound for 3SUM is only conjectured, it has been shown that under certain restricted models of computation, $\Omega(s^2)$ is a true lower bound (see [14, 15, 16]). Under models that allow more direct manipulation of numbers instead of just real arithmetic, such as the word-RAM model, an almost $\log^2 s$ factor improvement to the standard $O(s^2)$ solution has been shown to be possible under the Las Vegas model of randomization (see [7]). Nevertheless, a 3SUM-hardness result for a problem is a strong indication that finding an $O(s^{2-\varepsilon})$ time solution is going to be a challenging task.

Before we show that SHIFT-HAM and SHIFTSIZE-HAM are both 3SUM-hard, we provide a brief but formal discussion about reductions and define 3SUM-hardness.

3.1. 3SUM reductions. Following the definitions of [17] where 3SUM-hardness was first introduced, we say that a problem A is $g(s)$ -solvable using a problem B if and only if every instance of A of size s can be solved using a constant number of instances of B of at most $O(s)$ size and $O(g(s))$ additional time. We denote this as $A \lll_{g(s)} B$. When $g(s)$ is sufficiently small, lower bounds for A carry over to B . A problem B is 3SUM-hard if $3\text{SUM} \lll_{g(s)} B$ and $g(s) = o(s^{2-\varepsilon})$ for some constant $\varepsilon > 0$. In the definition of 3SUM-hardness of [17], the requirement was that $g(s) = o(s^2)$; however, to scale with more powerful models of computation, we require that $g(s) = o(s^{2-\varepsilon})$. If $A \lll_{g(s)} B$ and $B \lll_{g(s)} A$, then we say that A and B are $g(s)$ -equivalent.

In the following section we will show that $3\text{SUM} \lll_{s \log s} \text{SHIFT-HAM}$, where the instance size of SHIFT-HAM is a text of length $n = 5s$ and a pattern of length $m = 3s$.

In the literature there are a variety of definitions of the 3SUM problem. They differ only slightly in their formulations and are all equivalent. One common definition, used as the “base problem” in [17], is formulated as follows. Given a set of s integers, determine whether there are three elements a, b, c in the set such that $a + b + c = 0$. Without too much work, one can show that this definition is $O(s)$ -equivalent with Definition 3.1 of 3SUM above (small modifications of the proof of Theorem 3.1 in [17] can be used to prove this). Further, it was shown in [17] that GEOMBASE is $O(s)$ -equivalent to 3SUM.

3.2. 3SUM-hardness of SHIFT-HAM. In this section we show that SHIFT-HAM is 3SUM-hard.

LEMMA 3.3. $3\text{SUM} \lll_{s \log s} \text{SHIFT-HAM}$ where the instance size of SHIFT-HAM is a text of length $5s$ and a pattern of length $3s$.

Proof. Let the set S be an instance of 3SUM of size $s = |S|$. First, we sort all elements of S so that $S = \{x_1, \dots, x_s\}$, where $x_1 < x_2 < \dots < x_s$. Let $y_1 = 2x_s + 1$, and for $i \in \{2, \dots, 2s\}$, let $y_i = y_{i-1} + 1$. Thus, $x_s < y_1 < \dots < y_{2s}$. We define the following s -length strings over the alphabet $\{x_1, \dots, x_s\} \cup \{y_1, \dots, y_{2s}\} \cup \{0\}$:

$$\begin{aligned} S_0 &= 00 \cdots 0 \quad (s \text{ zeros}) & S_3 &= y_{s+1} y_{s+2} \cdots y_{2s} \\ S_1 &= x_1 x_2 \cdots x_s & S_4 &= x_s x_{s-1} \cdots x_1 \\ S_2 &= y_1 y_2 \cdots y_s \end{aligned}$$

We now construct an instance of SHIFT-HAM specified by

$$T = S_0 \| S_1 \| S_2 \| S_1 \| S_3 \quad \text{and} \quad P = S_4 \| S_0 \| S_0.$$

The text T has length $n = 5s$ and the pattern P has length $m = 3s$. First, we show that if there are elements $a, b, c \in S$ such that $a + b = c$, then there is a position i such that the shift-normalized Hamming distance between P and $T[i \dots i + m - 1]$ is

T: 0 0 0 0 0 0 0 x_1 x_2 x_3 x_4 x_5 $\boxed{x_6}$ x_7 y_1 y_2 y_3 y_4 y_5 y_6 y_7 x_1 x_2 $\boxed{x_3}$ x_4 x_5 x_6 x_7 y_8 y_9 y_{10} y_{11} y_{12} y_{13} y_{14}
P: x_7 x_6 x_5 $\boxed{x_4}$ x_3 x_2 x_1 0 0 0 0 0 0 0 $\boxed{0}$ 0 0 0 0 0 0

FIG. 3.1. Example of an alignment of P and T in the proof of Lemma 3.3.

at most $m - 2$. We will then show that if three such elements do not exist, then the shift-normalized Hamming distance between P and every m -length substring of T is strictly greater than $m - 2$.

As an illustrative example (see Figure 3.1), suppose that S contains seven elements and suppose that $x_4 + x_3 = x_6$. Consider the alignment of P and T , where x_4 in P is aligned with x_6 in T . We observe that shifting the pattern by x_3 will induce two matches, marked with the squares above. Thus, the shift-normalized Hamming distance is at most $m - 2$ (in fact, it is exactly $m - 2$). It should be easy to see how this generalizes to any size of S and any three elements $a, b, c \in S$ such that $a + b = c$. Namely, the alignment in which a is aligned with c has Hamming distance at most $m - 2$ since there must also be a match at the position where 0 aligned with b . The construction of P and T ensures that there is always an alignment that captures these matches.

Now suppose there are no elements $a, b, c \in S$ such that $a + b = c$. Consider a fixed alignment of P and T . We will show that there can be at most one match under any shift. By construction of P and T , the zeros in P are all aligned with distinct symbols in T . Hence for any shift, at most one of these zeros can be involved in a match. The nonzero symbols of P (i.e., the s -length prefix of P) appear in strictly decreasing order and are aligned with an s -length substring of T whose elements appear in nondecreasing order. Therefore, under any shift, at most one of the nonzero symbols in P can be involved in a match. It remains to show that there is no shift such that both a zero and a nonzero symbol in P are simultaneously involved in a match. First, we observe that if there is a match between a zero in P and some y_j in T , then there can be no other match, as every nonzero symbol in P is aligned with a value that is less than y_j . Suppose therefore that there is a match between a zero in P and some x_j in T (i.e., the shift is x_j). We need to consider three possible cases: there is also a match that involves some x_k in P aligned with either (i) a zero in T , (ii) some y_ℓ in T , or (iii) some x_ℓ in T . In case (i) the shift must be negative, and hence is not compatible with the shift x_j . In case (ii) we can see that the shift must be greater than x_s (the largest elements in the set S), and hence is not compatible with the shift x_j . In case (iii) we have that $x_k + x_j = x_\ell$, which contradicts the assumption that there are no elements $a, b, c \in S$ such that $a + b = c$. Thus, the shift-normalized Hamming distance is at least $m - 1$ for any alignment of P and T .

Finally, we observe that the most time consuming part of the reduction is the sorting of S which could take $O(s \log s)$ time. This concludes the proof. \square

THEOREM 3.4. SHIFT-HAM has no $O(nm^{1-\varepsilon})$ time algorithm, for any $\varepsilon > 0$, conditional on the hardness of the 3SUM problem.

Proof. Given a 3SUM instance of size s , by Lemma 3.3 we construct a SHIFT-HAM instance of size $n = 5s$ and $m = 3s$ in $O(s \log s)$ time. If SHIFT-HAM has an $O(nm^{1-\varepsilon})$ time algorithm, then 3SUM can be solved in $O(s^{2-\varepsilon})$ time. \square

Notice that SHIFT-HAM has an $O(nm \log m)$ time solution [24]. See section 4.1 for details.

3.3. 3SUM-hardness of SHIFTSIZE-HAM. In this section we show that SHIFTSIZE-HAM is 3SUM-hard.

LEMMA 3.5. $3\text{SUM} \lll_s \text{SHIFTSCALE-HAM}$, where the instance size of SHIFTSCALE-HAM is a pattern and text of length s each.

Proof. We reduce from the GEOMBASE problem which is $O(s)$ -equivalent to 3SUM . Before we describe the reduction we adopt a formulation of the GEOMBASE problem that differs slightly in notation. Instead of insisting on the points being on the horizontal lines $y = 0$, $y = 1$, and $y = 2$, we assume that the points are on the vertical lines $x = 0$, $x = 1$, and $x = 2$, and we want to determine whether there is a (nonvertical) line containing three points. Under this formulation, let S be an instance of GEOMBASE that contains the integer points $(x_1, y_1), (x_2, y_2), \dots, (x_s, y_s)$, where every $x_j \in \{0, 1, 2\}$.

We construct an instance of SHIFTSCALE-HAM that is specified by the text $T = y_1 y_2 \cdots y_s$ and the pattern $P = x_1 x_2 \cdots x_s$, both of length s . It should now be clear that SHIFTSCALE-HAM returns the shift-and-scale normalized Hamming distance $s-3$ (for the only alignment of P and T) if and only if there are two values α and β such that $\beta x_j + \alpha = y_j$ for three distinct positions j , which is equivalent to fitting a line through three points. Note that we minimize α and β over the rationals, and any line going through three points is indeed specified by rational values of α and β . Since the reduction takes linear time, we have proved the lemma. \square

THEOREM 3.6. SHIFTSCALE-HAM has no $O(nm^{1-\varepsilon})$ algorithm, for any $\varepsilon > 0$, conditional on the hardness of the 3SUM problem.

Proof. Given a 3SUM instance of size s , by Lemma 3.5 we construct a SHIFTSCALE-HAM instance of size $n = m = s$ in $O(s)$ time. If SHIFTSCALE-HAM has an $O(nm^{1-\varepsilon})$ algorithm, then 3SUM can be solved in $O(s^{2-\varepsilon})$ time. \square

4. Normalized k -mismatch under shifts. In this section we consider two versions of the normalized k -mismatch problem under shifts, defined as Problems 1.8 and 1.9 in the introduction. Both problems are parameterized by an integer k . In the first problem, $\text{SHIFT-}k\text{-MISMATCH}$, the output is the shift-normalized Hamming distance between P and T at every position for which the distance is k or less. Where the distance is larger than k , only $k+1$ is outputted. Recall from the introduction that the shift-normalized Hamming distance between P and $T[i \dots i + m - 1]$ is denoted $d_H^+(i)$ and defined by

$$d_H^+(i) \stackrel{\text{def}}{=} \min_{\alpha} |\{j \mid \alpha + P[j] \neq T[i+j]\}|.$$

In section 4.2 we give a deterministic algorithm that solves $\text{SHIFT-}k\text{-MISMATCH}$ in $O(nk \log k)$ time.

In section 4.3 we consider the second version of shift-normalized k -mismatch, $\text{SHIFT-}k\text{-DECISION}$, which unlike the previous problem indicates only by yes or no whether the shift-normalized Hamming distance is k or less. We give a randomized solution to this decision problem with the improved running time $O(cn\sqrt{k} \log k \log n)$. The parameter c is a constant that can be chosen arbitrarily to fine tune the error probability. Namely, the probability that our algorithm outputs the correct answer at every alignment is at least $(1 - 1/n^c)$. The errors are one-sided such the algorithm will never miss reporting an alignment for which the shift-normalized Hamming distance is indeed k or less. Our algorithm requires that $k < \sqrt{m/6}$, and hence it is suited to situations where the locations of text substrings similar to the pattern are required but the distances themselves are not needed.

4.1. The unbounded case. In [24], Mäkinen, Navarro, and Ukkonen gave an $O(nm \log m)$ time algorithm for the shift-normalized Hamming distance prob-

lem, SHIFT-HAM, which by definition solves the bounded, k -mismatch variant in $O(nm \log m)$ time also. We briefly recap their method by way of an introduction. First, observe that the maximum number of matches for any alignment is exactly

$$m - d_{\text{H}}^+(i) = \max_{\alpha} \{ j \mid T[i+j] - P[j] = \alpha \}.$$

For each alignment i , this value can be obtained by creating an m -length array A_i , which we refer to as the *shift array*, defined by

$$(4.1) \quad A_i[j] = T[i+j] - P[j]$$

for all $j \in [m]$. This shift array is then sorted to find the most frequent value, which is the α that minimizes $d_{\text{H}}^+(i)$. The number of times it occurs is $m - d_{\text{H}}^+(i)$. Computing this requires $O(m \log m)$ time per alignment and hence $O(nm \log m)$ time overall. In the next section we will reconsider A_i and demonstrate that it can be run-length encoded in $O(k)$ runs whenever $d_{\text{H}}^+ \leq k$.

4.2. A deterministic solution. The deterministic algorithm makes use of the notion of difference strings, which were introduced in [21] and are defined as follows.

DEFINITION 4.1. *Let S be a string of length s . The difference string of S , denoted S_{δ} , is defined by*

$$S_{\delta}[j] = S[j+1] - S[j]$$

for all $j \in [s-1]$. The length of S_{δ} is $s-1$.

We will also make use of a generalization of the difference string when we present our randomized algorithm in section 4.3. The core of our deterministic shift-normalized k -mismatch algorithm is the relationship between the number of mismatches between P_{δ} and $T_{\delta}[i \dots i+m-2]$ and the value of $d_{\text{H}}^+(i)$. We begin in Lemma 4.2 below by showing that if $d_{\text{H}}^+(i)$ is small, then the number of mismatches between the difference strings P_{δ} and T_{δ} is also small. In [24] a related result was used to reduce the shift-normalized exact matching problem to the conventional exact matching problem. Specifically, they observed that in the special case that $k=0$, the implication becomes an equivalence, i.e., $d_{\text{H}}^+(i) = 0$ if and only if $P_{\delta} = T_{\delta}[i \dots i+m-2]$. Unfortunately, this is not the case in general.

LEMMA 4.2. *Let P be a pattern and T a text. For all i ,*

$$d_{\text{H}}^+(i) \leq k \implies \text{Ham}(P_{\delta}, T_{\delta}[i \dots i+m-2]) \leq 2k.$$

Proof. Let i be such that $d_{\text{H}}^+(i) \leq k$, and therefore there exists an α such that for at most k distinct position $j \in [m]$ we have that $P[j] + \alpha \neq T[i+j]$. Further, at most $2k$ distinct positions $j \in [m-1]$ have either $P[j] + \alpha \neq T[i+j]$ or $P[j+1] + \alpha \neq T[i+j+1]$. This implies that there are at least $(m-1) - 2k$ distinct positions $j \in [m-1]$ such that $P[j] + \alpha = T[i+j]$ and $P[j+1] + \alpha = T[i+j+1]$. By rearranging these equations, for any such j we have that $P[j+1] - P[j] = T[i+j+1] - T[i+j]$, and hence by Definition 4.1, $P_{\delta}[j] = T_{\delta}[i+j]$. As required there are at most $2k$ mismatches between P_{δ} and $T_{\delta}[i \dots i+m-2]$ (recall that $|P_{\delta}| = m-1$). \square

Lemma 4.2 suggests the following strategy. First, we find the leftmost up to $2k+1$ mismatches between P_{δ} and $T_{\delta}[i \dots i+m-2]$ at each alignment i . By Lemma 4.2 we can disregard any alignments with more than $2k$ mismatches. Finally, we use the locations of these mismatches to infer $d_{\text{H}}^+(i)$ at the remaining alignments.

ALGORITHM 3. Overview of deterministic solution to SHIFT- k -MISMATCH.

1. Compute the difference strings P_δ and T_δ by scanning P and T .
2. Run a $2k$ -mismatch algorithm on P_δ and T_δ in order to find all alignments where the number of mismatches is at most $2k$. The $2k$ -mismatch algorithm must also return the locations of the mismatches at any alignment where there are at most $2k$ mismatches.
3. Discard all alignments with more than $2k$ mismatches.
4. For each undiscarded alignment i , decompose A_i into at most $4k + 1$ runs (substrings with a common value). The start and end points of the runs are determined by scanning the locations of the mismatches between P_δ and $T_\delta[i \dots i + m - 1]$.
5. Sort the runs in A_i by value in order to find the most frequent entry α in A_i . Then output $m - |\{j \mid A_i[j] = \alpha\}|$, which is the value $d_H^+(i)$.

The first step can be done using any k -mismatch (strictly $2k$ -mismatch) algorithm which returns the locations of the mismatches. The well-known “kangaroo” method of [20] achieves this in optimal $O(nk)$ time. The method is so named because it uses longest common extensions to “hop” between mismatches in constant time. The discarding phase is trivial and therefore we focus only on computing $d_H^+(i)$ from the locations of the (at most $2k$) mismatches between P_δ and $T_\delta[i \dots i + m - 2]$, where i is an arbitrary nondiscarded alignment.

Recall from section 4.1 the definition of the shift array A_i in (4.1), and recall that the value of $m - d_H^+(i)$ is the number of occurrences of the most frequent entry in A_i . We will now use the locations of the mismatches between P_δ and $T_\delta[i \dots i + m - 2]$ to obtain a run-length encoded version of A_i containing $O(k)$ runs. The key property we require is given in Lemma 4.3, which states that a matching substring in P_δ and $T_\delta[i \dots i + m - 2]$ corresponds to a run (a substring of equal values) in A_i . This immediately implies that A_i can be decomposed into at most $4k + 1$ runs, specifically, one run of length 1 for each mismatch and an additional run for each stretch between mismatches.

LEMMA 4.3. *If $P_\delta[\ell \dots r] = T_\delta[(i + \ell) \dots (i + r)]$, then $A_i[j] = A_i[\ell]$ for all $j \in \{\ell, \dots, r\}$.*

Proof. Suppose that $P_\delta[\ell \dots r] = T_\delta[(i + \ell) \dots (i + r)]$. We proceed by induction on $j \in \{\ell, \dots, r\}$. The base case $j = \ell$ is tautologically true. For the inductive step, let $j \in \{\ell + 1, \dots, r\}$. By the inductive hypothesis, we have that $A_i[j - 1] = T[i + j - 1] - P[j - 1] = A_i[\ell]$. As $P_\delta[j - 1] = T_\delta[i + j - 1]$, by Definition 4.1 (and rearranging the equation), we have $A_i[j] = T[i + j] - P[j] = T[i + j - 1] - P[j - 1] = A_i[\ell]$. \square

In section 4.1 we discussed that $d_H^+(i)$ equals $m - \max_\alpha |\{j \mid A_i[j] = \alpha\}|$, which could be found by sorting and scanning A_i in $O(m \log m)$ time. However, we now have A_i in run-length encoded form (with $O(k)$ runs), and therefore the time taken to find $d_H^+(i)$ is reduced to $O(k \log k)$. Over all alignments, this gives $O(nk \log k)$ time as desired.

We can now give an overview of our deterministic algorithm for SHIFT- k -MISMATCH. The steps are described in Algorithm 3, and the overall running time is given in Theorem 4.4.

THEOREM 4.4. *The shift-normalized k -mismatch problem (SHIFT- k -MISMATCH) can be solved deterministically in $O(nk \log k)$ time.*

Proof. The solution is outlined in Algorithm 3. Correctness follows directly from the discussion in this section. The time complexity of the five steps is as follows. By

inspection of the definition, the difference strings computed in step 1 require $O(n)$ time. Step 2 uses a $2k$ -mismatch algorithm as a black box and can be performed in $O(nk)$ time by using, for example, the algorithm in [20]. Step 3 makes a single pass of the output of the $2k$ -mismatch algorithm in $O(n)$ time. Step 4 constructs a run-length encoded version of A_i for each undiscarded i . This requires scanning the $O(k)$ mismatches at each undiscarded alignment. Therefore step 4 takes $O(nk)$ time. Step 5 scans and sorts each A_i , which takes $O(k \log k)$ time per alignment as A_i is encoded by $O(k)$ runs. Overall the algorithm requires $O(nk \log k)$ time as claimed. \square

4.3. An improved, randomized solution. We now present an improved solution to the shift-normalized k -mismatch problem which runs in $O(cn\sqrt{k \log k} \log n)$ time. The improved algorithm is for the case that $k < \sqrt{m/6}$ and is randomized. The errors are one-sided (false-positives) and it outputs the correct answer at all alignments with probability at least $1 - 1/n^c$ for any constant c . For each position i , the algorithm gives a yes/no answer to the question “is $d_H^+(i) \leq k$?” The algorithm does not output the actual distance $d_H^+(i)$. Throughout this section, we use T_i as shorthand for $T[i \dots i + m - 1]$.

In section 4.2 our deterministic algorithm made use of the locations of mismatches in the difference strings P_δ and $T_\delta[i \dots i + m - 1]$. Recall that the difference string S_δ was defined to give the differences between consecutive positions in a string S . That is, $S_\delta[j] = S[j + 1] - S[j]$ for all j . A key observation was that $P_\delta[j] = T_\delta[i + j]$ if and only if $P[j] - T[i + j] = P[j + 1] - T[i + j + 1] = -\alpha$; i.e., the positions of $P[j]$ and $P[j + 1]$ require the same shift α to match. However, there is no reason to consider only consecutive differences. In fact, as we will see, one may consider differences under any arbitrary permutation of the position set. This notion is formalized as follows.

DEFINITION 4.5. Let S be a string of length s and $\pi : [m] \rightarrow [m]$ be a permutation. The permuted difference string of S under π , denoted S_π , is defined by

$$S_\pi[j] = S[\pi(j)] - S[j]$$

for all $j \in [s]$. The length of S_π is s .

Note that the permuted difference string S_π has length $|S|$ in contrast to the difference string S_δ of Definition 4.1, which has length $|S| - 1$.

The central idea of our improved algorithm is to use the value of $\text{Ham}(P_\pi, (T_i)_\pi)$ to directly determine whether $d_H^+(i) \leq k$ at each alignment i . In Definition 4.6 we introduce the notion of a permutation being k -tight for some P, T_i . Intuitively, π is k -tight for P, T_i if we can infer directly from $\text{Ham}(P_\pi, (T_i)_\pi)$ whether $d_H^+(i) \leq k$.

DEFINITION 4.6. Let π be a permutation, P a pattern, and T_i a text substring. We say that π is k -tight for P, T_i if

$$d_H^+(i) \leq k \iff \text{Ham}(P_\pi, (T_i)_\pi) \leq 2k.$$

It would, of course, be highly desirable to find a permutation π which is k -tight for all P, T_i and any k . However, we will see that this is in general not possible.

We begin by showing that any π has the property that $d_H^+(i) \leq k$ implies that $\text{Ham}(P_\pi, (T_i)_\pi) \leq 2k$ for all P, T_i . To do so we first prove a general lemma which will also be useful later. Lemma 4.8 then gives the desired property and is a generalization of Lemma 4.2 to arbitrary permutations.

LEMMA 4.7. Let π be a permutation, P a pattern, and T_i a text substring. For all $j \in [m]$,

$$P[j] - T_i[j] = P[\pi(j)] - T_i[\pi(j)] \iff P_\pi[j] = (T_i)_\pi[j].$$

Proof. The left-hand side of the arrow is the same as $P[\pi(j)] - P[j] = T_i[\pi(j)] - T_i[j]$, which by Definition 4.5 is equivalent to the right-hand side of the arrow. \square

LEMMA 4.8. *Let π be a permutation, P a pattern, and T_i a text substring.*

$$d_H^+(i) \leq k \implies \text{Ham}(P_\pi, (T_i)_\pi) \leq 2k.$$

Proof. Let P and T_i be such that $d_H^+(i) \leq k$. By definition there exists an α such that the set $J = \{j \mid P[j] + \alpha \neq T_i[j]\}$ has size at most k . As π is a permutation, there are at most $2k$ positions $j \in [m]$ such that either $j \in J$ or $\pi(j) \in J$. Therefore, for all (at least) $m - 2k$ remaining positions $j' \in [m]$ we have that $P[j'] + \alpha = T_i[j']$ and $P[\pi(j')] + \alpha = T_i[\pi(j')]$. For each such position j' , by rearranging the two equations it follows from Lemma 4.7 that $P_\pi[j'] = (T_i)_\pi[j']$. Thus, there are at most $2k$ mismatches between P_π and $(T_i)_\pi$. \square

A logical next step would be to attempt to find a permutation π with the property that $\text{Ham}(P_\pi, (T_i)_\pi) \leq 2k$ implies that $d_H^+(i) \leq k$ for all P, T_i . Unfortunately, Lemma 4.9 shows that no such permutation can exist. As Corollary 4.10 states, this immediately implies that there is no permutation which is k -tight for all P, T_i . Instead we will select our permutation at random and show that we can obtain a permutation that is k -tight for a given P, T_i with constant probability.

LEMMA 4.9. *Let π be any permutation without fixed points and $6 \leq k < m/4$. There exist a pattern P and text substring T_i such that*

$$d_H^+(i) > k \quad \text{and} \quad \text{Ham}(P_\pi, (T_i)_\pi) \leq 2k.$$

Proof. We define P to be an m -length string of zeros. In order to define the m -length string T_i we first introduce some notation.

Let $k' = \lfloor k/2 \rfloor + 1$. We identify a set of k' locations $\ell_0, \ell_1, \dots, \ell_{k'-1} \in [m]$ defined as follows. Location $\ell_0 = 0$, and for $q \in \{1, \dots, k' - 1\}$, location ℓ_q is the smallest position in $[m]$ that is not any of the preceding locations $\ell_0, \dots, \ell_{q-1}$ or any location that is mapped to or from by any of these preceding locations (under π). Formally, ℓ_q is the smallest location which is not in the set $L_q = \{\ell_{q'}, \pi(\ell_{q'}), \pi^{-1}(\ell_{q'}) \mid q' \in [q]\}$. Observe that the set L_q has size at most $3k' \leq 3(k/2 + 1) < 3k < m$ (since $k < m/4$), and hence such a location always exists.

We can now define T_i as follows. For all $q \in [k']$, let $T_i[\ell_q] = 1$ and $T_i[\pi(\ell_q)] = 1$. At all other locations j , $T_i[j] = 0$. Observe that by construction and the assumption that π has no fixed points, the locations $\ell_0, \dots, \ell_{k'-1}$ and $\pi(\ell_0), \dots, \pi(\ell_{k'-1})$ are all distinct. Therefore, T_i contains exactly $2k'$ ones and $m - 2k'$ zeros. As $2k' \leq k + 2 < m/2$, more than half the locations have $T_i[j] = P[j] = 0$, and therefore $d_H^+(i)$ is minimized by the shift $\alpha = 0$. Thus, $d_H^+(i) = 2k' > k$.

We proceed by showing that the alignment of P_π and $(T_i)_\pi$ contains at least $m - 3k'$ matches. There are $m - 2k'$ locations j in T_i such that $T_i[j] = 0$. Of these locations, at most $2k'$ have $T_i[\pi(j)] = 1$. Therefore, there are at least $m - 4k'$ locations j such that $T_i[j] = T_i[\pi(j)] = 0$. As $P[j] = P[\pi(j)] = 0$, we have by Lemma 4.7 that $P_\pi[j] = (T_i)_\pi[j]$ at $m - 4k'$ locations. Now consider locations ℓ_q for $q \in [k']$. By construction, $T_i[\ell_q] = T_i[\pi(\ell_q)] = 1$, and therefore $P_\pi[\ell_q] = (T_i)_\pi[\ell_q]$ by Lemma 4.7. This implies further k' matching locations. There are therefore at least $m - 3k'$ matches or at most $3k'$ mismatches between P_π and $(T_i)_\pi$. Since $3k' \leq 3(k/2 + 1) \leq 2k$ for all $k \geq 6$ we have that $\text{Ham}(P_\pi, (T_i)_\pi) \leq 2k$. \square

COROLLARY 4.10. *Let π be any permutation and $6 \leq k < m/4$. There exist a pattern P and text substring T_i for which π is not k -tight.*

Proof. The proof is immediate from Definition 4.6 and Lemma 4.9. \square

4.3.1. Random permutations. We will choose a permutation uniformly at random from a simple family of permutations. On first inspection, we could have chosen from the family of all permutations. We claim without proof that a permutation chosen uniformly at random from the family of all permutations is k -tight for any P, T_i with constant probability. However, we must be able to efficiently compute $\text{Ham}(P_\pi, (T_i)_\pi)$ for all i under our chosen permutation. The key problem is that in general $(T_i)_\pi$ is not easily obtained from T . As i varies, $(T_i)_\pi$ could change drastically, even when i is incremented only by one. Therefore we must be careful in selecting our family of permutations.

We will use the family of cyclic permutations, denoted \mathbb{C}_m (for patterns of length m), defined as follows.

DEFINITION 4.11. *The set \mathbb{C}_m contains the $m - 1$ cyclic permutations $\pi_1, \pi_2, \dots, \pi_{m-1}$, where*

$$\pi_q(j) = j + q \bmod m.$$

We now show in Lemma 4.12 that \mathbb{C}_m has the desired property of k -tightness when $m > 6k^2$. There is a corner case when $k \in \{0, 1\}$, which is easily solved in $O(n)$ time using our deterministic algorithm from section 4.2. For Lemma 4.12 we require that $k \geq 2$.

LEMMA 4.12. *Let P be a pattern and T_i a text substring. When $m > 6k^2$ and $k \geq 2$,*

$$\frac{|\{ \pi \mid \pi \in \mathbb{C}_m \text{ is } k\text{-tight for } P, T_i \}|}{|\mathbb{C}_m|} \geq \frac{1}{6}.$$

Proof. Let $\rho = |\{ \pi \mid \pi \in \mathbb{C}_m \text{ is } k\text{-tight for } P, T_i \}| / |\mathbb{C}_m|$. We will show that $\rho \geq 1/6$. Note that $|\mathbb{C}_m| = m - 1$. We let $h = d_H^+(i)$ be the minimal number of mismatches between P and T_i , and let $\hat{\alpha}$ be the shift which minimizes $d_H^+(i)$.

Assume first that $h \leq k$. By Lemma 4.8 and Definition 4.6 we have that every $\pi \in \mathbb{C}_m$ is k -tight for P, T_i , and therefore $\rho = 1$. Assume second that $h > k$. We split the proof into three cases:

$$\text{Case 1. } k < h \leq 2k. \quad \text{Case 2. } 2k < h \leq \frac{m}{3}. \quad \text{Case 3. } \frac{m}{3} < h$$

First we introduce some notation. There are exactly $m - h$ positions j where $\hat{\alpha} + P[j] = T_i[j]$. We call such a position an $\hat{\alpha}$ -match. Similarly, any position with $\alpha + P[j] = T_i[j]$ for some α is called an α -match. Positions which are not α -matches are called α -mismatches. Hence there are h distinct $\hat{\alpha}$ -mismatches. We will refer to $\pi(j)$ as the position that j is mapped to (by π).

Case 1 ($k < h \leq 2k$). Let j be an arbitrary $\hat{\alpha}$ -mismatch. Position j is mapped to another $\hat{\alpha}$ -mismatch in exactly $h - 1$ distinct permutations of \mathbb{C}_m . This holds for each of the h distinct $\hat{\alpha}$ -mismatches. Hence there are at most $(h - 1)h$ permutations under which some $\hat{\alpha}$ -mismatch is mapped to another $\hat{\alpha}$ -mismatch. The remaining (at least) $(m - 1) - (h - 1)h$ permutations π in \mathbb{C}_m immediately have the following two properties:

- (i) if position j is an $\hat{\alpha}$ -mismatch, then $\pi(j)$ is an $\hat{\alpha}$ -match;
- (ii) if position $\pi(j)$ is an $\hat{\alpha}$ -mismatch, then position j is an $\hat{\alpha}$ -match.

There are h positions j with property (i) and another (disjoint) h positions j with property (ii). That is, for each $\hat{\alpha}$ -mismatch there are two positions j that meet one of

the two properties above. By Lemma 4.7, each such j implies that $P_\pi[j] \neq (T_i)_\pi[j]$. Therefore, in each of these $(m-1) - (h-1)h$ permutations π , $\text{Ham}(P_\pi, (T_i)_\pi) \geq 2h > 2k$, and so each such permutation is k -tight for P, T_i . By the assumption of Case 1, $h \leq 2k$, and the assumptions that $m \geq 6k^2$ and $k \geq 2$, we have that $(m-1) - (h-1)h \geq m - 4k^2 \geq m/3$. Thus, $\rho \geq (m/3)/(m-1) > 1/6$.

Case 2 ($2k < h \leq m/3$). Let K be an arbitrary set of $2k$ distinct $\hat{\alpha}$ -mismatches. For any permutation π , let

$$K_\pi^{-1} = \{ j \mid \pi(j) \in K \}.$$

We define

$$\text{Ham}_K(P_\pi, (T_i)_\pi) = |\{ j \mid j \in (K \cup K_\pi^{-1}) \wedge P_\pi[j] \neq (T_i)_\pi[j] \}|$$

to be the number of mismatch positions between P_π and $(T_i)_\pi$ that are also in K or K_π^{-1} . We now consider the total number of mismatches between P_π and $(T_i)_\pi$ (that are in K or K_π^{-1}) summed over all permutations in \mathbb{C}_m . Let

$$H_K(P, T_i) = \sum_{\pi \in \mathbb{C}_m} \text{Ham}_K(P_\pi, (T_i)_\pi).$$

Since $h \leq m/3$ by the assumption of Case 2, there are at least $2m/3$ $\hat{\alpha}$ -matches. A permutation π that maps a position $j \in K$ to an $\hat{\alpha}$ -match creates a mismatch $P_\pi[j] \neq (T_i)_\pi[j]$ by Lemma 4.7 (as j is an $\hat{\alpha}$ -mismatch). For a fixed $j \in K$, the number of permutations in \mathbb{C}_m that map j to an $\hat{\alpha}$ -match equals the number of $\hat{\alpha}$ -matches, which is at least $2m/3$. Thus, the set K of $2k$ $\hat{\alpha}$ -mismatches contributes at least $2k \cdot (2m/3)$ to $H_K(P, T_i)$.

Similarly, any position j which is an $\hat{\alpha}$ -match creates a mismatch $P_\pi[j] \neq (T_i)_\pi[j]$ by Lemma 4.7 if it is mapped to an $\hat{\alpha}$ -mismatch in K . This occurs under exactly $2k$ permutations. Recall that any j which is mapped to a position in K under π belongs to K_π^{-1} . Therefore, given that there are at least $2m/3$ $\hat{\alpha}$ -matches, the contribution is at least $2k \cdot (2m/3)$ further distinct mismatches to $H_K(P, T_i)$.

Summing up the previous two paragraphs, we have shown that $H_K(P, T_i) \geq (8/3)mk$. Each permutation π that is not k -tight for P, T_i has $\text{Ham}(P_\pi, (T_i)_\pi) \leq 2k$ (since $h > k$). Therefore, $m \cdot 2k$ is a generous upper bound on the number of mismatches across all permutations which are not k -tight. This leaves at least $(8/3)mk - 2mk = (2/3)mk$ mismatches among the k -tight permutations of \mathbb{C}_m . Since $|K \cup K_\pi^{-1}| \leq 4k$, we have that $\text{Ham}_K(P_\pi, (T_i)_\pi) \leq 4k$ for any π , and hence each permutation contributes at most $4k$ mismatches to $H_K(P, T_i)$. Therefore there are at least $(2/3)mk/(4k) = m/6$ distinct k -tight permutations. Thus, $\rho \geq (m/6)/(m-1) \geq 1/6$.

Case 3 ($m/3 < h$). Similarly to Case 2, we consider the total number of mismatches between P_π and $(T_i)_\pi$ summed over all permutations in \mathbb{C}_m . Let

$$H(P, T_i) = \sum_{\pi \in \mathbb{C}_m} \text{Ham}(P_\pi, (T_i)_\pi).$$

Since $h > m/3$, the number of α -mismatches is more than $m/3$ for all α . Fix an arbitrary position j and choose an α such that j is an α -match. There are at least $m/3$ permutations π in \mathbb{C}_m that map position j to an α -mismatch. By Lemma 4.7, $P_\pi[j] \neq (T_i)_\pi[j]$ for each of these permutations. Hence position j will contribute with

at least $m/3$ to $H(P, T_i)$. By considering all m positions j , we have that $H(P, T_i) \geq m \cdot (m/3)$.

Similarly to the reasoning in Case 2, each permutation π that is not k -tight for P, T_i has $\text{Ham}(P_\pi, (T_i)_\pi) \leq 2k$ (since $h > k$). Again, $m \cdot 2k$ is a generous upper bound on the number of mismatches across all permutations which are not k -tight. This leaves at least $m^2/3 - 2mk$ mismatches among the k -tight permutations of \mathbb{C}_m . As certainly $\text{Ham}(P_\pi, (T_i)_\pi) \leq m$, we have that there are at least $m/3 - 2k$ distinct k -tight permutations for P, T_i . Therefore,

$$\rho \geq \frac{m/3 - 2k}{m - 1} \geq \frac{k - 1}{3k} \geq \frac{1}{6},$$

where the second inequality follows from $m > 6k^2$ and the last inequality from $k \geq 2$, both assumptions in the statement of the lemma. \square

4.3.2. The algorithm. Before describing the randomized algorithm we turn our attention to the problem of finding all positions $i \in [n - m + 1]$ such that $\text{Ham}(P_\pi, (T_i)_\pi) \leq 2k$ under an arbitrary cyclic permutation $\pi \in \mathbb{C}_m$. We will describe a simple deterministic algorithm that computes $\text{Ham}(P_\pi, (T_i)_\pi)$ by reduction to the conventional k -mismatch problem.

Let $\pi_q \in \mathbb{C}_m$ be a fixed but arbitrary permutation ($q \in [1, \dots, m - 1]$). Recall that $\pi_q(j) = j + q \bmod m$. We define

$$\begin{aligned} P_q^+ &= P_{\pi_q}[0 \dots (m - q - 1)], \\ P_q^- &= P_{\pi_q}[(m - q) \dots (m - 1)]. \end{aligned}$$

Thus, $P_{\pi_q} = P_q^+ \| P_q^-$. We have $|P_q^+| = m - q$ and $|P_q^-| = q$. Now define T_q^+ and T_q^- such that

$$\begin{aligned} T_q^+[j] &= T[j + q] - T[j], \\ T_q^-[j] &= T[j + q - m] - T[j] \end{aligned}$$

for all $j \in [n]$ (except those that take the indices “out of range”). Observe that

$$(T_i)_{\pi_q} = T_q^+[i \dots (i + m - q - 1)] \| T_q^-[(i + m - q) \dots (i + m - 1)],$$

where the first substring has length $m - q$ and the second substring has length q . From these definitions it now follows directly that

$$\begin{aligned} (4.2) \quad \text{Ham}(P_{\pi_q}, (T_i)_{\pi_q}) &= \text{Ham}(P_q^+, T_q^+[i \dots (i + m - q - 1)]) \\ &\quad + \text{Ham}(P_q^-, T_q^-[(i + m - q) \dots (i + m - 1)]). \end{aligned}$$

Thus, in order to determine which positions i have $\text{Ham}(P_{\pi_q}, (T_i)_{\pi_q}) \leq 2k$, we first construct P_q^+ , P_q^- , T_q^+ , and T_q^- , and then we run a standard $2k$ -mismatch algorithm on the pairs (P_q^+, T_q^+) and (P_q^-, T_q^-) and use the previous formula.

We can now finally give an overview of our randomized algorithm for the SHIFT- k -DECISION problem. The steps are described in Algorithm 4. The overall running time and proof of correctness is given in Theorem 4.13. The algorithm makes one-sided errors and outputs a false match (incorrectly reports $d_H^+(i) \leq k$) with constant probability per alignment. As we will see in the proof of Theorem 4.13, by running the algorithm a logarithmic number of times drastically reduces the probability of an error occurring at one or more alignments.

ALGORITHM 4. Overview of randomized solution to SHIFT- k -DECISION.

1. Pick a cyclic permutation $\pi_q \in \mathbb{C}_m$ uniformly at random.
2. Construct the strings P_q^+ , P_q^- , T_q^+ , and T_q^- .
3. Run a $2k$ -mismatch algorithm on the pairs (P_q^+, T_q^+) and (P_q^-, T_q^-) as a black box.
4. Using the results from step 3 and (4.2), compute $\text{Ham}(P_{\pi_q}, (T_i)_{\pi_q})$ for all i .
5. Any alignment i with $\text{Ham}(P_{\pi_q}, (T_i)_{\pi_q}) \leq 2k$ is declared to have $d_H^+(i) \leq k$.

THEOREM 4.13. *For any choice of constant c , SHIFT- k -DECISION can be solved randomized in $O(cn\sqrt{k \log k} \log n)$ (deterministic) time when $k < \sqrt{m/6}$. The algorithm makes only false-positive errors (incorrectly declares that the Hamming distance is at most k). With probability at least $1 - 1/n^c$, the algorithm is correct at every alignment.*

Proof. As discussed in section 4.3.1, if $k \in \{0, 1\}$, then we can use the deterministic algorithm from section 4.2 and achieve time complexity of $O(n)$ and no errors. Therefore, we focus on the case that $k \geq 2$.

We first consider correctness. It follows from the discussion above that Algorithm 4 does indeed determine, for every alignment i , whether $\text{Ham}(P_{\pi_q}, (T_i)_{\pi_q}) \leq 2k$. We first show that

- (i) $\text{Ham}(P_{\pi_q}, (T_i)_{\pi_q}) \leq 2k$ when $d_H^+(i) \leq k$;
- (ii) the probability that $\text{Ham}(P_{\pi_q}, (T_i)_{\pi_q}) \leq 2k$ when $d_H^+(i) > k$ is at most $5/6$.

By Lemma 4.8 we have that if $d_H^+(i) \leq k$, then $\text{Ham}(P_{\pi_q}, (T_i)_{\pi_q}) \leq 2k$. This proves property (i). By Definition 4.6, $\text{Ham}(P_{\pi_q}, (T_i)_{\pi_q}) > 2k$ if $d_H^+(i) > k$ for all permutations π_q that are k -tight for P, T_i . The permutation π_q is selected uniformly at random from \mathbb{C}_m in step 1, and hence by Lemma 4.12 it is k -tight for P, T_i with probability at least $1/6$. This proves property (ii). Note that we can apply Lemma 4.12 since we have assumed that $m > 6k^2$ and $k \geq 2$.

As Algorithm 4 makes only false-positive errors, we can amplify the probability of giving correct outputs by repeating the algorithm. We repeat it $4(c+1)\lceil \log n \rceil$ times, where c is a constant, and output any alignment which is reported by all repeats. More precisely, let i be some alignment such that $d_H^+(i) > k$. The probability that one run of Algorithm 4 incorrectly reports position i as a match is at most $5/6$. Thus, the probability that all runs output i as a match is at most

$$(5/6)^{4(c+1)\lceil \log n \rceil} < (1/2)^{(c+1)\log n} < 1/n^{c+1}.$$

By the union bound over all positions i , the probability of the multirun algorithm outputting a false match in at least one alignment is at most $n \cdot 1/n^{c+1} = 1/n^c$ as required.

We now consider the time complexity of Algorithm 4 (without amplification). Step 1 requires only constant time to pick a permutation at random. Step 2 requires $O(n)$ time by inspection of the definitions. Step 3 makes two calls to a $2k$ -mismatch algorithm. For both calls the input is a pattern of length $O(m)$ and a text of length $O(n)$. Using the fastest known k -mismatch algorithm of Amir, Lewenstein, and Porat [4], this step takes $O(n\sqrt{k \log k})$ time. Steps 4 and 5 require only scanning the output of step 3 and therefore take $O(n)$ time. This gives a time complexity of $O(n\sqrt{k \log k})$ time. However, we repeat the algorithm $O(c \log n)$ times to reduce the error probability, and hence $O(cn\sqrt{k \log k} \log n)$ is the total time complexity. \square

5. Discussion. We have shown how to derive both new upper and lower bounds for a variety of pattern matching problems under polynomial transformations. In some cases we have improved on known results, and in others introduced new problem definitions and solutions. There remain, however, a number of open questions. First, we suspect that the true complexity of $\text{POLY-}r\text{-}L_2^*$ is unresolved, particularly for higher degree polynomial transformations. For example, when $r = m$, there exists a straightforward $O(nm)$ time solution. It is also still uncertain whether the normalized Hamming distance problem is 3SUM-hard for polynomials of degree greater than one. For $\text{SHIFT-}k\text{-DECISION}$, our fast randomized algorithm applies only when $k < \sqrt{m/6}$. However, our lower bound for the same problem applies to the case where we want to determine whether the Hamming distance is at most $m - 2$. This leaves a range of values of k , where the complexity is not yet determined. Another interesting question is whether our randomized solution can be efficiently modified to output the Hamming distance at each alignment rather than simply a decision about whether it is greater or less than k or, indeed, whether a new fast method can be found for this problem which will allow the presence of wildcards in the input.

Acknowledgment. The authors are grateful to Philip Bille for very helpful discussions on the topic of 3SUM and related problems and their application to lower bounds for pattern matching problems.

REFERENCES

- [1] K. ABRAHAMSON, *Generalized string matching*, SIAM J. Comput., 16 (1987), pp. 1039–1051.
- [2] A. AMIR, Y. AUMANN, M. LEWENSTEIN, AND E. PORAT, *Function matching*, SIAM J. Comput., 35 (2006), pp. 1007–1022.
- [3] A. AMIR AND M. FARACH, *Efficient 2-dimensional approximate matching of half-rectangular figures*, Inform. Comput., 118 (1995), pp. 1–11.
- [4] A. AMIR, M. LEWENSTEIN, AND E. PORAT, *Faster algorithms for string matching with k mismatches*, J. Algorithms, 50 (2004), pp. 257–275.
- [5] A. AMIR, O. LIPSKY, E. PORAT, AND J. UMANSKI, *Approximate matching in the L_1 metric*, in CPM '05: Proc. 16th Annual Symp. on Combinatorial Pattern Matching, Springer-Verlag, Berlin, 2005, pp. 91–103.
- [6] M. J. ATALLAH, *Faster image template matching in the sum of the absolute value of differences measure*, IEEE Trans. Image Process., 10 (2001), pp. 659–663.
- [7] I. BARAN, E. D. DEMAINE, AND M. PĂTRAȘCU, *Subquadratic algorithms for 3SUM*, in WADS '05: Proc. 9th International Conference on Algorithms and Data Structures, Springer-Verlag, Berlin, 2005, pp. 409–421.
- [8] P. CLIFFORD AND R. CLIFFORD, *Simple deterministic wildcard matching*, Inform. Proces. Lett., 101 (2007), pp. 53–54.
- [9] P. CLIFFORD, R. CLIFFORD, AND C. ILIOPOULOS, *Faster algorithms for δ, γ -matching and related problems*, in CPM '05: Proc. 16th Annual Symp. on Combinatorial Pattern Matching, Springer-Verlag, Berlin, 2005, pp. 68–78.
- [10] R. CLIFFORD AND C. ILIOPOULOS, *String algorithms in music analysis*, Soft Computing, 8 (2004), pp. 597–603.
- [11] R. COLE AND R. HARIHARAN, *Verifying candidate matches in sparse and wildcard matching*, in STOC '02: Proc. 34th Annual ACM Symp. Theory of Computing, ACM, New York, 2002, pp. 592–601.
- [12] D. COPPERSMITH AND S. WINOGRAD, *Matrix multiplication via arithmetic progressions*, J. Symbolic Computation, 9 (1990), pp. 251–280.
- [13] T. H. CORMEN, C. E. LEISERSON, AND R. L. RIVEST, *Introduction to Algorithms*, MIT Press, Cambridge, MA, 1990.
- [14] J. ERICKSON, *Lower bounds for linear satisfiability problems*, Chicago J. Theoret. Comput. Sci., 1999, 8 (electronic).
- [15] J. ERICKSON, *New lower bounds for convex hull problems in odd dimensions*, SIAM J. Comput., 28 (1999), pp. 1198–1214.

- [16] J. ERICKSON AND R. SEIDEL, *Better lower bounds on detecting affine and spherical degeneracies*, Discrete Comput. Geom., 13 (1995), pp. 41–57.
- [17] A. GAJENTAAN AND M. H. OVERMARS, *On a class of $O(n^2)$ problems in computational geometry*, Comput. Geom., 5 (1995), pp. 165–185.
- [18] J. KING, *A Survey of 3SUM-Hard Problems*, <http://www.cs.mcgill.ca/~jking/papers/3sumhard.pdf> (2004).
- [19] S. R. KOSARAJU, *Efficient String Matching*, manuscript.
- [20] G. M. LANDAU AND U. VISHKIN, *Efficient string matching with k mismatches*, Theoret. Comput. Sci. 43 (1986), pp. 239–249.
- [21] K. LEMSTRÖM AND E. UKKONEN, *Including interval encoding into edit distance based music comparison and retrieval*, in AISB '00: Proc. Symposium on Creative and Cultural Aspects and Applications of AI and Cognitive Science, Society for the Study of Artificial Intelligence and the Simulation of Behaviour, London, 2000, pp. 53–60.
- [22] J. P. LEWIS, *Fast template matching*, in Vision Interface '95, Canadian Image Processing and Pattern Recognition Society, Quebec, 1995, pp. 120–123.
- [23] O. LIPSKY AND E. PORAT, *Approximate matching in the L_∞ metric*, Inform. Process. Lett., 105 (2008), pp. 138–140.
- [24] V. MÄKINEN, G. NAVARRO, AND E. UKKONEN, *Transposition invariant string matching*, J. Algorithms, 56 (2005), pp. 124–153.