

# LOCATING ALL MAXIMAL APPROXIMATE RUNS IN A STRING

Mika Amit, Maxime Crochemore and Gad M. Landau


CPM 2013, Germany



# EXACT RUNS

- An exact run is a non-empty substring of  $T$ , that can be divided into identical adjacent non-overlapping substrings:  $u^t \cdot u_1$

$T = \dots$    $\dots$

  
 $u = abc$   
(period)

  
 $u_1 = ab$

period length =  $|u| = p$



# EXACT RUNS

- An exact run is a non-empty substring of  $\mathcal{T}$ , that can be divided into identical adjacent non-overlapping substrings:  $u^t \cdot u_1$

$T = \dots$    $\dots$

**period = abc**

- A maximal exact run is an exact run that cannot be extended



# RELATED WORK

- Given a string  $T$  of length  $n$ , detecting all occurrences of exact runs in  $T$  can be done in  $O(n \log n)$ -time
  - Crochemore, 1981
  - Apostolico & Preparata, 1983
  - Main & Lorentz, 1984



# RELATED WORK

- Given a string  $T$  of length  $n$ , detecting all exact runs in  $T$  can be done in  $O(n)$ -time
  - Crochemore, 1983
  - Main, 1989
  - Kolpakov and Kucherov, 1999



# APPROXIMATE RUNS

- Finding approximate runs is sometimes more sensible than finding exact runs
- There are many different measurements and definitions of approximate runs
- The difference is in measuring the differences between the periods



# RELATED WORK

- Landau and Schmidt, 1993
- Sim, Iliopoulos, Park and Smyth, 1999
- Landau, Schmidt and Sokol, 2001
- Kolpakov and Kucherov, 2001
- Amir, Eisenberg and Levy, 2012
- ...



# APPROXIMATE RUNS - V.1

- Given a string  $\mathcal{T}$  and a number  $k$ , find a substring of  $\mathcal{T}$  that can be divided into adjacent non-overlapping substrings:  $u^1, u_1$ , such that the difference between every two adjacent periods does not exceed  $k$

$T = \dots a \ b \ d \ a \ b \ c \ a \ d \ c \ d \ d \ c \ d \ e \ c \dots$

The diagram shows a string T with characters a, b, d, a, b, c, a, d, c, d, d, c, d, e, c. Above the string, four grey arcs connect the pairs (b,d), (d,a), (a,b), and (b,c). Below the string, four purple brackets connect the pairs (d,a), (a,b), (b,c), and (c,a). These arcs and brackets highlight a substring 'dabca' which is divided into four adjacent non-overlapping substrings 'd', 'a', 'b', and 'c'.

**1-MAR**

(Maximal Approximate Run)

**(4-MAR)**



# APPROXIMATE RUNS - V.2

- Given a string  $\mathbf{T}$  and a number  $\mathbf{k}$ , find a substring of  $\mathbf{T}$  that can be divided into adjacent non-overlapping substrings:  $\mathbf{u^t.u_1}$ , such that the difference between **every** two periods does not exceed  $\mathbf{k}$

$T = \dots a \ b \ d \ a \ b \ c \ a \ d \ c \ d \ b \ c \ d \ b \ c \dots$

**2-MAR**



# APPROXIMATE RUNS - V.3

- Given a string  $\mathbf{T}$  and a number  $\mathbf{k}$ , find a substring of  $\mathbf{T}$  that can be divided into adjacent non-overlapping substrings:  $\mathbf{u}^1, \mathbf{u}_1$ , such that the difference between each period and a **consensus** substring does not exceed  $\mathbf{k}$

$T = \dots a \ b \ d \ a \ b \ c \ a \ d \ c \ a \ d \ c \ d \ b \ c \dots$

The diagram shows a sequence of characters: a, b, d, a, b, c, a, d, c, a, d, c, d, b, c. Above the characters, there are four grey arcs connecting (a,b), (b,d), (d,a), (a,b), (b,c), (c,a), (a,d), (d,c), (c,a), (a,d), (d,c), (c,d), (d,b), and (b,c). Below the characters, a purple bracket highlights the substring 'a b c' starting from the third character 'd' and ending at the sixth character 'c'. Below this bracket, the text 'u = abc' is written in red.

1-MAR

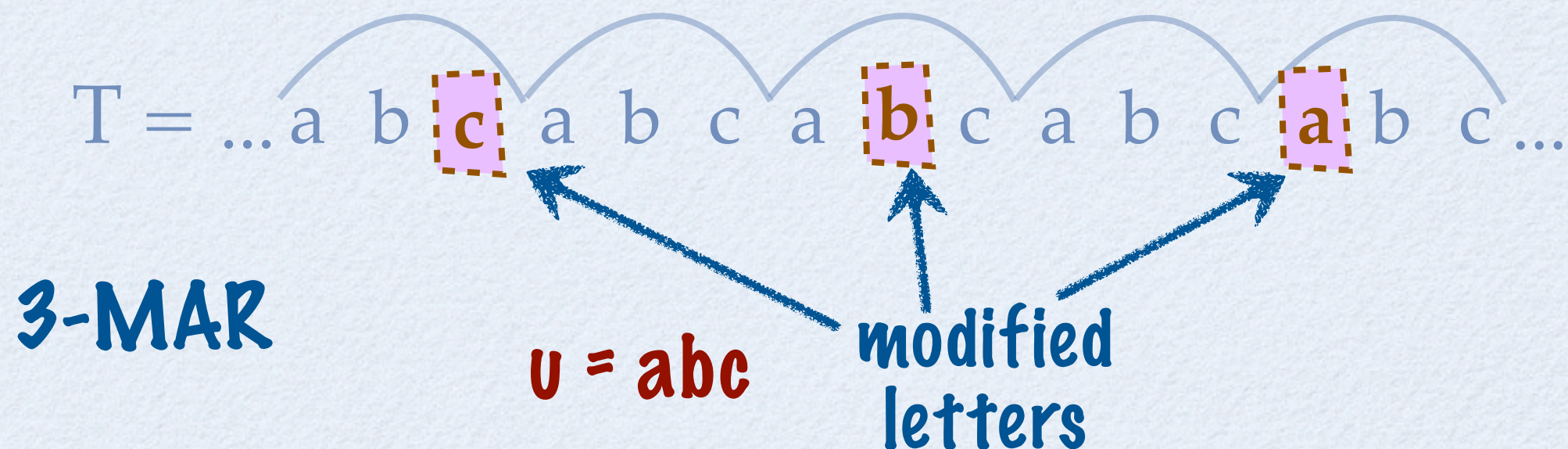
$u = abc$

(4-MAR)



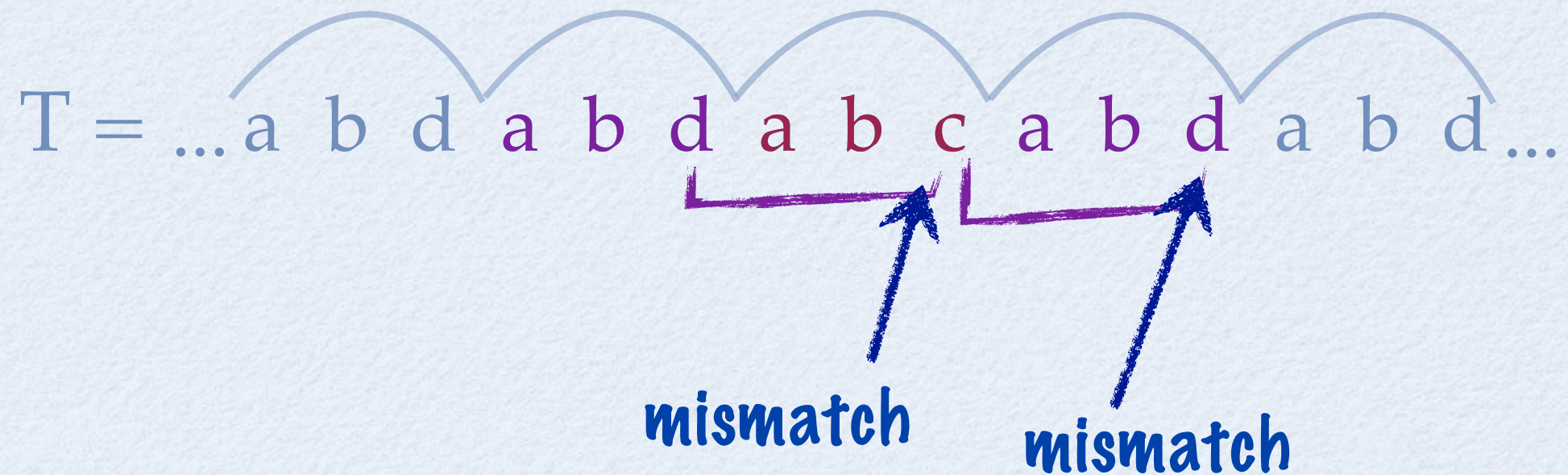
# APPROXIMATE RUNS - OUR VERSION

- Given a string  $\mathbf{T}$  and a number  $\mathbf{k}$ , find a substring of  $\mathbf{T}$  that can be divided into adjacent non-overlapping substrings:  $\mathbf{v}^t, \mathbf{v}_1$ , such that the **sum** of all differences between each period and a **consensus** substring  $\leq \mathbf{k}$





# MISMATCH VS. MODIFIED LETTER






# MISMATCH VS. MODIFIED LETTER

2 mismatches → 1 modified letter

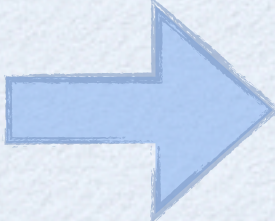
T = ... a b d a b d a b d a b d a b d ...

  
mismatch

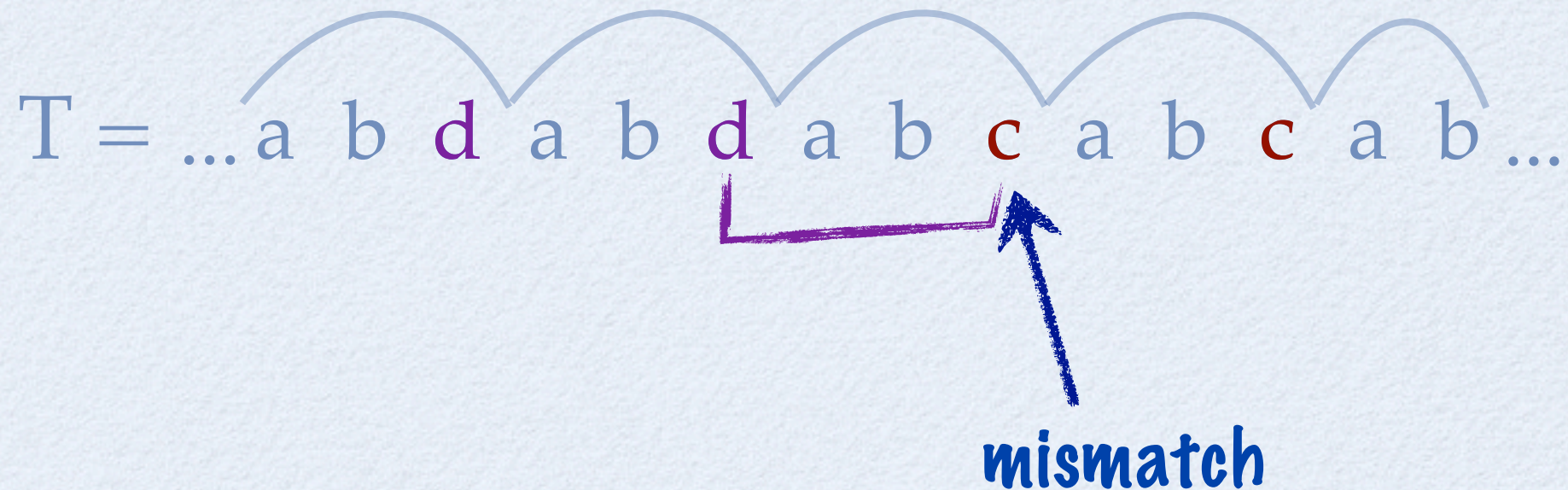


# MISMATCH VS. MODIFIED LETTER

2 mismatches  1 modified letter

1 mismatch   $\frac{n}{2p}$  modified letters

T = ... a b d a b d a b c a b c a b ...




mismatch



# MISMATCH VS. MODIFIED LETTER

2 mismatches  1 modified letter

1 mismatch   $\frac{n}{2p}$  modified letters

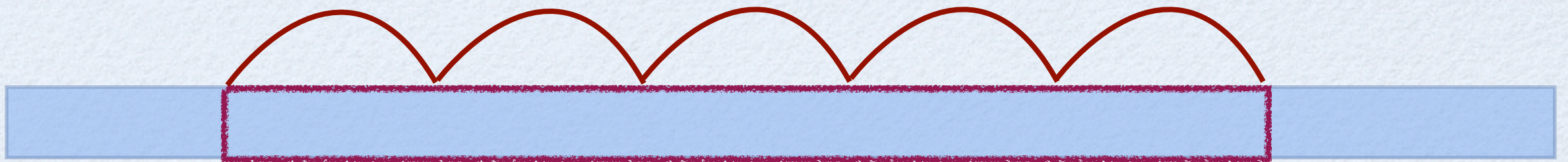
T = ... a b d a b d a b  a b  a b ...



# APPROXIMATE RUNS - PROBLEM DEFINITION

Input: a string  $\mathcal{T}$ , of length  $n$ , and a number  $k$

Output: for each period size  $1 \leq p \leq \frac{n}{2}$ ,  
find all occurrences of **k-MARs** with period  
length  $p$  in  $\mathcal{T}$

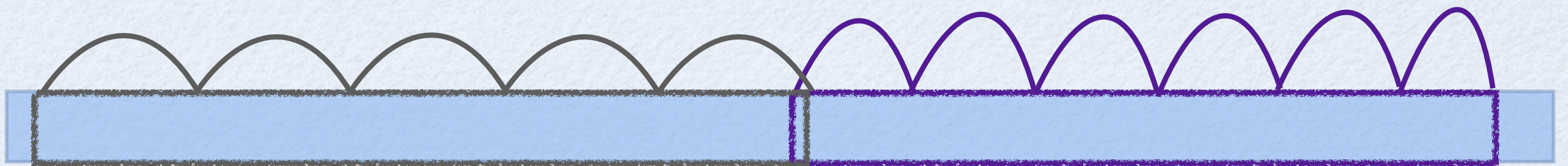




# APPROXIMATE RUNS - PROBLEM DEFINITION

Input: a string  $\tau$ , of length  $n$ , and a number  $k$

Output: for each period size  $1 \leq p \leq \frac{n}{2}$ ,  
find all occurrences of **k-MARs** with period  
length  $p$  in  $\tau$





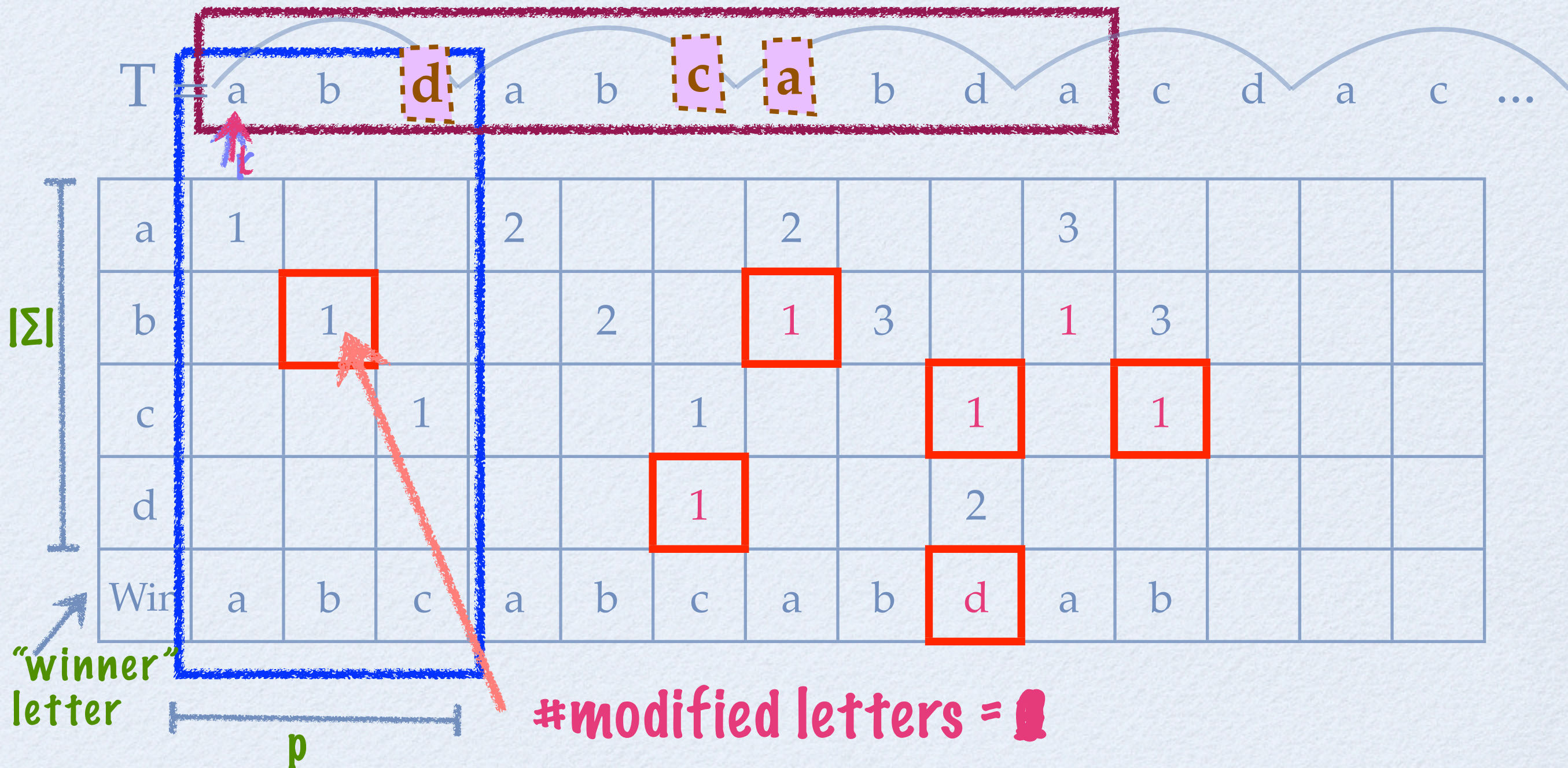
# LOCATING ALL K-MARS ALGORITHMS

- An  **$O(n^2)$**  algorithm  
(uses Parikh matrix)
- An improved  **$O(n \log n k^3)$**  algorithm  
(uses Main & Lorentz ('84) technique)
- An efficient  **$O(n \log n k^2 \log k)$**  algorithm



# AN $O(N^2)$ ALGORITHM

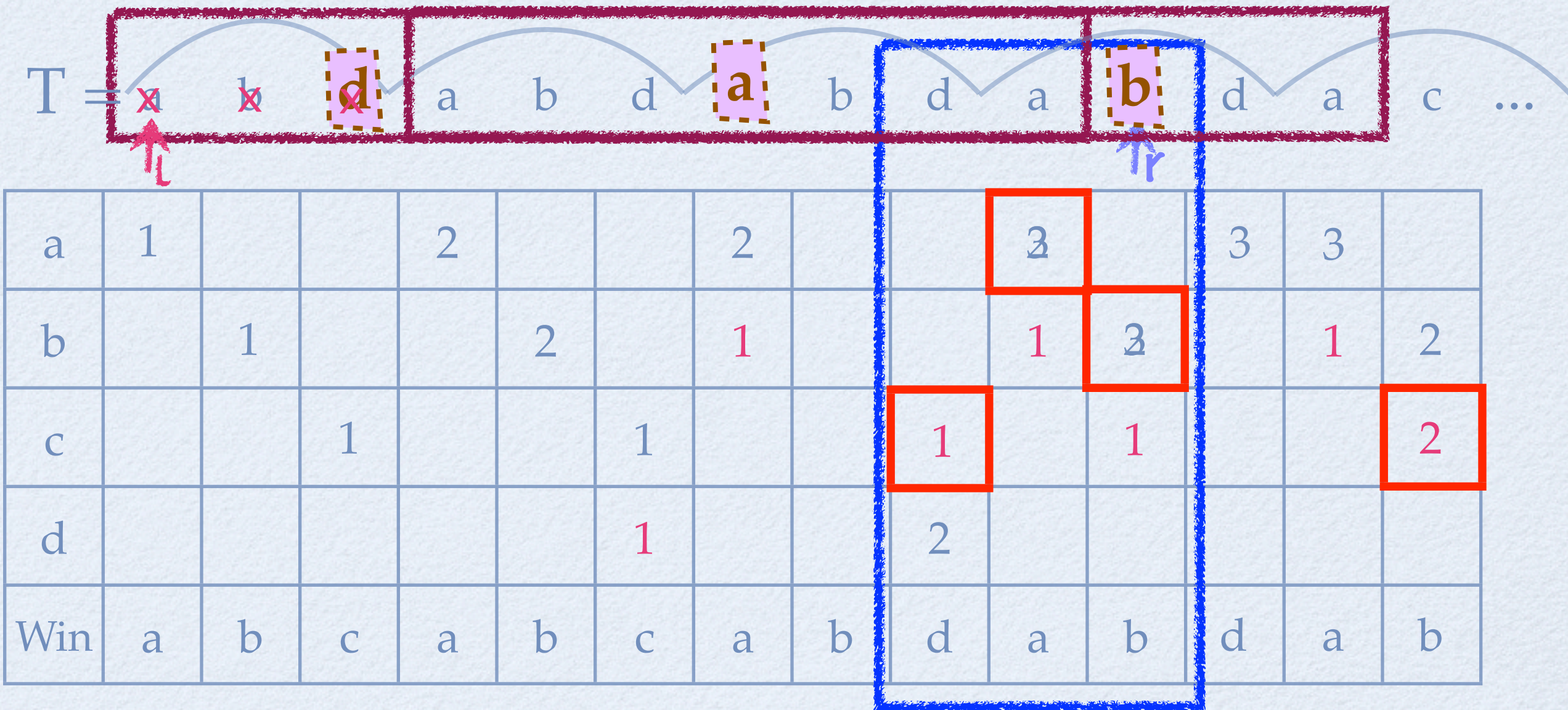
- Using Parikh matrix ( $p=3, k=2$ )





# AN $O(N^2)$ ALGORITHM

- Using Parikh matrix ( $p=3, k=2$ )



#modified letters = 3



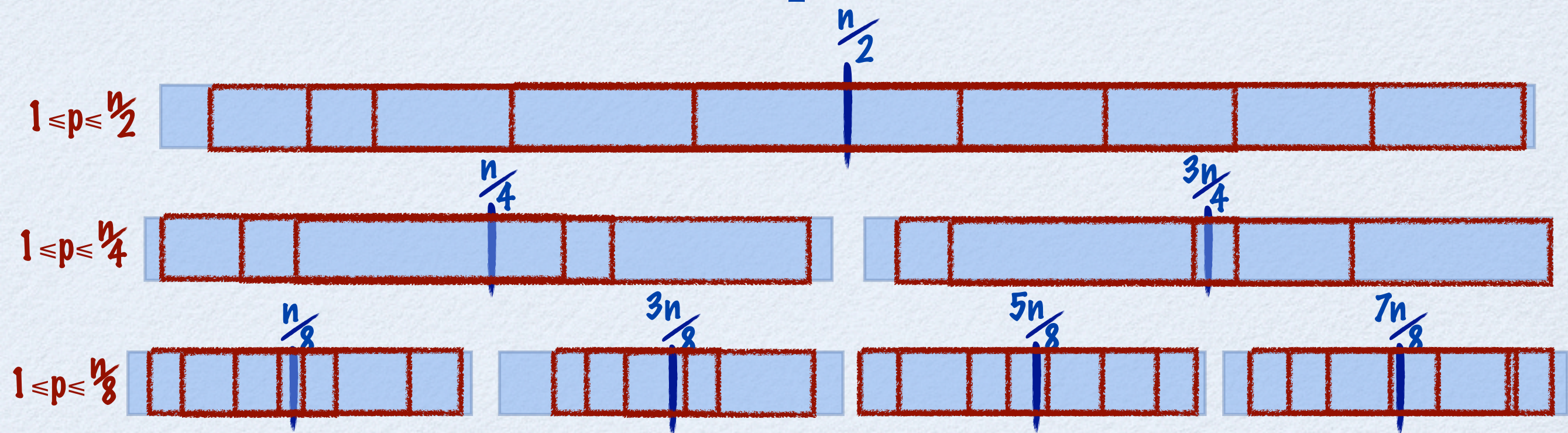
# AN $O(N^2)$ ALGORITHM

- Time complexity using Parikh matrix:
  - for each period size  $1 \leq p \leq \frac{n}{2}$
  - find all k-MARs of size  $p$  in  $O(n)$  time
  - Total:  $O(n^2)$  time (for constant alphabet)



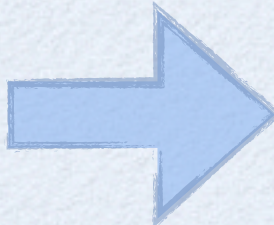
# AN IMPROVED ALGORITHM

- Divide and Conquer (Main & Lorentz):
  - \* for each period size  $p$
  - \* find all  $k$ -MARs with period size  $p$  that contain the middle position





# SOME OBSERVATIONS..

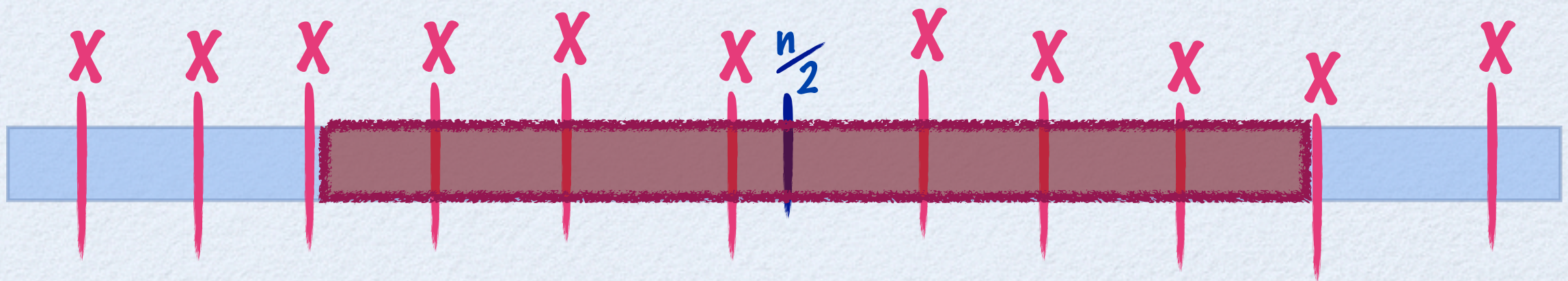
2 mismatches  1 modified letter

$T = \dots a \ b \ d \ a \ b \ d \ a \ b \ d \ a \ b \ d \dots$



Observation 1:

A **k-MAR** can contain at most  **$2k+1$**  mismatches



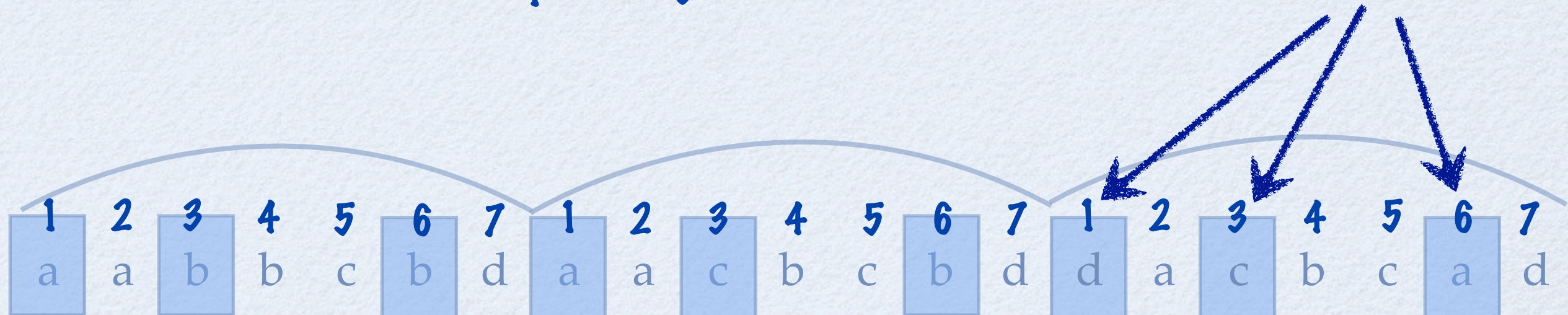
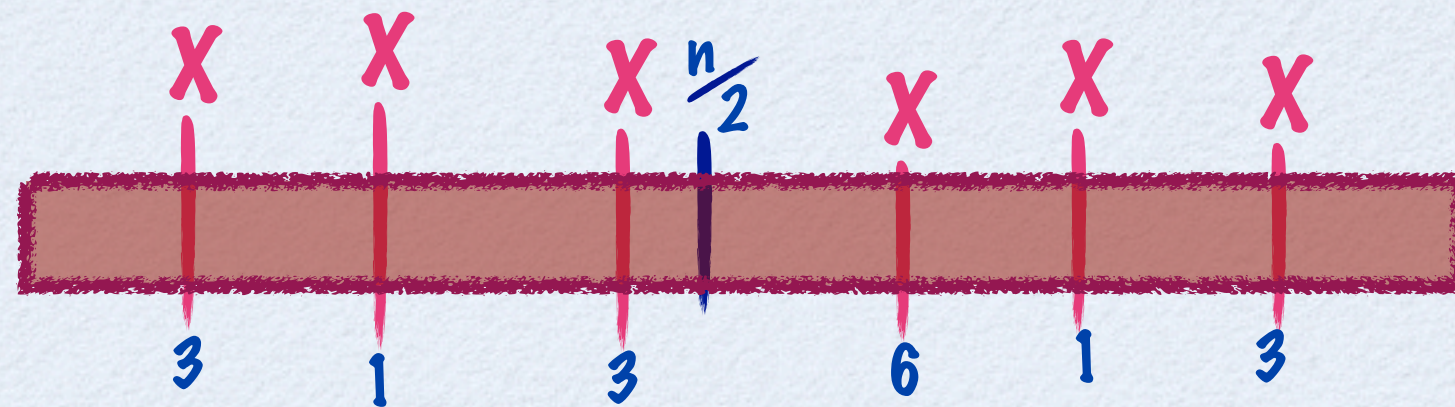


# SOME OBSERVATIONS..

## Observation 2:

A **k-MAR** can contain at most  $O(k)$  problematic columns for each period size

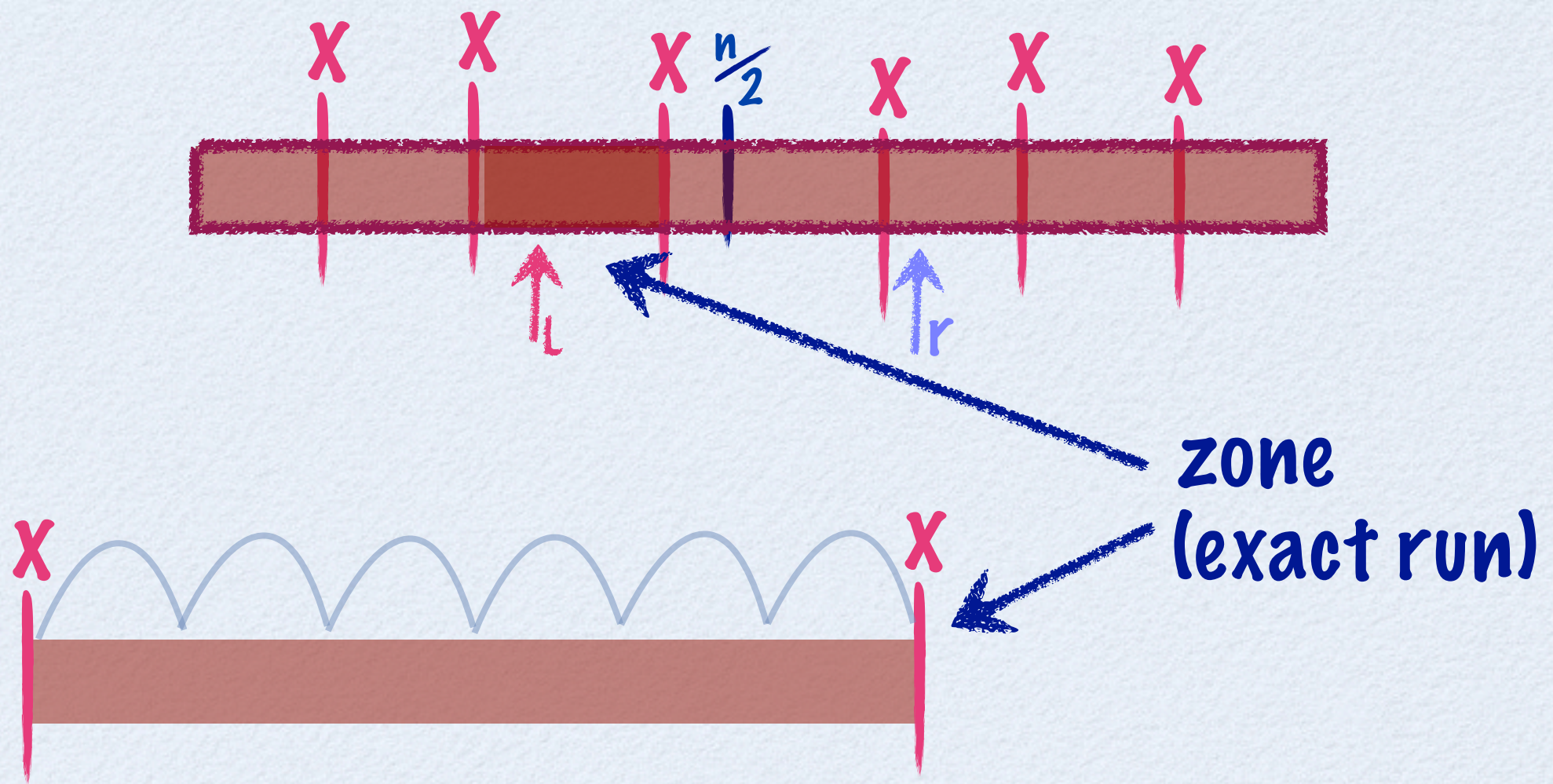
$p=7, k=1$



There are  $O(k)$   
problematic columns



# SOME OBSERVATIONS..



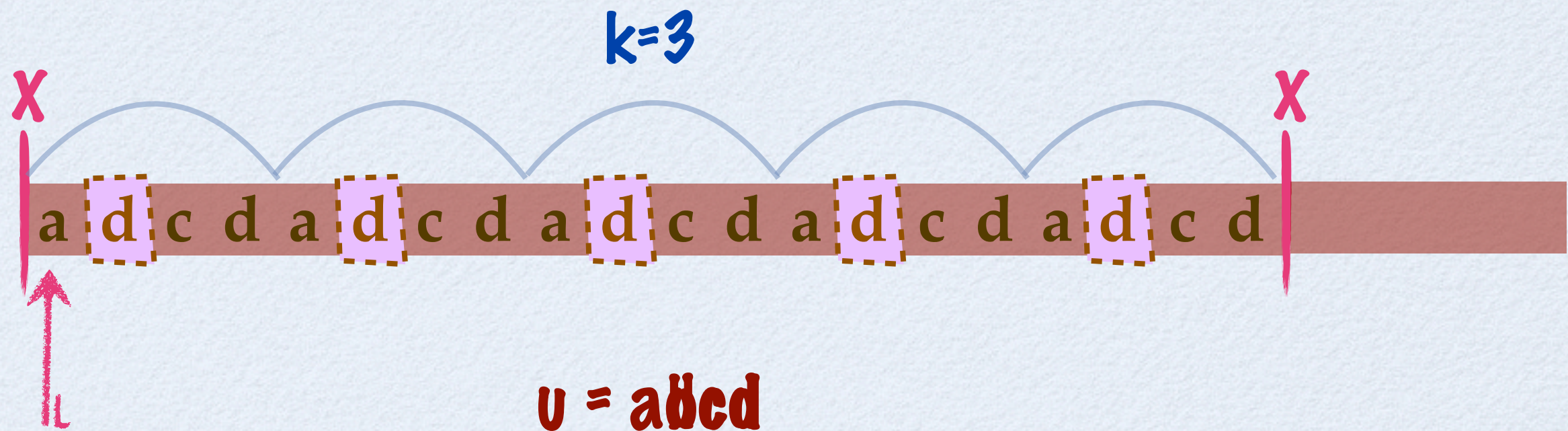
There are  $O(k)$  zones



# SOME OBSERVATIONS..

## Observation 3:

A **k-MAR** can either start on the leftmost position of a zone or on the rightmost  **$k+1$**  periods of it



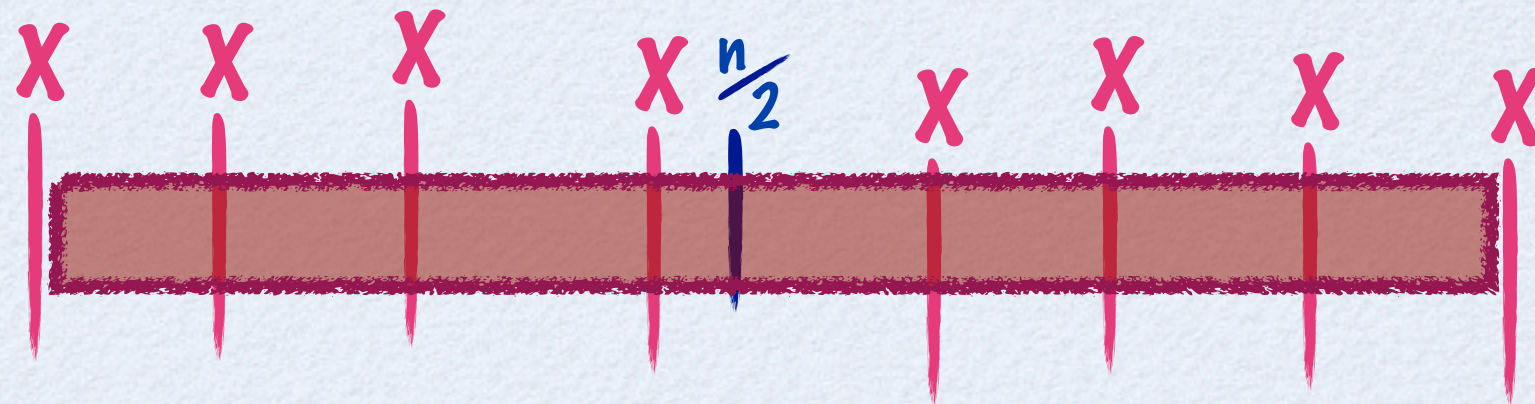
On each zone there are at most  $O(k^2)$  positions to visit



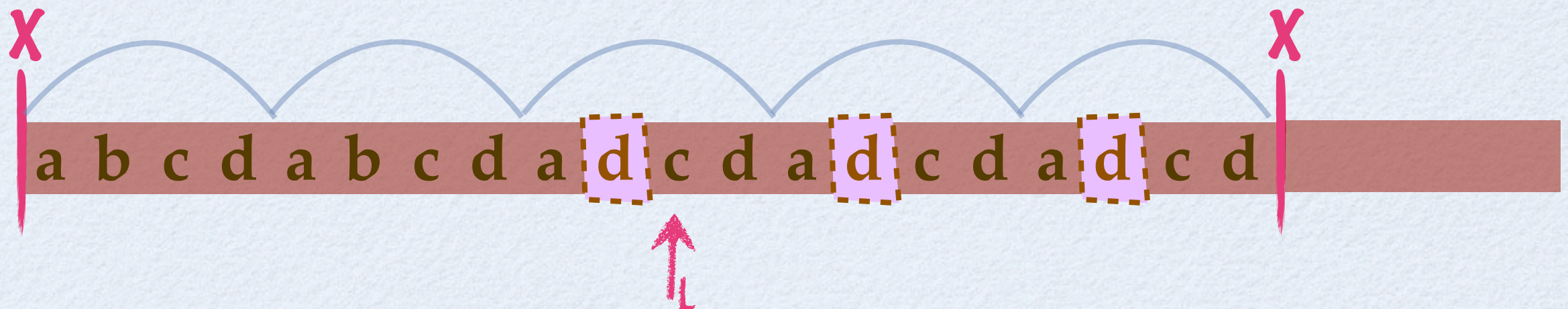
# AN IMPROVED ALGORITHM

Given a period length  $p$ , and a text  $T$ :

\* Find  $4k+2$  mismatch positions:



\* On each zone visit *only* problematic columns in the last  $k+1$  periods





# AN IMPROVED ALGORITHM

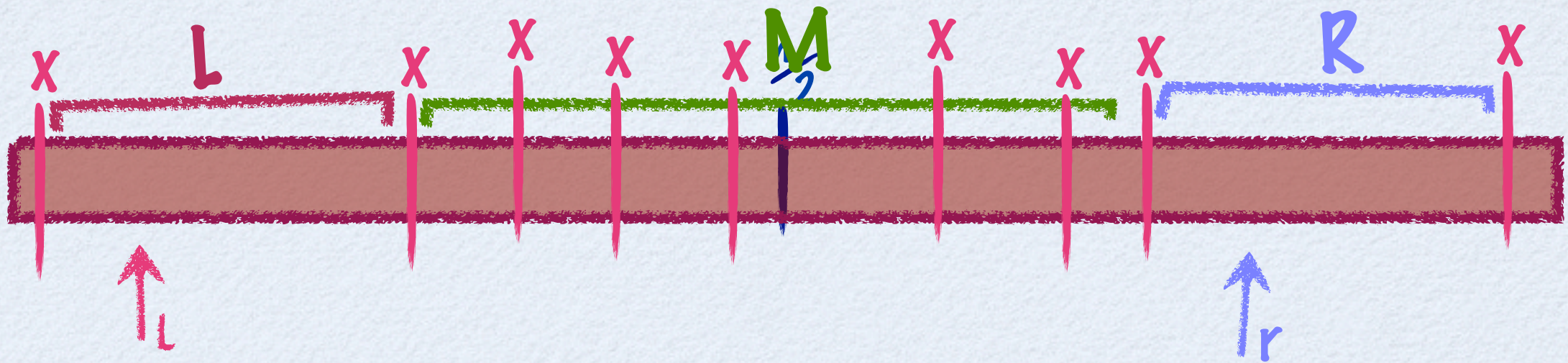
- Time complexity of iteration 1:
  - \* for all period length  $1 \leq p \leq \frac{n}{2}$
  - \* find all **k-MARs** that contain position  $\frac{n}{2}$ 
    1. find  **$O(k)$**  mismatch positions in  **$O(k)$**  time
    2. for each one of the  **$O(k)$**  zones:  
visit only  **$O(k^2)$**  positions in  **$O(k^3)$**  time
- Total of  **$O(nk^3)$**  for iteration 1
- Total of  **$O(n \log n k^3)$**  time for the entire algorithm



# THE EFFICIENT ALGORITHM

Observation:

Not all  $O(k^2)$  positions in a zone need to be visited

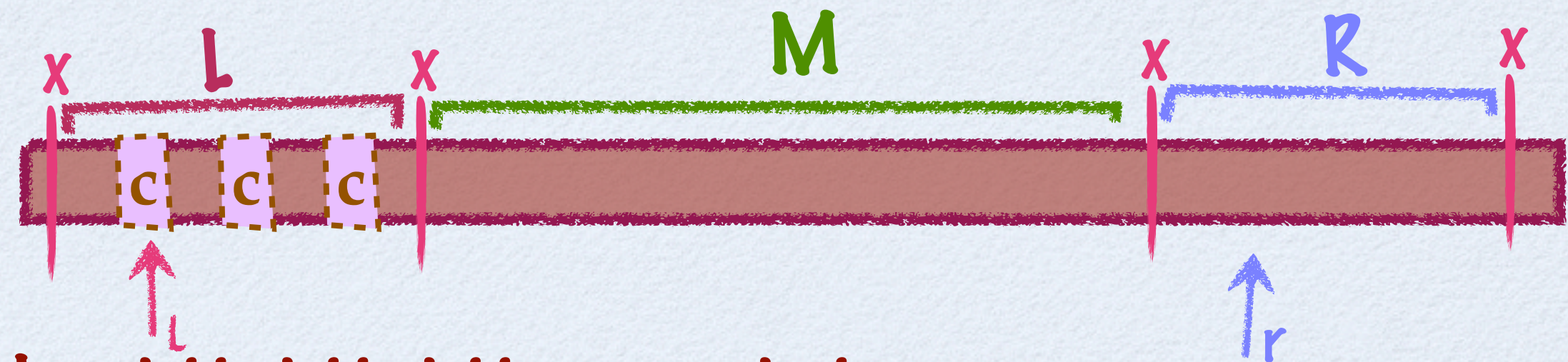




# THE EFFICIENT ALGORITHM

## Observation 3:

Not all  $O(k^2)$  positions in a zone need to be visited



$L = abdd\ abdd\ abdd$

problematic  
positions

$U = abcd$

problematic  
columns

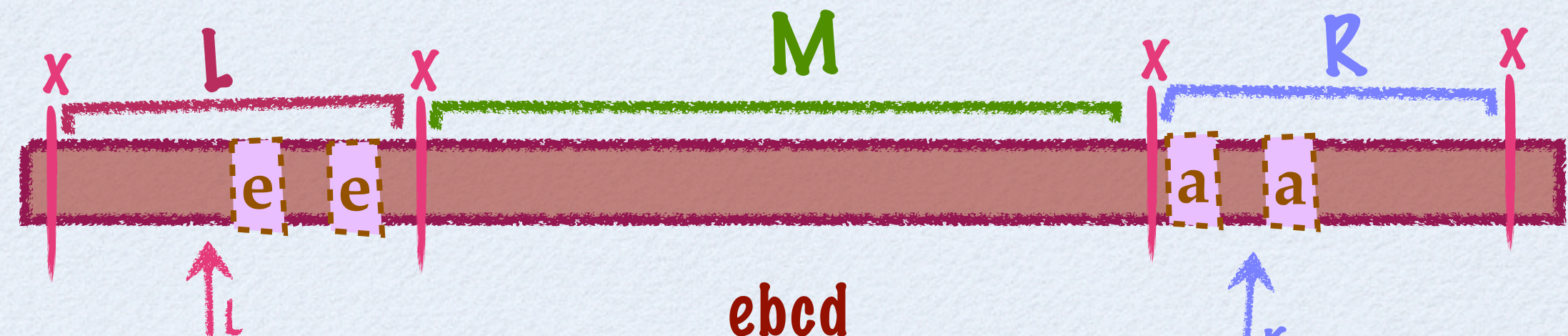
initial number of  
problematic positions:  $O(k)$



# THE EFFICIENT ALGORITHM

## Observation 3:

Not all  $O(k^2)$  positions in a zone need to be visited



$L = abdd\ abdd\ abdd$

“new”  
problematic  
positions

$u = \begin{matrix} ebcd \\ \cancel{abcd} \end{matrix}$

problematic  
columns

$R = ebcd\ ebcd\ ebcd\ ebcd$

number of added  
problematic  
positions:  $O(k \log k)$



# THE EFFICIENT ALGORITHM

- Time complexity of iteration 1:
  - \* for all period length  $1 \leq p \leq \frac{n}{2}$
  - \* find all **k-MARs** that contain position  $\frac{n}{2}$ 
    1. find  **$O(k)$**  mismatch positions in  **$O(k)$**  time
    2. for each one of the  **$O(k)$**  zones:  
visit  **$O(k \log k)$**  positions in  **$O(k^2 \log k)$**  time
- Total  **$O(nk^2 \log k)$**  for iteration 1
- Total  **$O(n \log n k^2 \log k)$**  time for the entire algorithm



Thank you!