

Two-dimensional pattern matching with k mismatches

Jonas Ellert¹, Paweł Gawrychowski², Adam Górkiewicz³, and Tatiana Starikovskaya⁴

¹?

²?

³?

⁴?

Abstract

1 Related Work

J: work in progress ↓, unfinished

Hamming-Threshold PM in 1D

•

J: ...

Exact PM in 2D Assuming $n \times n$ text and $m \times m$ pattern, $m \leq n$.

- Two-dimensional pattern matching dates back (at least) to the 1970s, when Richard Bird generalized the Knuth-Morris-Pratt algorithm [38] for one-dimensional pattern matching, achieving $\mathcal{O}(n^2 + m^2)$ time [17].
- By using multiple pattern matching on only $\mathcal{O}(n/m)$ rows of the text, Baeza-Yates and Régnier achieve average time complexity $\mathcal{O}(n^2/m)$, using $\mathcal{O}(m^2)$ space and still running in $\mathcal{O}(n^2)$ time in the worst case [14]. This was followed by more results optimizing the average time complexity [34, 48].
- Crochremore et al. [26] solve two-dimensional pattern matching space efficiently, using an $\mathcal{O}(m^2)$ time and $\mathcal{O}(\log m)$ space preprocessing followed by an $\mathcal{O}(n^2)$ time and constant space search phase.
- On a PRAM, the problem can be solved in constant time and $\mathcal{O}(n^2)$ work [25], after preprocessing the pattern in $\mathcal{O}(\log \log m)$ time and $\mathcal{O}(m^2)$ work [23].
- Other works studying the problem are, e.g., [3, 15, 29, 37, 49, 50]

Hamming-Threshold PM in 2D Assuming $n \times n$ text and $m \times m$ pattern, $m \leq n$, and threshold $k < m^2$.

- First result in 1987 by Krithivasan and Sitalakshmi [39], who show that $\tilde{\mathcal{O}}(kmn^2)$ time can be achieved.
- Faster algorithm by Rank and Heywood runs in $\tilde{\mathcal{O}}((k + m)n^2)$ time [46].

- Current state of the art is $\mathcal{O}(kn^2)$ [6]
- More work [13, 34, 44] is aimed at optimizing the average time complexity. This has also been done for the setting in which arbitrary rotations of the pattern are allowed [28].
- Dynamic and online algorithms [22]

J: check

2D Periodicity

- Two-dimensional periodicity was introduced in [1, 2] and refined in [29]. All periods (and witnesses for all non-periods) can be computed in $\mathcal{O}(n^2)$ time [24]. Witnesses can be used to efficiently eliminate potential occurrences of periodic patterns. Since witness computation is costly for run-length encoded strings, Amir et al. [4] show additional properties of two-dimensional periodicities that allow the elimination of occurrences without witnesses.
- Famous periodicity theorems have been generalized to two dimensions, for example the Fine-Wilf theorem [42] and the Lyndon-Schützenberger theorem [31].
- Smallest rectangular cover in $\mathcal{O}(n^2)$ time, all aperiodic covers in $\mathcal{O}(n^2 \log n)$ time [20]. Covers of two-dimensional strings were also considered in [27, 30].
- All tile covers in $\mathcal{O}(n^{2+\epsilon})$ time [45].
- There can be $\Omega(n^3 \log n)$ tandem substrings (repeated rectangle, like square in a normal string) [10].
- All tandems can be computed in $\mathcal{O}(n^3 \log n)$ time [11].
- There can be $\Omega(n^2 \log n)$ two-dimensional runs (maximal repetitions) [32].
- There are at most $\mathcal{O}(n^2 \log^2 n)$ two-dimensional runs [19].
- All runs can be computed in $\mathcal{O}(n^2 \log^2 n)$ time [7].

Other 2D Shannanigans

- A survey by Rytter [47] discusses compressed representations of two-dimensional strings and the computation on such representations without explicit decompression. If the text is compressed as a two-dimensional SLP, then merely deciding if a pattern is present is NP-complete with respect to the compressed size [16]. If both text and pattern are compressed using a two-dimensional version of Lempel-Ziv compression, then exact matching can be done in $\mathcal{O}(n^2 + m^3)$ time, using additional working space proportional to the compressed representation of the pattern [8]. If text and pattern are given in two-dimensional run-length encoding, then the time for exact matching is linear in the size of the compressed text, and the space is linear in the size of the compressed pattern [9]. If multiple patterns are given, then they can be compressed into a space efficient self-index that allows fast dictionary matching [43].
- More work on structural elements in two-dimensional strings, e.g., frames (rectangles with matching outer columns and rows) [18], Lyndon words [41], motif patterns (allowing “don’t cares”, and satisfying special properties) [12], patterns with at least two occurrences [36], and palindromic substrings [40].
- There is an efficient algorithm for matching a half-rectangular pattern (fixed height, but irregular shape) in a text, allowing a fixed number of mismatches, insertions and deletions [5].

J: work in progress, unfinished ↑

2 Introduction

We consider the one-dimensional all-substring Hamming distance problem (HD1D), where for a given text string T of length n and a string P of length m ($m < n$), we want to calculate the Hamming distance between P and every fragment T of length m .

We consider the two-dimensional all-substring Hamming distance problem (HD2D), where for a given 2D string T of size $n \times n$ and a string P of size $m \times m$ ($m < n$), we want to calculate the Hamming distance between P and every $m \times m$ fragment of T .

We also consider the bounded variants of HD1D and HD2D, where we are only required to calculate the distances which are not greater than k , for some parameter $k \in \mathbb{Z}^+$.

Theorem 1 (Main result). *Bounded HD2D can be solved in $\tilde{O}((m^2 + mk^{5/4})n^2/m^2)$ time.*

3 Preliminaries

For our purposes we will not use the standard definition of a two-dimensional string, where we associate it with a two-dimensional array of characters, and instead we will define it more broadly. Although we will occasionally use the array notation, we will do it exclusively for $n \times m$ strings. For any $n \in \mathbb{Z}^+$ we will denote $[n] = \{0, \dots, n-1\}$. We will use the terms point and vector interchangeably. Our results hold under word-RAM model of computation.

Definition 1 (Two-dimensional string). We define a **string** S as a partial function $\mathbb{Z}^2 \rightarrow \Sigma$ which maps some arbitrary set of integer points, denoted as $\text{dom}(S)$, to characters. For simplicity we will write $u \in S$ to denote that $u \in \text{dom}(S)$. We say that a string S is **partitioned** into strings R_1, \dots, R_ℓ when the sets $\text{dom}(R_1), \dots, \text{dom}(R_\ell)$ partition $\text{dom}(S)$ and $R_i(u) = S(u)$ for all $u \in R_i$. We call a string S **monochromatic** when $S(u) = \sigma$ for every $u \in S$ for some $\sigma \in \Sigma$ and we will write $C(S)$ to denote the value σ . We say that S is $n \times m$ for some $n, m \in \mathbb{Z}^+$ when $\text{dom}(S) = [n] \times [m]$. Physically we represent a string as a list of point-character pairs.

Definition 2 (Shifting). For a set of points $V \subseteq \mathbb{Z}^2$ and a vector $u \in \mathbb{Z}^2$, we denote $V + u$ as $\{v + u : v \in V\}$. For a string S and a vector $u \in \mathbb{Z}^2$ we denote $S + u$ as a string R such that $\text{dom}(R) = \text{dom}(S) + u$ and $R(v) = S(v - u)$ for $v \in \text{dom}(R)$. Intuitively, we shift the set of points while maintaining their character values.

Definition 3 (Hamming distance). For a pair of strings S, R we define

$$\text{Ham}(S, R) = |\{u : u \in \text{dom}(S) \cap \text{dom}(R), S(u) \neq R(u)\}|,$$

which corresponds to the number of mismatches between S and R .

Under such notation, the HD2D problem is equivalent to calculating the (bounded or unbounded) values of $\text{Ham}(P + q, T)$ for all $q \in \mathbb{Z}^2$ such that $\text{dom}(P + q) \subseteq \text{dom}(T)$ (so for $q \in [n - m + 1]^2$).

Definition 4 (Don't care symbol). We define the **don't care** symbol as a special character which matches with every character. We will denote it with $?$. Unless stated otherwise, we assume it is not allowed in Σ and in both HD1D and HD2D every character present in T and P matches only with itself.

Definition 5 (Vector operators). For any $u \in \mathbb{Z}^2$ we refer to its coordinates as $u.x, u.y$. For $u, v \in \mathbb{Z}^2$ we denote $u \cdot v = u.x \cdot v.x + u.y \cdot v.y$ and $u \times v = u.x \cdot v.y - u.y \cdot v.x$. Note that alternatively $u \cdot v = |u||v| \cos \alpha$ and $u \times v = |u||v| \sin \alpha$ where α is the angle between u and v .

Definition 6 (Quadrants). We define the four **quadrants** as

$$\begin{aligned}\mathcal{Q}_1 &= (0, +\infty) \times [0, +\infty), \\ \mathcal{Q}_2 &= (-\infty, 0] \times (0, +\infty), \\ \mathcal{Q}_3 &= (-\infty, 0) \times (-\infty, 0], \\ \mathcal{Q}_4 &= [0, +\infty) \times (-\infty, 0).\end{aligned}$$

4 One-dimensional generalizations

In this section we explore some of the methods used for one-dimensional strings. Specifically, as our goal is to generalize the solution for pattern matching with k mismatches described in [33], we are especially interested in two-dimensional variants of the techniques that were used to solve the one-dimensional case.

Theorem 2 (Instancing). *Consider an algorithm \mathcal{A} which solves HD2D (bounded or unbounded), but only when $2|n$ and $n \leq \frac{3}{2}m$. If its running time is $\mathcal{T}(m)$, then the general case can be solved in $\mathcal{O}(\mathcal{T}(m)n^2/m^2)$.*

Proof. Let $r = \lfloor m/2 \rfloor$ and let $n' = r + m - 1$ or $r + m$ if $r + m - 1$ is odd. We see that the set $N = [n']^2$ satisfies the conditions for the text domain. For any vector $q \in [n - m]^2$ we can find a vector u such that $r|u.x, r|u.y$ and $q - u \in [r]^2$, so we have $\text{Ham}(P + q, T) = \text{Ham}(P + q - u, T_u)$ where T_u is the restriction of $T - u$ to N . If $T - u$ is not defined for some $v \in N$, we can pad $T_u(v)$ with any character. We see that $\text{dom}(P + q - u) \subseteq N = \text{dom}(T_u)$. There are $\mathcal{O}(n^2/m^2)$ possible vectors u and we run \mathcal{A} for every pair of T_u and P . \square

Theorem 3 (Kangaroo jumps). *Consider an $n \times n$ string T , $m \times m$ string P and set of vectors Q such that $\text{dom}(P + q) \subseteq \text{dom}(T)$ for every $q \in Q$. There exists an algorithm which calculates $d_q = \text{Ham}(P + q, T)$ for every $q \in Q$ in total time $\tilde{\mathcal{O}}(n^2 + \sum_{q \in Q} d_q)$.*

Proof. For the sake of clarity, we will temporarily switch to the classical array notation for strings. Let T_0, \dots, T_{n-m} denote an array of two-dimensional strings (arrays) such that $T_k[0..n-1, 0..m-1] = T[0..n-1, k..k+m-1]$. For every row $P[0], \dots, P[m-1]$ of P and every row $T_k[0], \dots, T_k[n-1]$ of every T_k we assign an integer identifier so that $\text{Id}(P[i]) = \text{Id}(T_k[j]) \Leftrightarrow P[i] = T_k[j]$ by using the KMR algorithm ([reference]) in $\tilde{\mathcal{O}}(n^2)$.

We use the approach described in [kangaroo reference]. There exists a data structure (suffix array) which for a given one-dimensional array S allows us to detect all mismatches between any given two of its subarrays of equal length. It can be built in $\tilde{\mathcal{O}}(|S|)$ and the query time is $\tilde{\mathcal{O}}(d+1)$ where d is the number of mismatches. We construct the suffix array for the concatenation of the following arrays:

- the rows $P[i]$ for every i ,
- the rows $T[i]$ for every i ,
- the array $\text{Id}(P[0]) \text{Id}(P[1]) \dots \text{Id}(P[m-1])$,
- the arrays $\text{Id}(T_k[0]) \text{Id}(T_k[1]) \dots \text{Id}(T_k[n-1])$ for every k ,

the total length of which is $\mathcal{O}(n^2)$. Let us consider a problem of detecting mismatches between P and some $T' = T[j..j+m-1, k..k+m-1]$. We can first find all row indices i for which $P[i] \neq T'[i]$ by finding all mismatches between $\text{Id}(P[0]) \dots \text{Id}(P[m-1])$ and $\text{Id}(T_k[j]) \dots \text{Id}(T_k[j+m-1])$, which we do with query to the data structure. For every such i we can then find all mismatches between $P[i]$ and $T'[i]$ by querying $P[i]$ and $T[i+j][k..k+m-1]$. If the distance between P and T' is d , the first query takes $\tilde{\mathcal{O}}(d+1)$ operations and all subsequent queries take $\tilde{\mathcal{O}}(d+1)$ operations in total. \square

Lemma 1. HD1D with don't care symbols can be solved in $\tilde{O}(n|\Sigma|)$ by running $|\Sigma|$ instances of FFT.

Lemma 2. There exists a $(1 + \varepsilon)$ -approximate algorithm (introduced in [35]) which solves HD1D with don't care symbols in $\tilde{O}(n)$.

Theorem 4. HD2D with don't care symbols can be solved in $\tilde{O}(n^2|\Sigma|)$.

Proof. We will again use the array notation. We construct one-dimensional strings \bar{T} and \bar{P} by concatenating subsequent rows $T[0], \dots, T[n-1]$ of T and rows $P[0], \dots, P[m-1]$ of P padded with don't care symbols:

$$\begin{aligned}\bar{T} &= T[0] \ T[1] \ \dots \ T[n-1], \\ \bar{P} &= P[0] \ ?^{n-m} \ P[1] \ ?^{n-m} \ \dots \ ?^{n-m} \ P[m-1].\end{aligned}$$

We run the algorithm from Lemma 1. The distance between $T[i \dots i+m-1, j \dots j+m-1]$ and P is equal to the distance between $\bar{T}[in+j \dots in+j+nm-n+m-1]$ and \bar{P} . \square

Theorem 5. There exists a $(1 + \varepsilon)$ -approximate algorithm which solves HD2D with don't care symbols in $\tilde{O}(n^2)$.

Proof. Identical to Theorem 4, but we use the algorithm from Lemma 2 instead of Lemma 1. \square

The same reduction as in Theorem 4 can be applied for every HD1D solution which allows don't care symbols. Unfortunately, the most effective known algorithms for bounded HD1D rely on periodicity ([21], [33]) and inherently do not allow don't care symbols, thus, they cannot be easily generalized.

Observation 1 (Don't care padding). *Every HD2D solution which allows don't care symbols (eg. the algorithms from Theorem 4 and Theorem 5) can be extended to also calculate the Hamming distance for occurrences of P which are not entirely contained in T . It can be done by padding the text with don't care symbols and it does not change the complexity of the solution.*

5 Proof of Theorem 1

We show an algorithm which works in time $\tilde{O}(m^2 + mk^{5/4})$, assuming $2|n$ and $m < n \leq \frac{3}{2}m$. By Theorem 2, our main result follows.

We start by running the algorithm from Theorem 5 with $\varepsilon = 1$. We construct the set Q as the set of such vectors $q \in \mathbb{Z}^2$ for which the estimated value of $\text{Ham}(P+q, T)$ is at most $2k$. For every $q \in \{0, \dots, n-m\}^2 \setminus Q$ we say that $\text{Ham}(P+q, T)$ equals ∞ . The next step is to calculate the exact value of $\text{Ham}(P+q, T)$ for every $q \in Q$.

Let us consider the case when $|Q| \leq 2m + m^2/k$. We can run the algorithm from Theorem 3 and by the fact that $\text{Ham}(P+q, T) \leq 4k$ for every $q \in Q$, it will perform $\tilde{O}(m^2 + mk)$ operations. We are left with the case when $|Q| > 2m + m^2/k$, in which we take advantage of the fact that some strings $P+q$ for $q \in Q$ must have a large overlap and small Hamming distance from each other, and thus P must be periodic.

5.1 Two-dimensional periodicity

In this section we introduce a range of new tools related to two-dimensional periodicity. We then select some special periods of the pattern and show how to decompose it into some regularly structured monochromatic strings.

Definition 7 (Periodicity). Consider any vector $\delta \in \mathbb{Z}^2$. We say that a string S has an ℓ -period δ when

$$\text{Ham}(S + \delta, S) \leq \ell.$$

Lemma 3. For every $u, v \in Q$, the vector $u - v$ is an $8k$ -period of P .

Proof. $\text{Ham}(P+u-v, P) = \text{Ham}(P+u, P+v) \leq \text{Ham}(P+u, T) + \text{Ham}(P+v, T) \leq 4k + 4k$. \square

Theorem 6. For a given $\ell \in \mathbb{Z}^+$ and a set of points $U \subseteq [\ell + 1]^2$ such that $|U| > 4\ell$ there exist $s, t, s', t' \in U$ such that the following conditions hold for $w = t - s$ and $w' = t' - s'$:

- $0 < |w||w'| \leq 22 \frac{\ell^2}{|U|}$,
- $|\sin \alpha| \geq \frac{1}{2}$ where α is the angle between w and w' ,
- $w, w', -w, -w'$ are all contained in different quadrants.

Such w, w' can be found in $\tilde{\mathcal{O}}(|U|)$ operations.

Proof. See Section 5.4. \square

We run the algorithm from Theorem 6 on the set Q (where $\ell = n - m \leq m/2$, thus $|Q| > 2m + m^2/k \geq 4\ell$). We obtain vectors $\varphi \in \mathcal{Q}_4$ and $\psi \in \mathcal{Q}_1$ which by Lemma 3 are $\mathcal{O}(k)$ -periods of P . We will refer to those vectors throughout the rest of the description and we define $p = \varphi \times \psi$. Note that because $|Q| > 2m + m^2/k$, we have $p \leq |\varphi||\psi| = \mathcal{O}(\min\{m, k\})$.

Definition 8 (Lattice congruency). We define $\mathcal{L} = \{s\varphi + t\psi : s, t \in \mathbb{Z}\}$. We say that two vectors $u, v \in \mathbb{Z}^2$ are **lattice-congruent** and denote $u \equiv v$ when $u - v \in \mathcal{L}$ [Galil citation].

Lemma 4. There exists a set of points $\Gamma \subseteq \mathbb{Z}^2$ such that $|\Gamma| = p$ and every point $u \in \mathbb{Z}^2$ is lattice-congruent to exactly one point $\gamma \in \Gamma$.

Proof. TODO \square

Definition 9 (Parquet). We call a non-empty set $U \subseteq \mathbb{Z}^2$ a **parquet** when there exist some values $x_0, x_1, y_0, y_1, \varphi_0, \varphi_1, \psi_0, \psi_1 \in \mathbb{Z}$, which we will call its **signature**, such that

$$U = [x_0, x_1] \times [y_0, y_1] \cap \{u : u \in \mathbb{Z}^2, \varphi \times u \in [\varphi_0, \varphi_1], \psi \times u \in [\psi_0, \psi_1]\}.$$

- If additionally $x_1 - x_0 \geq |\varphi.x| + |\psi.x|$ and $y_1 - y_0 \geq |\varphi.y| + |\psi.y|$, then U is a **spacious** parquet.
- If additionally $x_0, y_0 = -\infty$ and $x_1, y_1 = +\infty$, then U is a **simple** parquet.

Note that every simple parquet is spacious.

Definition 10 (Subparquet). We call a non-empty set $V \subseteq \mathbb{Z}^2$ a **subparquet** when there exists a parquet U and a point $\gamma \in \mathbb{Z}^2$ such that

$$V = \{u : u \in U, u \equiv \gamma\}.$$

We call V a spacious/simple subparquet when there exists U which is (correspondingly) a spacious/simple parquet. We say that V is lattice-congruent to some $v \in \mathbb{Z}^2$ (denoted as $V \equiv v$) when $v \equiv \gamma$. We similarly define the lattice congruency between two subparquets.

Theorem 7 (Subparquet convolution). For a given list of signatures of simple subparquets $U_0, \dots, U_{\ell-1}$, list of signatures of subparquets $V_0, \dots, V_{\ell-1}$ and a set of vectors Q we can calculate

$$\sum_{i=0}^{\ell-1} |(U_i + q) \cap V_i|$$

for every $q \in Q$ in total time $\tilde{\mathcal{O}}(m^2 + \ell + |Q|)$, assuming that the subparquets only contain vectors of length $\mathcal{O}(m)$.

Proof. See Section 5.5. \square

Definition 11 (Parquet string). We call a string S a spacious/simple (sub-)parquet string when $\text{dom}(S)$ is a spacious/simple (sub-)parquet.

Theorem 8 (Periodic string decomposition). *A given spacious/simple parquet string R with $\mathcal{O}(k)$ -periods φ and ψ can be partitioned in time $\tilde{\mathcal{O}}(|\text{dom}(R)|)$ into $\mathcal{O}(k)$ monochromatic spacious/simple subparquet strings, correspondingly.*

Proof. See Section 5.6. \square

Since $|\varphi.x|, |\varphi.y|, |\psi.x|, |\psi.y| \leq n - m \leq m/2$, the $m \times m$ string P is a spacious parquet string and satisfies the assumptions of Theorem 8. We partition P into a set of strings \mathcal{V} . We then group the strings based on the single character they contain. Specifically, we construct $\mathcal{V}_\sigma = \{V : V \in \mathcal{V}, C(V) = \sigma\}$ for every character $\sigma \in \Sigma$ present in P .

Theorem 9. *For a given set of monochromatic simple subparquet strings \mathcal{S} we can calculate*

$$\sum_{S \in \mathcal{S}} \text{Ham}(P + q, S)$$

for every $q \in Q$ in total time $\tilde{\mathcal{O}}(m^2 + \sum_{S \in \mathcal{S}} |\mathcal{V}_{C(S)}|)$, assuming that the sets $\text{dom}(S)$ for $S \in \mathcal{S}$ are some pairwise disjoint subsets of $\text{dom}(T)$.

Proof. Let $U = \bigcup_{S \in \mathcal{S}} \text{dom}(S)$. Observe that

$$\sum_{S \in \mathcal{S}} \text{Ham}(P + q, S) = |(P + q) \cap U| - \sum_{S \in \mathcal{S}} \sum_{V \in \mathcal{V}_{C(S)}} |\text{dom}(V + q) \cap \text{dom}(S)|.$$

We can calculate $|(P + q) \cap U|$ for every $q \in Q$ with a single instance of FFT (see Theorem 4) or by using prefix sums in time $\tilde{\mathcal{O}}(m^2)$. To calculate the values

$$\sum_{S \in \mathcal{S}} \sum_{V \in \mathcal{V}_{C(S)}} |\text{dom}(V + q) \cap \text{dom}(S)|$$

we use the algorithm from Theorem 7 (where $\ell = \sum_{S \in \mathcal{S}} |\mathcal{V}_{C(S)}|$). \square

5.2 Text decomposition

Because the text is not necessarily periodic, we unfortunately cannot use the same approach as for the pattern. In this section we show how to decompose T using a similar, but more nuanced approach.

Definition 12 (Active text). We define the **active text** $T_{\mathbf{a}}$ as the restriction of T to $\bigcup_{q \in Q} \text{dom}(P) + q$. For a point $u \in \mathbb{Z}^2$ we define its **border distance** as $\min \{ |u - v| : v \in \mathbb{Z}^2 \setminus \text{dom}(T_{\mathbf{a}}) \}$, which we will denote as $\text{BorderDis}(u)$. For a set of points $U \subseteq \mathbb{Z}^2$, we define $\text{BorderDis}(U) = \max \{ \text{BorderDis}(u) : u \in U \}$ and for a string S , we define $\text{BorderDis}(S) = \text{BorderDis}(\text{dom}(S))$. Note that we consider the maximum distance, not minimum.

Observation 2. $\text{Ham}(P + q, T) = \text{Ham}(P + q, T_{\mathbf{a}})$ for every $q \in Q$.

Theorem 10 (Active text decomposition). *For a given positive integer $\ell \leq m$, we can partition the active text in time $\tilde{\mathcal{O}}(m^2 + \ell k)$ into a set of $\mathcal{O}(\ell k)$ monochromatic simple subparquet strings and a string F such that $\text{BorderDis}(F) = \mathcal{O}(m/\ell)$.*

Proof. See Section 5.7. \square

We partition $T_{\mathbf{a}}$ using the algorithm from Theorem 10 with $\ell = mk^{-3/4}$ into a set of simple subparquet strings \mathcal{S} and a string F , where $|\mathcal{S}| = \mathcal{O}(mk^{1/4})$ and $\text{BorderDis}(F) = \mathcal{O}(k^{3/4})$. For every $q \in Q$ we then have

$$\text{Ham}(P + q, T_{\mathbf{a}}) = \text{Ham}(P + q, F) + \sum_{S \in \mathcal{S}} \text{Ham}(P + q, S).$$

By Theorem 9, we can calculate $\sum_{S \in \mathcal{S}} \text{Ham}(P + q, S)$ for every $q \in Q$ in time $\tilde{\mathcal{O}}(m^2 + mk^{5/4})$, since $\sum_{S \in \mathcal{S}} |\mathcal{V}_{C(S)}| \leq |\mathcal{S}| |\mathcal{V}| = \mathcal{O}(mk^{5/4})$. In the next section we will introduce Theorem 12, which states that we can calculate $\text{Ham}(P + q, F)$ for every $q \in Q$ in total time $\tilde{\mathcal{O}}(m^2 + mk^{1/2} \text{BorderDis}(F))$. By substituting $\text{BorderDis}(F) = \mathcal{O}(k^{3/4})$, we get the complexity of $\tilde{\mathcal{O}}(m^2 + mk^{5/4})$, which ends the main proof.

5.3 Border proximity

In this section we explore the properties of strings defined only for the points close to the border. We consider any non-empty string S , such that $\text{dom}(S) \subseteq \text{dom}(T_{\mathbf{a}})$ and let $d = \text{BorderDis}(S)$. We define a partitioning of S into strings S_1, \dots, S_4 , by splitting it through the middle with a horizontal and vertical line. Specifically

- S_1 is the restriction of S to $\{n/2, \dots, n-1\} \times \{n/2, \dots, n-1\}$ (upper right quarter),
- S_2 is the restriction of S to $\{0, \dots, n/2-1\} \times \{n/2, \dots, n-1\}$ (upper left quarter),
- S_3 is the restriction of S to $\{0, \dots, n/2-1\} \times \{0, \dots, n/2-1\}$ (lower left quarter),
- S_4 is the restriction of S to $\{n/2, \dots, n-1\} \times \{0, \dots, n/2-1\}$ (lower right quarter).

We will demonstrate some characteristics of S_1 , and by symmetry, generalize them to S .

Lemma 5. *Assuming $d \leq m/4$, there does not exist $u \in S_1$ and $v \in T_{\mathbf{a}}$ such that $v.x - u.x \geq d$ and $v.y - u.y \geq d$.*

Proof. Assume the contrary. Since $u \in S_1$, we have $\text{BorderDis}(u) \leq d$, so there exists $w \in \mathbb{Z}^2 \setminus \text{dom}(T_{\mathbf{a}})$ such that $u.x - d \leq w.x \leq u.x + d$ and $u.y - d \leq w.y \leq u.y + d$. Since $v \in T_{\mathbf{a}}$, there exists $q \in Q$ such that $v \in [m]^2 + q$. We have

$$w.x \geq u.x - d \geq n/2 - m/4 \geq n - m > q.x$$

and

$$w.x \leq u.x + d \leq v.x \leq q.x + m - 1.$$

Similarly we can show that $q.y \leq w.y \leq q.y + m - 1$, and thus $w \in [m]^2 + q$. Since $[m]^2 + q \subseteq \text{dom}(T_{\mathbf{a}})$ and $w \notin T_{\mathbf{a}}$, we get a contradiction. \square

Theorem 11. *We can calculate $\text{Ham}(P + q, S)$ for every $q \in Q$ in total time $\tilde{\mathcal{O}}(m^2 + md|\Sigma|)$, where $|\Sigma|$ is the number of different characters present in both P and S .*

Proof. See Section 5.3.1. \square

Theorem 12. *We can calculate $\text{Ham}(P + q, S)$ for every $q \in Q$ in total time $\tilde{\mathcal{O}}(m^2 + mdk^{1/2})$.*

Proof. See Section 5.3.2. \square

5.3.1 Proof of Theorem 11

We base our approach on the simple method of calculating the Hamming distance by running an instance of FFT for each unique character. We will again utilize partitioning to reduce the problem to some smaller ones and then solve them naively.

Definition 13. (width & height) For a non-empty set $U \subseteq \mathbb{Z}^2$ we define its **width** as $\max\{u.x - v.x + 1 : u, v \in U\}$ and its **height** as $\max\{u.y - v.y + 1 : u, v \in U\}$. For a non-empty string R we define the width and height as the width and height of $\text{dom}(R)$.

Theorem 13. Given two non-empty strings P and T of widths w_P, w_T and heights h_P, h_T , we can calculate $\text{Ham}(P + q, T)$ for every $q \in \mathbb{Z}^2$, for which the result is non-zero, in total time $\tilde{\mathcal{O}}((|\Sigma| + 1)(w_P + w_T)(h_P + h_T))$, where $|\Sigma|$ denotes the number of different characters present in both P and T .

Proof. We can prove it by slightly generalizing Theorem 4, although following the same method, and utilizing Observation 1. \square

We will assume that $d \leq m/4$, since for $d = \Omega(m)$ we can, by Theorem 13, calculate the results in time $\tilde{\mathcal{O}}(m^2 + m^2|\Sigma|)$, which is sufficient.

Recall that $\text{Ham}(P + q, S) = \text{Ham}(P + q, S_1) + \dots + \text{Ham}(P + q, S_4)$. We will only show how to calculate $\text{Ham}(P + q, S_1)$ for every $q \in Q$, since the other cases are symmetric.

We will take advantage of the fact that the points close to the border can overlap only with a small subset of points from the pattern when considering the occurrences fully contained in the active text. Specifically, let us consider a string P_0 , defined as the restriction of P to $[m - d]^2$ and a string P_1 , defined as the restriction of P to $\text{dom}(P) \setminus \text{dom}(P_0)$. Since the strings P_0 and P_1 partition P , we have

$$\text{Ham}(P + q, S_1) = \text{Ham}(P_0 + q, S_1) + \text{Ham}(P_1 + q, S_1).$$

Lemma 6. $\text{dom}(P_0 + q) \cap \text{dom}(S_1) = \emptyset$ for every $q \in Q$.

Proof. Let us assume the contrary. Select any $q \in Q$ such that $\text{dom}(P_0 + q) \cap \text{dom}(S_1)$ contains some point u and consider the point $v = (u.x + d, u.y + d)$. Since $u \in [m - d]^2 + q$, we have $v \in [m]^2 + q \subseteq \text{dom}(T_a)$, thus the points $u \in S_1$ and $v \in T_a$ contradict Lemma 5. \square

Observation 3. P_1 can be partitioned into two strings P_2 and P_3 such that the width of P_2 and the height of P_3 are equal to d .

By Lemma 6, $\text{Ham}(P_0 + q, S_1) = 0$ for every $q \in Q$ and by Observation 3 we then have

$$\text{Ham}(P + q, S_1) = \text{Ham}(P_1 + q, S_1) = \text{Ham}(P_2 + q, S_1) + \text{Ham}(P_3 + q, S_1).$$

for some newly constructed strings P_2 and P_3 , such that the width of P_2 and the height of P_3 equals d . We calculate $\text{Ham}(P_2 + q, S_1)$ and $\text{Ham}(P_3 + q, S_1)$ for every Q independently and sum the results. We only show how to calculate $\text{Ham}(P_2 + q, S_1)$, since the other case is symmetric.

We will now partition S_1 . Consider an array of strings $U_0, \dots, U_{\lceil n/d \rceil - 1}$, where U_i is the restriction of S_1 to $\{id, \dots, id + d - 1\} \times [n] \cap \text{dom}(S_1)$. For the sake of formality (since the maximum/minimum of an empty set is undefined), let $V_0, \dots, V_{\ell-1}$ consist of all non-empty strings U_i , given in the increasing order of i . Observe that $V_0, \dots, V_{\ell-1}$ partition S_1 and their width is not greater than d .

For each $i \in [\ell]$ we find $h_i \in \mathbb{Z}^+$, which we define as the minimal number such that $(u.x, u.y + h_i) \notin T_a$ for every $u \in V_i$.

Lemma 7. The sum of all h_i is $\mathcal{O}(m)$.

Proof. Since for $\ell < 2$, the proof is trivial, we assume $\ell \geq 2$. For every i (since h_i is minimal) there exists a point $u_i \in V_i$, such that $(u_i.x, u_i.y + h_i - 1) \in T_a$. It can be shown that for all $i \geq 2$ we have

$$h_i \leq u_{i-2}.y - u_i.y + d,$$

since if that was not the case for some i , then the points u_{i-2} and $v = (u_i.x, u_i.y + h_i - 1)$ would contradict Lemma 5. We can conclude that

$$\sum_{i=0}^{\ell-1} h_i \leq h_0 + h_1 + \sum_{i=2}^{\ell-1} (u_{i-2}.y - u_i.y + d) = h_0 + h_1 + u_0.y + u_1.y - u_{\ell-2}.y - u_{\ell-1}.y + (\ell-2)d = \mathcal{O}(m).$$

□

Observation 4. For every i , the height of V_i is not greater than h_i . By Lemma 7 we have $|S_1| = \mathcal{O}(md)$ and by extension $|S| = \mathcal{O}(md)$.

For every $i \in [l]$ we construct the string L_i as the restriction of P_2 to $[m] \times [m - h_i] \cap \text{dom}(P_2)$ and the string H_i as the restriction of P_2 to $\text{dom}(P_2) \setminus \text{dom}(L_i)$. Since L_i and H_i partition V_i , we have

$$\text{Ham}(P_2 + q, S_1) = \sum_{i=0}^{\ell-1} \text{Ham}(P_2 + q, V_i) = \sum_{i=0}^{\ell-1} \text{Ham}(L_i + q, V_i) + \sum_{i=0}^{\ell-1} \text{Ham}(H_i + q, V_i).$$

Lemma 8. $\text{dom}(L_i + q) \cap \text{dom}(V_i) = \emptyset$ for every $q \in Q$ and $i \in [\ell]$.

Proof. Let us assume the contrary. Select any $q \in Q$ and $i \in [\ell]$, such that $\text{dom}(L_i + q) \cap \text{dom}(V_i)$ contains some point u and consider the point $v = (u.x, u.y + h_i)$. Since $u \in [m] \times [m - h_i] + q$, we have $v \in [m]^2 + q \subseteq \text{dom}(T_a)$, thus $v \in T_a$, which contradicts the definition of h_i . □

By Lemma 8, for every $q \in Q$ we have $\sum_{i=0}^{\ell-1} \text{Ham}(L_i + q, V_i) = 0$, thus our result is equal to $\sum_{i=0}^{\ell-1} \text{Ham}(H_i + q, V_i)$. We run the algorithm from Theorem 13 for every pair of H_i and V_i and, since both H_i and V_i have widths not greater than d and heights not greater than h_i , we obtain the total complexity of $\tilde{\mathcal{O}}(\sum_{i=0}^{\ell-1} (|\Sigma| + 1)dh_i)$, which, by Lemma 7, is $\tilde{\mathcal{O}}(m^2 + md|\Sigma|)$.

5.3.2 Proof of Theorem 12

Recall the construction of the sets \mathcal{V}_σ described in Section 5.1. We define $\sigma \in \Sigma$ to be a **frequent** character if $|\mathcal{V}_\sigma| \geq \sqrt{k}$ and if $|\mathcal{V}_\sigma| < \sqrt{k}$, we call it an **infrequent** character.

Observation 5. The number of different frequent characters is $\mathcal{O}(\sqrt{k})$.

We partition S into two strings F and I , based on character frequency, so that F consists of only the frequent characters and I consists of only the infrequent ones. For every $q \in Q$ we then have

$$\text{Ham}(P + q, S) = \text{Ham}(P + q, F) + \text{Ham}(P + q, I).$$

By Observation 5 and Theorem 11, we can calculate $\text{Ham}(P + q, F)$ for every $q \in Q$ in total time $\tilde{\mathcal{O}}(m^2 + mdk^{1/2})$, since $\text{BorderDis}(F) \leq d$.

We partition I into $|\text{dom}(I)|$ strings, one per every $u \in I$. Specifically, let I_u be the restriction of I to $\{u\}$ for every $u \in I$. We have $\text{Ham}(P + q, I) = \sum_{u \in I} \text{Ham}(P + q, I_u)$ for every $q \in Q$. By Definition 10, I_u are simple subarquet strings, and thus, we can by Theorem 9 calculate the results in $\tilde{\mathcal{O}}(m^2 + \sum_{u \in I} |\mathcal{V}_{I(u)}|)$. Since $I(u)$ is an infrequent character for every $u \in I$, we have $|\mathcal{V}_{I(u)}| < k^{1/2}$ for every $u \in I$. By Observation 4 we have $|\text{dom}(I)| = \mathcal{O}(md)$, and thus the total complexity is $\tilde{\mathcal{O}}(m^2 + mdk^{1/2})$.

5.4 Proof of Theorem 6

First, we find any closest pair of vectors $s, t \in U$ by running the standard $\tilde{O}(|U|)$ time algorithm and denote $w = t - s$. We define a partial order \leq_w where $v \leq_w u$ for some $u, v \in U$ when at least one condition holds:

- (a) $u = v$,
- (b) $u - v$ and w belong to the same quadrant,
- (c) $\alpha \in (-\frac{\pi}{6}, \frac{\pi}{6})$ where α is the angle between w and $u - v$.

We find the longest chain C and the longest antichain A using dynamic programming in $\tilde{O}(|U|)$ operations. We then find any closest pair of vectors $s', t' \in A$ and denote $w' = t' - s'$. We have the following inequalities:

- (i) $|U| \leq |C||A|$ (by Dilworth's theorem),
- (ii) $(|C| - 1)|w| \leq (1 + \sqrt{3})\ell$ (roughly by the fact that vectors in C must be increasing in a certain direction),
- (iii) $(|A| - 1)|w'| \leq 2\ell$ (by using a similar argument for vectors in A).

By considering the assumption $|U| > 4\ell$ it can be proven that $|w||w'| \leq 22\frac{\ell^2}{|U|}$ and the other conditions also hold.

5.5 Proof of Theorem 7

Throughout this section we will denote $D = \{u : u \in \mathcal{L}, \varphi \times u \geq 0, \psi \times u \geq 0\}$, where \mathcal{L} is the set introduced in Definition 8.

Lemma 9. *Given a set of subparquets \mathcal{V} and a set of points Q , we can calculate*

$$\sum_{V \in \mathcal{V}} |(D + q) \cap V|$$

for every $q \in Q$ in total time $\tilde{O}(n^2 + |Q| + |\mathcal{V}|)$, assuming that every $V \in \mathcal{V}$ consists of vectors of length $\mathcal{O}(n)$.

Proof. For every $u \in \mathbb{Z}^2$ let us define $\text{score}(u) = |\{V : V \in \mathcal{V}, u \in V\}|$. Observe that

$$\sum_{V \in \mathcal{V}} |(D + q) \cap V| = \sum_{u \in D+q} \text{score}(u).$$

We start by explicitly calculating the scores. We find the maximum length of a vector that some $V \in \mathcal{V}$ is defined for, which we denote ℓ . We construct the set $U \subseteq \mathbb{Z}^2$ of all vectors of length at most ℓ . By the assumption, we have $\ell = \mathcal{O}(n)$, and thus $|U| = \mathcal{O}(\ell^2) = \mathcal{O}(n^2)$. We observe that since all the scores are zero for points outside of U , we can only calculate them for $u \in U$.

We find the set Γ introduced in Lemma 4 and for every $\gamma \in \Gamma$ we construct $U_\gamma = U \cap (\mathcal{L} + \gamma)$. Consider any $u \in U_\gamma$ for some fixed $\gamma \in \Gamma$ and any $V \in \mathcal{V}$. We observe that if $V \not\equiv \gamma$, then $u \notin V$ and thus V does not contribute to $\text{score}(u)$. If $V \equiv \gamma$, then we can find a parquet W such that $V = W \cap (\mathcal{L} + \gamma)$ and we have $u \in V \Leftrightarrow u \in W \cap (\mathcal{L} + \gamma) \Leftrightarrow u \in W$. Thus, if we denote \mathcal{W}_γ as the set of parquets W obtained for every $V \in \mathcal{V}$ such that $V \equiv \gamma$, then $\text{score}(u)$ for $u \in U_\gamma$ is the number of parquets $W \in \mathcal{W}_\gamma$ such that $u \in W$. We calculate $\text{score}(u)$ for every $u \in U_\gamma$ by sweeping U_γ and \mathcal{W}_γ in time $\tilde{O}(|U_\gamma| + |\mathcal{W}_\gamma|)$. We do it independently for every $\gamma \in \Gamma$, performing $\tilde{O}(|U| + |\mathcal{V}|) = \tilde{O}(n^2 + |\mathcal{V}|)$ operations in total.

Now consider a query vector $q \in Q$. Let $\gamma \in \Gamma$ be such that $q \equiv \gamma$. We have already showed that the sum of scores for $u \in D + q$ is equal to the sum of scores for $u \in (D + q) \cap U$. Since $(D + q) \cap U = (D + q) \cap U_\gamma$, we see that the result is the sum of scores for such $u \in U_\gamma$, for which $\varphi \times u \geq \varphi \times q$ and $\psi \times u \geq \psi \times q$. If we denote $Q_\gamma = Q \cap (\mathcal{L} + \gamma)$, we see that we can calculate the results for all $q \in Q_\gamma$ by sweeping Q_γ and U_γ in time $\tilde{O}(|Q_\gamma| + |U_\gamma|)$. We do it independently for every $\gamma \in \Gamma$, performing $\tilde{O}(|Q| + |U|) = \tilde{O}(n^2 + |Q|)$ operations in total. \square

Lemma 10. *For any simple subparquet U we can find $w_0, \dots, w_3 \in \mathbb{Z}^2$ such that*

$$|U \cap X| = \sum_{j=0}^3 (-1)^j |(D + w_j) \cap X|$$

for every $X \subseteq \mathbb{Z}^2$. If U consists of vectors of length $\mathcal{O}(n)$, then w_0, \dots, w_3 are of length $\mathcal{O}(n)$.

Proof. TODO \square

We apply Lemma 10 to every U_i and find $w_{i,0}, \dots, w_{i,3}$ so that we have

$$\begin{aligned} \sum_{i=0}^{\ell-1} |(U_i + q) \cap V_i| &= \sum_{i=0}^{\ell-1} |U_i \cap (V_i - q)| = \sum_{i=0}^{\ell-1} \sum_{j=0}^3 (-1)^j |(D + w_{i,j}) \cap (V_i - q)| = \\ &= \sum_{j=0}^3 (-1)^j \sum_{i=0}^{\ell-1} |(D + q) \cap (V_i - w_{i,j})|. \end{aligned}$$

By Lemma 9 we can independently calculate the values $\sum_{i=0}^{\ell-1} |(D + q) \cap (V_i - w_{i,j})|$ for every j by running the algorithm for $\mathcal{V}_j = \{V_i - w_{i,j} : i \in [\ell]\}$ and Q .

5.6 Proof of Theorem 8

Definition 14 (Lattice graph). For a set $U \subseteq \mathbb{Z}^2$ we define its **lattice graph** $G_U = (U, E_U)$ where

$$E_U = \{ \{u, u + \delta\} : \delta \in \{\varphi, \psi\}, u \in U, u + \delta \in U \}$$

so every vector is connected with its translations by $\varphi, \psi, -\varphi, -\psi$.

Lemma 11. *If U is a spacious subparquet, then G_U is connected.*

Firstly, we partition R into a set of subparquet strings \mathcal{S} . For every non-empty $S \in \mathcal{S}$ we consider a lattice graph $G_{\text{dom}(S)}$. If S is not monochromatic, then since $G_{\text{dom}(S)}$ is connected, there must exist a pair of neighboring vectors v, w such that $S(v) \neq S(w)$. We select any such pair and partition S into spacious (or simple if S is simple) subparquet strings S' and S'' such that $v \in S'$ and $w \in S''$. For example if $v = w + \varphi$, then $S' = \{u : u \in S, \psi \times u \leq \psi \times v\}$ and $S'' = \{u : u \in S, \psi \times u > \psi \times v\}$. In the cases when $v = w + \delta$ for $\delta \in \{-\varphi, \psi, -\psi\}$ the construction is similar.

We can recursively partition S' and S'' further until we obtain monochromatic strings. Because R has $\mathcal{O}(k)$ -periods φ and ψ , the total number of neighbor pairs v, w such that $S(v) \neq S(w)$ is $\mathcal{O}(k)$ throughout all $S \in \mathcal{S}$. Thus the total number of recursive calls is $\mathcal{O}(k)$ and because $|\mathcal{S}| = \mathcal{O}(k)$, the total number of constructed strings is $\mathcal{O}(k)$. The algorithm can be implemented to work in time $\tilde{O}(|\text{dom}(R)|)$.

5.7 Proof of Theorem 10

In this section we will consider an empty set to be a valid parallelogram.

Definition 15. For a set of points $U \subseteq \mathbb{R}^2$ we will denote

$$X(U) = \{u.x : u \in U\}, \quad Y(U) = \{u.y : u \in U\}.$$

Observation 6. For any given $\ell \in \mathbb{Z}^+$ and $v \in \mathbb{Z}^2$ we can find an array of parallel lines f_0, f_1, \dots, f_ℓ , where $f_i = \{u : u \in \mathbb{R}^2, v \times u = c_i\}$ for some $c_i \in \mathbb{R} \setminus \mathbb{Q}$, such that

- $c_0 < v \times u < c_\ell$ for every $u \in [n]^2$, or namely, the set $[n]^2$ is between f_0 and f_ℓ ,
- $0 < c_{i+1} - c_i = \mathcal{O}(n|v|/\ell)$ for every $i \in [\ell]$, or namely, the distance between every two consecutive lines is $\mathcal{O}(n/\ell)$.

We use Observation 6 with $v = \varphi$ to construct the lines h_0, \dots, h_ℓ and with $v = \psi$ to construct the lines s_0, \dots, s_ℓ . For every $i, j \in [\ell + 1]$ we construct a point $w_{i,j}$ as the intersection of h_i and s_j (it is easy to see that since φ and ψ are not collinear, h_i and s_j are not parallel). For every $i, j \in [\ell]$ we construct a parallelogram $p_{i,j}$ defined as the area between s_i and s_{i+1} intersected with the area between h_j and h_{j+1} . Specifically,

$$p_{i,j} = \{u : u \in \mathbb{R}^2, \varphi \times u \in [\varphi \times w_{i,j}, \varphi \times w_{i+1,j+1}], \psi \times u \in [\psi \times w_{i,j}, \psi \times w_{i+1,j+1}]\}.$$

For a better reference, the vertices of $p_{i,j}$ are $w_{i,j}, w_{i+1,j}, w_{i+1,j+1}, w_{i,j+1}$. Observe that every $u \in [n]^2$ is contained strictly inside exactly one parallelogram $p_{i,j}$.

Lemma 12. For every $i \in [\ell - 1]$ and $j \in [\ell]$ we have

$$\begin{aligned} \min X(p_{i,j}) &< \min X(p_{i+1,j}), & \min Y(p_{i,j}) &\leq \min Y(p_{i+1,j}), \\ \max X(p_{i,j}) &< \max X(p_{i+1,j}), & \max Y(p_{i,j}) &\leq \max Y(p_{i+1,j}) \end{aligned}$$

and for every $i \in [\ell]$ and $j \in [\ell - 1]$ we have

$$\begin{aligned} \min X(p_{i,j}) &\geq \min X(p_{i,j+1}), & \min Y(p_{i,j}) &< \min Y(p_{i,j+1}), \\ \max X(p_{i,j}) &\geq \max X(p_{i,j+1}), & \max Y(p_{i,j}) &< \max Y(p_{i,j+1}). \end{aligned}$$

Proof. It follows from the fact that we selected $\varphi \in [0, +\infty) \times (-\infty, 0)$ and $\psi \in (0, +\infty) \times [0, +\infty)$. For example, to prove the first inequality, we can consider a point $u \in p_{i+1,j}$, such that $u.x = \min X(p_{i+1,j})$ and then construct a point $v \in p_{i,j}$, such that $v = u - t\psi$ for some $t > 0$, and thus $\min X(p_{i,j}) \leq v.x \leq u.x = \min X(p_{i+1,j})$. The other inequalities can be proven analogously. \square

Lemma 13. For every $i, j \in [\ell]$ and every $u, v \in X(p_{i,j}) \times Y(p_{i,j})$, we have $|u - v| = \mathcal{O}(n/\ell)$.

Proof. See Section 5.7.2. \square

Consider the case when $\max X(p_{i,j}) - \min X(p_{i,j}) \geq m/4$ for some $i, j \in [\ell]$. By Lemma 13, we would have $m/4 \leq \max X(p_{i,j}) - \min X(p_{i,j}) = \mathcal{O}(n/\ell)$, and thus $\ell = \mathcal{O}(1)$. In that case we can return a trivial partitioning where $F = T_{\mathbf{a}}$ and the set of monochromatic strings is empty, since $\text{BorderDis}(T_{\mathbf{a}}) = \mathcal{O}(m)$. We can use the same argument if we have $\max Y(p_{i,j}) - \min Y(p_{i,j}) \geq m/4$ for some $i, j \in [\ell]$. Thus, from now on we will assume that $\max X_{i,j} - \min X_{i,j} < m/4$ and $\max Y_{i,j} - \min Y_{i,j} < m/4$ for every $i, j \in [\ell]$.

Let $z = \frac{n-1}{2}$. We split the plain with two lines $x = z$ and $y = z$ into four quarters:

- 1) $K_1 = (z, +\infty) \times (z, +\infty)$,

- 2) $K_2 = (-\infty, z) \times (z, +\infty)$,
- 3) $K_3 = (-\infty, z) \times (-\infty, z)$,
- 4) $K_4 = (z, +\infty) \times (-\infty, z)$.

Let us denote by \mathcal{I} the set of all parallelograms $p_{i,j}$, such that they intersect with the line $x = z$ or with the line $y = z$ (or both). Observe that every parallelogram $p_{i,j} \notin \mathcal{I}$ must be fully contained in one of the quarters, meaning $p_{i,j} \subseteq K_d$ for some $d \in \{1, \dots, 4\}$.

Lemma 14. $|\mathcal{I}| = \mathcal{O}(\ell)$.

Proof. Consider the line $x = z$, denoted f . It intersects with every line h_0, \dots, h_ℓ at most once (and does not overlap with any of them). Similarly, it intersects with every line s_0, \dots, s_ℓ at most once. Denote the set of such intersections as U . For every parallelogram $p \in \mathcal{I}$, there must exist $u \in U$, such that $u \in p$. For every $u \in U$, there are at most four parallelograms $p \in \mathcal{I}$, such that $u \in p$, thus $|\mathcal{I}| \leq 4|U| = \mathcal{O}(\ell)$. \square

Now consider any $j \in [\ell]$. By Lemma 12, we can find $s, t \in [\ell + 1]$, such that the array $p_{0,j}, \dots, p_{\ell-1,j}$ is split into three groups:

- a) $p_{0,j}, \dots, p_{s-1,j}$, which includes only parallelograms fully contained in K_3 ,
- b) $p_{s,j}, \dots, p_{t-1,j}$, which does not include any parallelogram fully contained in K_1 or K_3 ,
- c) $p_{t,j}, \dots, p_{\ell-1,j}$, which includes only parallelograms fully contained in K_1 .

We now "merge together" the parallelograms from group (a) and from group (c). Specifically, we construct

$$g_j^3 = \bigcup_{i=0}^{s-1} p_{i,j}, \quad g_j^1 = \bigcup_{i=t}^{\ell-1} p_{i,j}.$$

We do it for every $j \in [\ell]$. Observe that the sets $g_0^1, \dots, g_{\ell-1}^1$ are parallelograms (possibly empty) and that they cover the same area as all the fully contained in K_1 parallelograms $p_{i,j}$. The same is true for $g_0^3, \dots, g_{\ell-1}^3$ and the parallelograms in K_3 .

Now for every $i \in [\ell]$ we similarly find $s, t \in [\ell + 1]$, such that the array $p_{i,0}, \dots, p_{i,\ell-1}$ is split into three groups:

- a) $p_{i,0}, \dots, p_{i,s-1}$, which includes only parallelograms fully contained in K_4 ,
- b) $p_{i,s}, \dots, p_{i,t-1}$, which does not include any parallelogram fully contained in K_2 or K_4 ,
- c) $p_{i,t}, \dots, p_{i,\ell-1}$, which includes only parallelograms fully contained in K_2 ,

and then construct

$$g_i^4 = \bigcup_{j=0}^{s-1} p_{i,j}, \quad g_i^2 = \bigcup_{j=t}^{\ell-1} p_{i,j}.$$

We denote

$$\mathcal{G} = \{g_0^1, \dots, g_{\ell-1}^1\} \cup \{g_0^2, \dots, g_{\ell-1}^2\} \cup \{g_0^3, \dots, g_{\ell-1}^3\} \cup \{g_0^4, \dots, g_{\ell-1}^4\}.$$

Again observe that for every $u \in [n]^2$ there exists exactly one parallelogram $p \in \mathcal{G} \cup \mathcal{I}$, such that $u \in p$, and since the sides of p do not contain integer points, u lays strictly inside p .

Definition 16 (Coverability). We say that a set $U \subseteq \mathbb{Z}^2$ is **coverable** if $U \subseteq \text{dom}(P + q)$ for some $q \in \mathcal{Q}$.

Lemma 15. *For every $p \in \mathcal{I}$, the set $p \cap \mathbb{Z}^2$ is coverable.*

FINISH. □

Lemma 16. *For every $g \in \mathcal{G}$ we can construct two parallelograms c and b , such that*

- $c \cap \mathbb{Z}^2$ is coverable,
- $b \cap \mathbb{Z}^2$ is a $\mathcal{O}(n/\ell)$ -periphery,
- $g \cap \mathbb{Z}^2$ is partitioned into $b \cap \mathbb{Z}^2$ and $c \cap \mathbb{Z}^2$.

Proof. See Section 5.7.1. □

Lemma 17. *The restriction of T to a coverable set has $\mathcal{O}(k)$ -periods φ and ψ .*

Proof. Let R denote the restriction. For $q \in Q$, such that $\text{dom}(R) \subseteq \text{dom}(P + q)$ we have

$$\begin{aligned} \text{Ham}(R + \varphi, R) &\leq \text{Ham}(R + \varphi, P + q + \varphi) + \text{Ham}(P + q + \varphi, P + q) + \text{Ham}(P + q, R) \leq \\ &\leq \text{Ham}(T, P + q) + \text{Ham}(P + \varphi, P) + \text{Ham}(P + q, T) = \mathcal{O}(k) \end{aligned}$$

and identically $\text{Ham}(R + \psi, R) = \mathcal{O}(k)$. □

[FINISH]

5.7.1 Proof of Lemma 16

For an empty parallelogram, the proof is trivial. Consider a non-empty set g_j^1 for some $j \in [\ell]$. We denote $B = K_1 \cap \mathbb{Z}^2 \setminus \text{dom}(T_{\mathbf{a}})$.

Lemma 18. *There does not exist a pair of points $u \in T_{\mathbf{a}}$ and $v \in B$, such that $u.x \geq v.x$ and $u.y \geq v.y$.*

Proof. □

Lemma 19. *Every set $U \subseteq K_1 \cap \mathbb{Z}^2$, such that $(\max X(U), \max Y(U)) \in T_{\mathbf{a}}$ is coverable.*

Proof. □

Consider a non-empty set g_j^1 for some $j \in [\ell]$. Recall that there exists $t \in [\ell + 1]$, such that $g_j^1 = \bigcup_{i=t}^{\ell-1} p_{i,j}$. We find

$$f = \min \{ i : i \in \{ t, \dots, \ell - 1 \}, (\lfloor \max X(p_{i,j}) \rfloor, \lfloor \max Y(p_{i,j}) \rfloor) \in B \}.$$

If the minimum does not exist, we set $f = \ell$. We then construct the parallelograms

$$c = \bigcup_{i=0}^{f-1} p_{i,j}, \quad b = \bigcup_{i=f}^{\ell-1} p_{i,j}.$$

We now show that the set $c \cap \mathbb{Z}^2$ is coverable. If $c \cap \mathbb{Z}^2$ is empty, then it is coverable, so let us assume it is not. In that case $f > 0$. It is clear that $c \cap \mathbb{Z}^2 \subseteq K_1$. By Lemma 12, we have

$$\begin{aligned} \max X(c) &= \max X(p_{f-1,j}), \\ \max Y(c) &= \max Y(p_{f-1,j}), \end{aligned}$$

and thus

$$\begin{aligned}\max X(c \cap \mathbb{Z}^2) &\leq \lfloor \max X(c) \rfloor = \lfloor \max X(p_{f-1,j}) \rfloor, \\ \max Y(c \cap \mathbb{Z}^2) &\leq \lfloor \max Y(c) \rfloor = \lfloor \max Y(p_{f-1,j}) \rfloor.\end{aligned}$$

We see that $(\max X(c \cap \mathbb{Z}^2), \max Y(c \cap \mathbb{Z}^2)) \in T_{\mathbf{a}}$, since it would otherwise contradict Lemma 18, considering that $(\lfloor \max X(p_{f-1,j}) \rfloor, \lfloor \max Y(p_{f-1,j}) \rfloor) \in T_{\mathbf{a}}$. We see that $c \cap \mathbb{Z}^2$ satisfies the conditions of Lemma 19, thus $c \cap \mathbb{Z}^2$ is coverable.

We now show that the set $b \cap \mathbb{Z}^2$ is a $\mathcal{O}(n/\ell)$ -periphery. If it is empty, then the proof is trivial, so let us assume it is not. In that case $f < \ell$. Denote $v = (\lfloor \max X(p_{f,j}) \rfloor, \lfloor \max Y(p_{f,j}) \rfloor)$. By definition, $v \in B$. Consider any point $u \in b \cap \mathbb{Z}^2$. There exists exactly one $i \in \{f, \dots, \ell-1\}$, such that u lays strictly inside $p_{i,j}$. [FINISH]

5.7.2 Proof of Lemma 13

Consider a parallelogram $p_{i,j}$ for some $i, j \in [\ell]$.

Lemma 20. *For every $u, v \in p_{i,j}$, we have $|u - v| = \mathcal{O}(n/\ell)$.*

Proof. Consider any $u, v \in p_{i,j}$ and denote $w = u - v$. By definition of $p_{i,j}$, we have

$$|\varphi \times w| = |\varphi \times (u - v)| = |\varphi \times u - \varphi \times v| = \mathcal{O}(n|\varphi|/\ell)$$

and similarly $|\psi \times w| = \mathcal{O}(n|\psi|/\ell)$. Since φ and ψ are not collinear, there exist $s, t \in \mathbb{R}$, such that $w = s\varphi + t\psi$. Recall that by Theorem 6 we have $|\varphi \times \psi| \geq \frac{1}{2}|\varphi||\psi|$ (since $|\sin \alpha| \geq 1/2$), thus

$$\frac{1}{2}|t||\varphi||\psi| \leq |t||\varphi \times \psi| = |\varphi \times (s\varphi + t\psi)| = |\varphi \times w| = \mathcal{O}(n|\varphi|/\ell),$$

which gives us $|t\varphi| = \mathcal{O}(n/\ell)$. We can similarly prove that $|s\psi| = \mathcal{O}(n/\ell)$ and finally

$$|w| = |s\varphi + t\psi| \leq |s\varphi| + |t\psi| = \mathcal{O}(n/\ell).$$

□

Lemma 21. *For every point $u \in X(p_{i,j}) \times Y(p_{i,j})$ there exists a point $v \in p_{i,j}$, such that $|u - v| = \mathcal{O}(n/\ell)$.*

Proof. There exists a point $w \in p_{i,j}$, such that $u.x = w.x$, and a point $v \in p_{i,j}$, such that $u.y = v.y$. By Lemma 20, we have

$$|u - v| = |u.x - v.x| = |w.x - v.x| \leq |w - v| = \mathcal{O}(n/\ell).$$

□

Consider any $u, v \in p_{i,j}$. By Lemma 21, there exist $u', v' \in p_{i,j}$ such that $|u - u'| = \mathcal{O}(n/\ell)$ and $|v - v'| = \mathcal{O}(n/\ell)$. By Lemma 20, we have

$$|u - v| \leq |u - u'| + |u' - v'| + |v' - v| = \mathcal{O}(n/\ell).$$

References

- [1] Amihoud Amir and Gary Benson. Two-dimensional periodicity and its applications. In Greg N. Frederickson, editor, *Proceedings of the Third Annual ACM/SIGACT-SIAM Symposium on Discrete Algorithms, 27-29 January 1992, Orlando, Florida, USA*, pages 440–452. ACM/SIAM, 1992.
- [2] Amihoud Amir and Gary Benson. Two-dimensional periodicity in rectangular arrays. *SIAM J. Comput.*, 27(1):90–106, 1998.
- [3] Amihoud Amir, Gary Benson, and Martin Farach. An alphabet independent approach to two-dimensional pattern matching. *SIAM J. Comput.*, 23(2):313–323, 1994.
- [4] Amihoud Amir, Gary Benson, and Martin Farach. Optimal two-dimensional compressed matching. *J. Algorithms*, 24(2):354–379, 1997.
- [5] Amihoud Amir and Martin Farach. Efficient 2-dimensional approximate matching of half-rectangular figures. *Inf. Comput.*, 118(1):1–11, 1995.
- [6] Amihoud Amir and Gad M. Landau. Fast parallel and serial multidimensional approximate array matching. *Theor. Comput. Sci.*, 81(1):97–115, 1991.
- [7] Amihoud Amir, Gad M. Landau, Shoshana Marcus, and Dina Sokol. Two-dimensional maximal repetitions. *Theor. Comput. Sci.*, 812:49–61, 2020.
- [8] Amihoud Amir, Gad M. Landau, and Dina Sokol. Inplace 2d matching in compressed images. *J. Algorithms*, 49(2):240–261, 2003.
- [9] Amihoud Amir, Gad M. Landau, and Dina Sokol. Inplace run-length 2d compressed search. *Theor. Comput. Sci.*, 290(3):1361–1383, 2003.
- [10] Alberto Apostolico and Valentin E. Brimkov. Fibonacci arrays and their two-dimensional repetitions. *Theor. Comput. Sci.*, 237(1-2):263–273, 2000.
- [11] Alberto Apostolico and Valentin E. Brimkov. Optimal discovery of repetitions in 2d. *Discret. Appl. Math.*, 151(1-3):5–20, 2005.
- [12] Alberto Apostolico, Laxmi Parida, and Simona E. Rombo. Motif patterns in 2d. *Theor. Comput. Sci.*, 390(1):40–55, 2008.
- [13] Ricardo A. Baeza-Yates and Gonzalo Navarro. Fast two-dimensional approximate pattern matching. In Claudio L. Lucchesi and Arnaldo V. Moura, editors, *LATIN '98: Theoretical Informatics, Third Latin American Symposium, Campinas, Brazil, April, 20-24, 1998, Proceedings*, volume 1380 of *Lecture Notes in Computer Science*, pages 341–351. Springer, 1998.
- [14] Ricardo A. Baeza-Yates and Mireille Régnier. Fast two-dimensional pattern matching. *Inf. Process. Lett.*, 45(1):51–57, 1993.
- [15] Theodore P. Baker. A technique for extending rapid exact-match string matching to arrays of more than one dimension. *SIAM J. Comput.*, 7(4):533–541, 1978.
- [16] Piotr Berman, Marek Karpinski, Lawrence L. Larmore, Wojciech Plandowski, and Wojciech Rytter. On the complexity of pattern matching for highly compressed two-dimensional texts. *J. Comput. Syst. Sci.*, 65(2):332–350, 2002.
- [17] Richard S. Bird. Two dimensional pattern matching. *Inf. Process. Lett.*, 6(5):168–170, 1977.

- [18] Itai Boneh, Dvir Fried, Shay Golan, Matan Kraus, Adrian Miclaus, and Arseny Shur. Searching 2d-strings for matching frames. *CoRR*, abs/2310.02670, 2023.
- [19] Panagiotis Charalampopoulos, Jakub Radoszewski, Wojciech Rytter, Tomasz Walen, and Wiktor Zuba. The number of repetitions in 2d-strings. In Fabrizio Grandoni, Grzegorz Herman, and Peter Sanders, editors, *28th Annual European Symposium on Algorithms, ESA 2020, September 7-9, 2020, Pisa, Italy (Virtual Conference)*, volume 173 of *LIPIcs*, pages 32:1–32:18. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2020.
- [20] Panagiotis Charalampopoulos, Jakub Radoszewski, Wojciech Rytter, Tomasz Walen, and Wiktor Zuba. Computing covers of 2d-strings. In Pawel Gawrychowski and Tatiana Starikovskaya, editors, *32nd Annual Symposium on Combinatorial Pattern Matching, CPM 2021, July 5-7, 2021, Wrocław, Poland*, volume 191 of *LIPIcs*, pages 12:1–12:20. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2021.
- [21] Raphaël Clifford, Allyx Fontaine, Ely Porat, Benjamin Sach, and Tatiana Starikovskaya. The k-mismatch problem revisited. *CoRR*, abs/1508.00731, 2015.
- [22] Raphaël Clifford, Allyx Fontaine, Tatiana Starikovskaya, and Hjalte Wedel Vildhøj. Dynamic and approximate pattern matching in 2d. In Shunsuke Inenaga, Kunihiko Sadakane, and Tetsuya Sakai, editors, *String Processing and Information Retrieval - 23rd International Symposium, SPIRE 2016, Beppu, Japan, October 18-20, 2016, Proceedings*, volume 9954 of *Lecture Notes in Computer Science*, pages 133–144, 2016.
- [23] Richard Cole, Maxime Crochemore, Zvi Galil, Leszek Gasieniec, Ramesh Hariharan, S. Muthukrishnan, Kunsoo Park, and Wojciech Rytter. Optimally fast parallel algorithms for preprocessing and pattern matching in one and two dimensions. In *34th Annual Symposium on Foundations of Computer Science, Palo Alto, California, USA, 3-5 November 1993*, pages 248–258. IEEE Computer Society, 1993.
- [24] Richard Cole, Zvi Galil, Ramesh Hariharan, S. Muthukrishnan, and Kunsoo Park. Parallel two dimensional witness computation. *Inf. Comput.*, 188(1):20–67, 2004.
- [25] Maxime Crochemore, Leszek Gasieniec, Ramesh Hariharan, S. Muthukrishnan, and Wojciech Rytter. A constant time optimal parallel algorithm for two-dimensional pattern matching. *SIAM J. Comput.*, 27(3):668–681, 1998.
- [26] Maxime Crochemore, Leszek Gasieniec, Wojciech Plandowski, and Wojciech Rytter. Two-dimensional pattern matching in linear time and small space. In Ernst W. Mayr and Claude Puech, editors, *STACS 95, 12th Annual Symposium on Theoretical Aspects of Computer Science, Munich, Germany, March 2-4, 1995, Proceedings*, volume 900 of *Lecture Notes in Computer Science*, pages 181–192. Springer, 1995.
- [27] Maxime Crochemore, Costas S. Iliopoulos, and Maureen Korda. Two-dimensional prefix string matching and covering on square matrices. *Algorithmica*, 20(4):353–373, 1998.
- [28] Kimmo Fredriksson, Gonzalo Navarro, and Esko Ukkonen. Optimal exact and fast approximate two dimensional pattern matching allowing rotations. In Alberto Apostolico and Masayuki Takeda, editors, *Combinatorial Pattern Matching, 13th Annual Symposium, CPM 2002, Fukuoka, Japan, July 3-5, 2002, Proceedings*, volume 2373 of *Lecture Notes in Computer Science*, pages 235–248. Springer, 2002.
- [29] Zvi Galil and Kunsoo Park. Alphabet-independent two-dimensional witness computation. *SIAM J. Comput.*, 25(5):907–935, 1996.

- [30] Guilhem Gamard and Gwénaél Richomme. Coverability and multi-scale coverability on infinite pictures. *J. Comput. Syst. Sci.*, 104:258–277, 2019.
- [31] Guilhem Gamard, Gwénaél Richomme, Jeffrey O. Shallit, and Taylor J. Smith. Periodicity in rectangular arrays. *Inf. Process. Lett.*, 118:58–63, 2017.
- [32] Pawel Gawrychowski, Samah Ghazawi, and Gad M. Landau. Lower bounds for the number of repetitions in 2d strings. In Thierry Lecroq and Hélène Touzet, editors, *String Processing and Information Retrieval - 28th International Symposium, SPIRE 2021, Lille, France, October 4-6, 2021, Proceedings*, volume 12944 of *Lecture Notes in Computer Science*, pages 179–192. Springer, 2021.
- [33] Pawel Gawrychowski and Przemyslaw Uznanski. Optimal trade-offs for pattern matching with k mismatches. *CoRR*, abs/1704.01311, 2017.
- [34] Juha Kärkkäinen and Esko Ukkonen. Two- and higher-dimensional pattern matching in optimal expected time. *SIAM J. Comput.*, 29(2):571–589, 1999.
- [35] Howard J. Karloff. Fast algorithms for approximately counting mismatches. *Inf. Process. Lett.*, 48(2):53–60, 1993.
- [36] Richard M. Karp, Raymond E. Miller, and Arnold L. Rosenberg. Rapid identification of repeated patterns in strings, trees and arrays. In Patrick C. Fischer, H. Paul Zeiger, Jeffrey D. Ullman, and Arnold L. Rosenberg, editors, *Proceedings of the 4th Annual ACM Symposium on Theory of Computing, May 1-3, 1972, Denver, Colorado, USA*, pages 125–136. ACM, 1972.
- [37] Richard M. Karp and Michael O. Rabin. Efficient randomized pattern-matching algorithms. *IBM J. Res. Dev.*, 31(2):249–260, 1987.
- [38] Donald E. Knuth, James H. Morris Jr., and Vaughan R. Pratt. Fast pattern matching in strings. *SIAM J. Comput.*, 6(2):323–350, 1977.
- [39] Kamala Krithivasan and R. Sitalakshmi. Efficient two-dimensional pattern matching in the presence of errors. *Inf. Sci.*, 43(3):169–184, 1987.
- [40] Kalpana Mahalingam and Palak Pandoh. On the maximum number of distinct palindromic sub-arrays. In Carlos Martín-Vide, Alexander Okhotin, and Dana Shapira, editors, *Language and Automata Theory and Applications - 13th International Conference, LATA 2019, St. Petersburg, Russia, March 26-29, 2019, Proceedings*, volume 11417 of *Lecture Notes in Computer Science*, pages 434–446. Springer, 2019.
- [41] Shoshana Marcus and Dina Sokol. 2d lyndon words and applications. *Algorithmica*, 77(1):116–133, 2017.
- [42] Filippo Mignosi, Antonio Restivo, and Pedro V. Silva. On fine and wilf’s theorem for bidimensional words. *Theor. Comput. Sci.*, 292(1):245–262, 2003.
- [43] Shoshana Neuburger and Dina Sokol. Succinct 2d dictionary matching. *Algorithmica*, 65(3):662–684, 2013.
- [44] Kunsoo Park. Analysis of two-dimensional approximate pattern matching algorithms. *Theor. Comput. Sci.*, 201(1-2):263–273, 1998.
- [45] Jakub Radoszewski, Wojciech Rytter, Juliusz Straszynski, Tomasz Walen, and Wiktor Zuba. Rectangular tile covers of 2d-strings. In Hideo Bannai and Jan Holub, editors, *33rd Annual Symposium on Combinatorial Pattern Matching, CPM 2022, June 27-29, 2022, Prague*,

Czech Republic, volume 223 of *LIPICs*, pages 23:1–23:14. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2022.

- [46] Sanjay Ranka and Todd Heywood. Two-dimensional pattern matching with k mismatches. *Pattern Recognit.*, 24(1):31–40, 1991.
- [47] Wojciech Rytter. Compressed and fully compressed pattern matching in one and two dimensions. *Proc. IEEE*, 88(11):1769–1778, 2000.
- [48] Jorma Tarhio. A sublinear algorithm for two-dimensional string matching. *Pattern Recognit. Lett.*, 17(8):833–838, 1996.
- [49] Rui Feng Zhu, Masayuki Nakajima, and Takeshi Agui. The extension of the aho-corasick algorithm to multiple rectangular patterns of different sizes and n-dimensional patterns and text. In *Proceedings of IAPR Workshop on Computer Vision - Special Hardware and Industrial Applications, MVA 1988, Tokyo, Japan, October 12-14, 1988*, pages 185–188, 1988.
- [50] Rui Feng Zhu and Tadao Takaoka. A technique for two-dimensional pattern matching. *Commun. ACM*, 32(9):1110–1120, 1989.