

# The $k$ -mismatch problem revisited

R. Clifford<sup>1</sup>, A. Fontaine<sup>1</sup>, E. Porat<sup>2</sup>, B.Sach<sup>1</sup> and T. Starikovskaya<sup>1</sup>

<sup>1</sup>University of Bristol, Department of Computer Science, Bristol, U.K. (@bristolcstheory)

<sup>2</sup>Bar-Ilan University, Department of Computer Science, Israel

allyx.fontaine@bristol.ac.uk

## The $k$ -mismatch problem

### The Hamming distance

The Hamming distance between two strings of the same length is the number of mismatches between them.

Let  $\text{Ham}(P, T)[i]$  be the Hamming distance between a  $m$ -length pattern  $P$  and a subtext  $T[j - m, i]$ .

### The $k$ -mismatch problem

Input:  $\begin{cases} \text{A pattern } P \text{ of length } m \\ \text{A streaming text } T \text{ of length } n \end{cases}$

Output: For all positions  $m \leq i < n$ ,  
 $\begin{cases} \text{Ham}(P, T[i]) & \text{if } \text{Ham}(P, T[i]) \leq k \\ \text{No} & \text{otherwise} \end{cases}$

### 3-mismatch example

```

T  b b a c c b c a b a c ...
   | | | | | | | | |
P  a b c a b | | | | |
   a b c a b | | | | |
     a b c a b | | | | |
       a b c a b | | | | |
         a b c a b | | | | |
           a b c a b | | | | |
             a b c a b

```

x: 3-mismatch position - x: No 3-mismatch - x: Mismatch

## Algorithm

### Approximate period

The  $3k/2$ -period of  $P$  is the smallest integer  $\ell$ , such that  $\text{HAM}(P[\ell, m-1], P[0, m-1-\ell]) \leq 3k/2$ .

### $3k/2$ -period example

Let  $P = \text{abcabaacabcaccacca}$  and  $k = 4$ .  
 The  $3k/2$ -period of  $P$  is  $\ell = 3$ .

```

1 shift  abcabaacabcaccacca
         abcabaacabcaccacca
2 shifts abcabaacabcaccacca
         abcabaacabcaccacca
3 shifts abcabaacabcaccacca
         abcabaacabcaccacca

```

### Case 1: Small approximate period ( $\ell \leq k$ )

1. Identify a compressible region of the text which contains all the  $k$ -mismatches.
2. Partition this region into  $O(k)$  subtexts and the pattern into  $O(k)$  subpatterns.
3. Run length encode all the subpatterns and subtexts.
4. Compute run length encoded Hamming distances for each subpattern/subtext pair.
5. Sum the Hamming distances from Step 4.

## Faster solutions using less space

### Space complexity

Problem	Previous	Ours
Deterministic online	$\tilde{O}(m)$	–
Randomised online	–	$\tilde{O}(k^2)$
$(1 + \epsilon)$ -approximation	$\tilde{O}(m/\epsilon^2)$	$\tilde{O}(k^2/\epsilon^2)$
Randomised online worst case	$\tilde{O}(k^3)$	$\tilde{O}(k^2)$

### Time complexity

Problems	Previous	Ours
Deterministic offline	$\tilde{O}(nk^3/m)$	$\tilde{O}(nk^2/m)$
Randomised online	–	$\tilde{O}(nk^2/m)$
$(1 + \epsilon)$ -approximation	–	$\tilde{O}(1/\epsilon^2)$
Randomised online worst case	$\tilde{O}(k^2)$	$\tilde{O}(\sqrt{k})$

### Run-length encoding example

Let  $P = \text{abcabaacabcaccacca}$ .

Partition	Encoding
$P^0 = \text{aaabcc}$	$(a, 3)(b, 1)(c, 2)$
$P^1 = \text{bbcccc}$	$(b, 2)(c, 4)$
$P^2 = \text{caaaaa}$	$(c, 1)(a, 5)$

All required information from  $P$  and  $T$  encoded in only  $O(k)$  space!

### Case 2: Large approximate period ( $\ell > k$ )

1. Filter out all alignments of the pattern and text with Hamming distance greater than  $3k/2$ .
2. Verify whether the Hamming distance is at most  $k$  at those positions.