

# **Jumbled String Matching: Online, offline, binary or not**

**Zsuzsanna Lipták**

University of Verona (Italy)

McMaster University, 10 Oct. 2014

# Jumbled String Matching

Parikh vectors:

Given string  $t$  over constant-size ordered alphabet  $\Sigma$ , with  $|\Sigma| = \sigma$ .

The Parikh vector  $p(t)$  counts the multiplicity of characters in  $t$ .

Ex.:  $p(aabbcac) = (3, 1, 2)$ .

# Jumbled String Matching

Parikh vectors:

Given string  $t$  over constant-size ordered alphabet  $\Sigma$ , with  $|\Sigma| = \sigma$ .

The **Parikh vector  $p(t)$**  counts the multiplicity of characters in  $t$ .

Ex.:  $p(aabdac) = (3, 1, 2)$ .

(a.k.a. jumbled string = compomer = composition = multiplicity vector = statistics = histogram . . . . . )

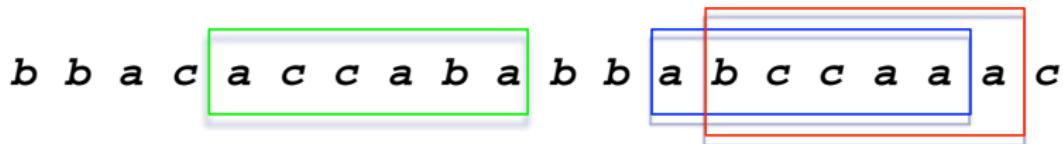
# Problem Statement

## JUMBLED STRING MATCHING

Given string  $s$  of length  $n$ , and query Parikh vector  $q \in \mathbb{N}^\sigma$ .

Find all occurrences of substrings  $t$  of  $s$  s.t.  $p(t) = q$ .

Ex.:  $\Sigma = \{a, b, c\}$ , query  $q = (3, 1, 2)$



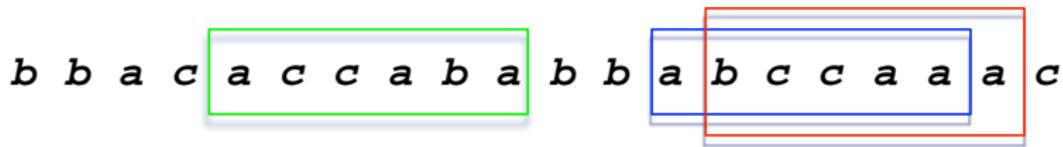
## Problem Statement

### JUMBLED STRING MATCHING

Given string  $s$  of length  $n$ , and query Parikh vector  $q \in \mathbb{N}^\sigma$ .

Find all occurrences of substrings  $t$  of  $s$  s.t.  $p(t) = q$ .

Ex.:  $\Sigma = \{a, b, c\}$ , query  $q = (3, 1, 2)$



= any permutation of  $aabcac$

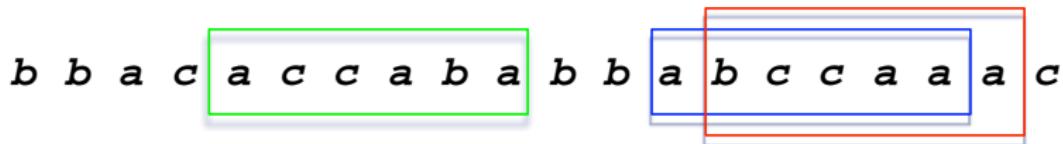
## Problem Statement

### JUMBLED STRING MATCHING

Given string  $s$  of length  $n$ , and query Parikh vector  $q \in \mathbb{N}^\sigma$ .

Find all occurrences of substrings  $t$  of  $s$  s.t.  $p(t) = q$ .

Ex.:  $\Sigma = \{a, b, c\}$ , query  $q = (3, 1, 2)$



= any permutation of  $aabcac$  = any jumble of  $aabcac$

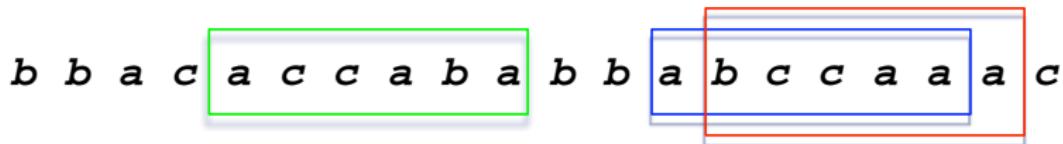
## Problem Statement

### JUMBLED STRING MATCHING

Given string  $s$  of length  $n$ , and query Parikh vector  $q \in \mathbb{N}^\sigma$ .

Find all occurrences of substrings  $t$  of  $s$  s.t.  $p(t) = q$ .

Ex.:  $\Sigma = \{a, b, c\}$ , query  $q = (3, 1, 2)$



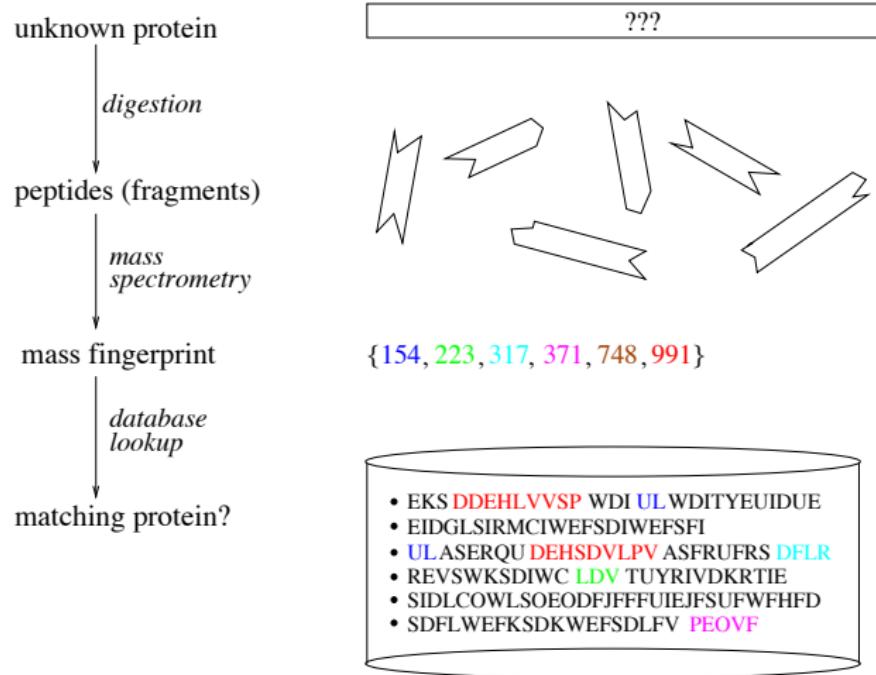
= any permutation of  $aabcac$  = any jumble of  $aabcac$

a.k.a. permutation matching, Parikh vector matching, abelian matching, JPM =  
jumbled pattern matching, jumbled indexing - see later

# Motivations, Applications

- Mass spectrometry
- Gene clusters
- Motif search in graphs and trees
- Filter for exact pattern matching

# Protein identification with mass spectrometry (here: PMF)

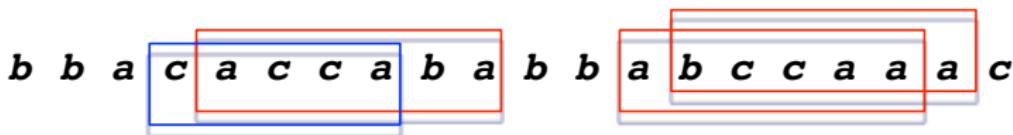


Takes advantage of different molecular masses of “characters” (AAs, nucleotides, . . . )

## Modelling sample identification with MS

Every character has a mass:  $\mu : \Sigma \rightarrow \mathbb{R}^+$  mass function,  $\mu(t) = \sum_i \mu(t_i)$ .

Ex:  $\Sigma = \{a, b, c\}$  with  $\mu(a) = 2, \mu(b) = 3, \mu(c) = 5$ . Query  $M = 19$



Actually we can also look for all Parikh vectors with query mass  $M$ !  
⇒ Jumbled String Matching!

## Application 2: Gene clusters

Given:  $k$  genomes

Find: maximal blocks consisting of same genes

1 2 **3 4 3 5 1 5 4** 7 2 5

**5 3 1 4 4** 7 1 7 2 2 1 3

1 2 6 7 **5 1 3 5 4** 2 2 5

gene cluster: {1, 3, 4, 5}

## Application 2: Gene clusters

Given:  $k$  genomes

Find: maximal blocks consisting of same genes

1 2 **3 4 3 5 1 5 4** 7 2 5

**5 3 1 4 4** 7 1 7 2 2 1 3

1 2 6 7 **5 1 3 5 4** 2 2 5

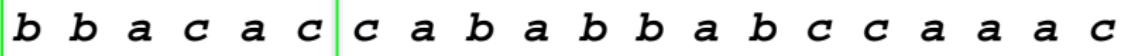
gene cluster: {1, 3, 4, 5}

Caveat: Problem slightly different (so far).

# 1. Simple solutions

# Window algorithm

Jumbled string matching query  $q = (3, 1, 2)$



b b a c a c | c a b a b b a b c c a a a c

- sliding window (of fixed-size  $m = |q| = \sum_i q_i$ )

# Window algorithm

Jumbled string matching query  $q = (3, 1, 2)$



b b a c a c c a b a b b a b c c a a a c

- sliding window (of fixed-size  $m = |q| = \sum_i q_i$ )

# Window algorithm

Jumbled string matching query  $q = (3, 1, 2)$

b b a c a c c a b a b b a b c c a a a a c

- sliding window (of fixed-size  $m = |q| = \sum_i q_i$ )

# Window algorithm

Jumbled string matching query  $q = (3, 1, 2)$

b b a c a c c a b a b b a b c c a a a a c



- sliding window (of fixed-size  $m = |q| = \sum_i q_i$ )

# Window algorithm

Jumbled string matching query  $q = (3, 1, 2)$

b b a c a c c a b a b b a b c c a a a c

- sliding window (of fixed-size  $m = |q| = \sum_i q_i$ )

# Window algorithm

Jumbled string matching query  $q = (3, 1, 2)$

$b \ b \ a \ c \ a \ c \ c \ a \ b \ a \ b \ b \ a \ b \ c \ c \ a \ a \ a \ c$

- sliding window (of fixed-size  $m = |q| = \sum_i q_i$ )
- worst-case optimal for one query:  $O(n)$  time,  $O(\sigma)$  additional space.

## Indexed jumbled string matching

What about many queries?  $\rightsquigarrow$  indexed version of problem

simple solutions ( $K$  queries):

1. no index:  $O(Kn)$  query time
2. store all:  $O(n^2)$  index size,  $O(K \log n)$  query time.

## Indexed jumbled string matching

What about many queries?  $\rightsquigarrow$  indexed version of problem

simple solutions ( $K$  queries):

1. no index:  $O(Kn)$  query time
2. store all:  $O(n^2)$  index size,  $O(K \log n)$  query time.

For exact string matching, elaborate solutions exist:  
suffix trees, suffix arrays, ...

- index has  $O(n)$  size (size of text)
- construction time and space  $O(n)$
- query time  $O(m)$  (size of query) – so  $O(Km)$  for  $K$  queries

## 2. General alphabets: Jumping Algorithm

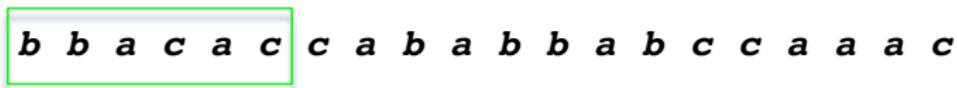
# Jumping Algorithm for Jumbled String Matching

Cicalese, Fici, L. (PSC 2009, FUN 2010)

Recall the window algorithm.

fixed size window → variable size window

$$q = (312)$$



The diagram shows a horizontal sequence of characters: b, b, a, c, a, c, c, a, b, a, b, b, a, b, c, c, c, a, a, a, c. A window of length 6 is highlighted with a green border, starting at the second character ('b') and ending at the seventh character ('c').

# Jumping Algorithm for Jumbled String Matching

Cicalese, Fici, L. (PSC 2009, FUN 2010)

Recall the window algorithm.

fixed size window → variable size window

$$q = (312)$$

The diagram illustrates a window of size 7 being滑动 over a string of length 15. The window is highlighted with a green border and contains the characters b, b, a, c, a, c, c. The string below it is a, b, a, b, b, a, b, c, c, a, a, a, c.

# Jumping Algorithm for Jumbled String Matching

Cicalese, Fici, L. (PSC 2009, FUN 2010)

Recall the window algorithm.

fixed size window → variable size window

$$q = (312)$$

b b a c a c c a      b a b b a b c c a a a c

# Jumping Algorithm for Jumbled String Matching

Cicalese, Fici, L. (PSC 2009, FUN 2010)

Recall the window algorithm.

fixed size window → variable size window

$$q = (312)$$

b a c a c c a      b a b b a b c c a a a c

# Jumping Algorithm for Jumbled String Matching

Cicalese, Fici, L. (PSC 2009, FUN 2010)

Recall the window algorithm.

fixed size window → variable size window

$$q = (312)$$

b b a c a c c c a b a b b a b c c c a a a a c

# Jumping Algorithm for Jumbled String Matching

Cicalese, Fici, L. (PSC 2009, FUN 2010)

Recall the window algorithm.

fixed size window → variable size window

$$q = (312)$$

b b a c a c c a b a b b a b c c a a a c

# Jumping Algorithm for Jumbled String Matching

Cicalese, Fici, L. (PSC 2009, FUN 2010)

Recall the window algorithm.

fixed size window → variable size window

$$q = (312)$$



A sequence of characters: **b** **b** **a** **c** **a** **c** **c** **a** **b** **a** **b** **b** **a** **b** **c** **c** **a** **a** **a** **c**. A green rectangular box highlights the substring **a** **c** **c** **a**.

# Jumping Algorithm for Jumbled String Matching

Cicalese, Fici, L. (PSC 2009, FUN 2010)

Recall the window algorithm.

fixed size window → variable size window

$$q = (312)$$



b b a c a c c a b a b b a b c c a a a c

The Jumping Algorithm simulates these moves by jumps.

## Jumping algorithm: update rules

$q = (312)$

|          | 0        | 1        | 2        | 3        | 4        | 5        | 6        | 7        | 8        | 9        | 10       | 11       | 12       | 13       | 14       | 15       | 16       | 17       | 18       | 19       | 20 |
|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----|
| <b>L</b> | <b>b</b> | <b>b</b> | <b>a</b> | <b>c</b> | <b>a</b> | <b>c</b> | <b>c</b> | <b>a</b> | <b>b</b> | <b>a</b> | <b>b</b> | <b>b</b> | <b>a</b> | <b>b</b> | <b>c</b> | <b>c</b> | <b>a</b> | <b>a</b> | <b>a</b> | <b>c</b> |    |
| <b>a</b> | 0        | 0        | 1        | 1        | 2        | 2        | 2        | 3        | 3        | 4        | 4        | 4        | 4        | 5        | 5        | 5        | 5        | 6        | 7        | 8        | 8  |
| <b>b</b> | 1        | 2        | 2        | 2        | 2        | 2        | 2        | 2        | 3        | 3        | 3        | 4        | 5        | 5        | 6        | 6        | 6        | 6        | 6        | 6        | 6  |
| <b>c</b> | 0        | 0        | 0        | 1        | 1        | 2        | 3        | 3        | 3        | 3        | 3        | 3        | 3        | 3        | 4        | 5        | 5        | 5        | 5        | 6        |    |

## Jumping algorithm: update rules

$q = (312)$

|          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |    |
|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----|
|          | 0        | 1        | 2        | 3        | 4        | 5        | 6        | 7        | 8        | 9        | 10       | 11       | 12       | 13       | 14       | 15       | 16       | 17       | 18       | 19       | 20 |
| <b>L</b> |          |          |          |          |          |          |          |          | <b>R</b> |          |          |          |          |          |          |          |          |          |          |          |    |
|          | <b>b</b> | <b>b</b> | <b>a</b> | <b>c</b> | <b>a</b> | <b>c</b> | <b>c</b> | <b>a</b> |          | <b>b</b> | <b>a</b> | <b>b</b> | <b>b</b> | <b>a</b> | <b>b</b> | <b>c</b> | <b>c</b> | <b>a</b> | <b>a</b> | <b>a</b> |    |
| <b>a</b> | 0        | 0        | 1        | 1        | 2        | 2        | 2        | 2        | 3        | 3        | 4        | 4        | 4        | 5        | 5        | 5        | 5        | 6        | 7        | 8        | 8  |
| <b>b</b> | 1        | 2        | 2        | 2        | 2        | 2        | 2        | 2        | 3        | 3        | 3        | 4        | 5        | 5        | 6        | 6        | 6        | 6        | 6        | 6        | 6  |
| <b>c</b> | 0        | 0        | 0        | 1        | 1        | 2        | 3        | 3        | 3        | 3        | 3        | 3        | 3        | 3        | 4        | 5        | 5        | 5        | 5        | 6        |    |

## Jumping algorithm: update rules

$q = (312)$

|          | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7        | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |
|----------|---|---|---|---|---|---|---|----------|---|---|----|----|----|----|----|----|----|----|----|----|----|
| <b>L</b> | b | b | a | c | a | c | c | <b>a</b> | b | a | b  | b  | a  | b  | c  | c  | a  | a  | a  | c  |    |
| <b>R</b> |   |   |   |   |   |   |   |          |   |   |    |    |    |    |    |    |    |    |    |    |    |
| <b>a</b> | 0 | 0 | 1 | 1 | 2 | 2 | 2 | 3        | 3 | 4 | 4  | 4  | 5  | 5  | 5  | 5  | 6  | 7  | 8  | 8  |    |
| <b>b</b> | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 2        | 3 | 3 | 4  | 5  | 5  | 6  | 6  | 6  | 6  | 6  | 6  | 6  |    |
| <b>c</b> | 0 | 0 | 0 | 1 | 1 | 2 | 3 | 3        | 3 | 3 | 3  | 3  | 3  | 3  | 4  | 5  | 5  | 5  | 5  | 6  |    |

## Jumping algorithm: update rules

$$q = (312)$$

|          | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |
|----------|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|
| <b>L</b> | b | b | a | c | a | c | c | a |   | b | a  | b  | b  | a  | b  | c  | c  | a  | a  | a  | c  |
| <b>R</b> |   |   |   |   |   |   |   |   |   |   |    |    |    |    |    |    |    |    |    |    |    |
| <b>a</b> | 0 | 0 | 1 | 1 | 2 | 2 | 2 | 2 | 3 | 3 | 4  | 4  | 4  | 5  | 5  | 5  | 5  | 6  | 7  | 8  | 8  |
| <b>b</b> | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 3 | 3 | 3  | 4  | 5  | 5  | 6  | 6  | 6  | 6  | 6  | 6  | 6  |
| <b>c</b> | 0 | 0 | 0 | 1 | 1 | 2 | 3 | 3 | 3 | 3 | 3  | 3  | 3  | 3  | 4  | 5  | 5  | 5  | 5  | 6  |    |

update  $R$ :  $R \leftarrow \max_{a \in \Sigma} \underbrace{\text{select}_a(\text{rank}_a(L) + q_a)}_{\text{necessary char's}}$

## Jumping algorithm: update rules

$q = (312)$

|          |          |          |          |          |          |          |          |          |   |          |          |          |          |          |          |          |          |          |          |          |          |
|----------|----------|----------|----------|----------|----------|----------|----------|----------|---|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|
|          | 0        | 1        | 2        | 3        | 4        | 5        | 6        | 7        | 8 | 9        | 10       | 11       | 12       | 13       | 14       | 15       | 16       | 17       | 18       | 19       | 20       |
|          | <b>b</b> | <b>b</b> | <b>a</b> | <b>c</b> | <b>a</b> | <b>c</b> | <b>c</b> | <b>a</b> |   | <b>b</b> | <b>a</b> | <b>b</b> | <b>b</b> | <b>a</b> | <b>b</b> | <b>c</b> | <b>c</b> | <b>a</b> | <b>a</b> | <b>a</b> | <b>c</b> |
| <b>a</b> | 0        | 0        | 1        | 1        | 2        | 2        | 2        | 2        | 3 | 3        | 4        | 4        | 4        | 5        | 5        | 5        | 5        | 6        | 7        | 8        | 8        |
| <b>b</b> | 1        | 2        | 2        | 2        | 2        | 2        | 2        | 2        | 3 | 3        | 3        | 4        | 5        | 5        | 6        | 6        | 6        | 6        | 6        | 6        | 6        |
| <b>c</b> | 0        | 0        | 0        | 1        | 1        | 2        | 3        | 3        | 3 | 3        | 3        | 3        | 3        | 3        | 3        | 4        | 5        | 5        | 5        | 5        | 6        |

update  $R$ :  $R \leftarrow \max_{a \in \Sigma} \underbrace{\text{select}_a(\text{rank}_a(L) + q_a)}_{\text{necessary char's}}$

## Jumping algorithm: update rules

$q = (312)$

|          |          |          |          |          |          |          |          |          |   |          |          |          |          |          |          |          |          |          |          |          |          |
|----------|----------|----------|----------|----------|----------|----------|----------|----------|---|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|
|          | 0        | 1        | 2        | 3        | 4        | 5        | 6        | 7        | 8 | 9        | 10       | 11       | 12       | 13       | 14       | 15       | 16       | 17       | 18       | 19       | 20       |
|          | <b>b</b> | <b>b</b> | <b>a</b> | <b>c</b> | <b>a</b> | <b>c</b> | <b>c</b> | <b>a</b> |   | <b>b</b> | <b>a</b> | <b>b</b> | <b>b</b> | <b>a</b> | <b>b</b> | <b>c</b> | <b>c</b> | <b>a</b> | <b>a</b> | <b>a</b> | <b>c</b> |
| <b>a</b> | 0        | 0        | 1        | 1        | 2        | 2        | 2        | 2        | 3 | 3        | 4        | 4        | 4        | 5        | 5        | 5        | 5        | 6        | 7        | 8        | 8        |
| <b>b</b> | 1        | 2        | 2        | 2        | 2        | 2        | 2        | 2        | 3 | 3        | 3        | 4        | 5        | 5        | 6        | 6        | 6        | 6        | 6        | 6        | 6        |
| <b>c</b> | 0        | 0        | 0        | 1        | 1        | 2        | 3        | 3        | 3 | 3        | 3        | 3        | 3        | 3        | 3        | 4        | 5        | 5        | 5        | 5        | 6        |

update  $R$ :  $R \leftarrow \max_{a \in \Sigma} \underbrace{\text{select}_a(\text{rank}_a(L) + q_a)}_{\text{necessary char's}}$

update  $L$ :  $L \leftarrow \max_{a \in \Sigma} \underbrace{\text{select}_a(\text{rank}_a(R) - q_a)}_{\text{unnecessary char's}}$  (no match),

## Jumping algorithm: update rules

$q = (312)$

|          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |    |   |
|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----|---|
|          | 0        | 1        | 2        | 3        | 4        | 5        | 6        | 7        | 8        | 9        | 10       | 11       | 12       | 13       | 14       | 15       | 16       | 17       | 18       | 19       | 20 |   |
|          | <b>b</b> | <b>b</b> | <b>a</b> | <b>c</b> | <b>a</b> | <b>c</b> | <b>c</b> | <b>a</b> | <b>b</b> | <b>a</b> | <b>b</b> | <b>b</b> | <b>a</b> | <b>b</b> | <b>c</b> | <b>c</b> | <b>a</b> | <b>a</b> | <b>a</b> | <b>c</b> |    |   |
| <b>a</b> | 0        | 0        | 1        | 1        | 2        | 2        | 2        | 2        | 3        | 3        | 4        | 4        | 4        | 4        | 5        | 5        | 5        | 5        | 6        | 7        | 8  | 8 |
| <b>b</b> | 1        | 2        | 2        | 2        | 2        | 2        | 2        | 2        | 3        | 3        | 3        | 4        | 5        | 5        | 6        | 6        | 6        | 6        | 6        | 6        | 6  | 6 |
| <b>c</b> | 0        | 0        | 0        | 1        | 1        | 2        | 3        | 3        | 3        | 3        | 3        | 3        | 3        | 3        | 3        | 4        | 5        | 5        | 5        | 5        | 6  |   |

update  $R$ :  $R \leftarrow \max_{a \in \Sigma} \underbrace{\text{select}_a(\text{rank}_a(L) + q_a)}_{\text{necessary char's}}$

update  $L$ :  $L \leftarrow \max_{a \in \Sigma} \underbrace{\text{select}_a(\text{rank}_a(R) - q_a)}_{\text{unnecessary char's}}$  (no match),

## Jumping algorithm: update rules

$q = (312)$

|          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |
|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|
|          | 0        | 1        | 2        | 3        | 4        | 5        | 6        | 7        | 8        | 9        | 10       | 11       | 12       | 13       | 14       | 15       | 16       | 17       | 18       | 19       | 20       |
|          | <b>b</b> | <b>b</b> | <b>a</b> | <b>c</b> | <b>a</b> | <b>L</b> | <b>c</b> | <b>c</b> | <b>a</b> | <b>b</b> | <b>a</b> | <b>R</b> | <b>b</b> | <b>b</b> | <b>a</b> | <b>b</b> | <b>c</b> | <b>c</b> | <b>a</b> | <b>a</b> | <b>a</b> |
| <b>a</b> | 0        | 0        | 1        | 1        | 2        | 2        | 2        | 2        | 3        | 3        | 4        | 4        | 4        | 4        | 5        | 5        | 5        | 5        | 6        | 7        | 8        |
| <b>b</b> | 1        | 2        | 2        | 2        | 2        | 2        | 2        | 2        | 3        | 3        | 3        | 4        | 5        | 5        | 6        | 6        | 6        | 6        | 6        | 6        | 6        |
| <b>c</b> | 0        | 0        | 0        | 1        | 1        | 2        | 3        | 3        | 3        | 3        | 3        | 3        | 3        | 3        | 3        | 4        | 5        | 5        | 5        | 5        | 6        |

update  $R$ :  $R \leftarrow \max_{a \in \Sigma} \underbrace{\text{select}_a(\text{rank}_a(L) + q_a)}_{\text{necessary char's}}$

update  $L$ :  $L \leftarrow \max_{a \in \Sigma} \underbrace{\text{select}_a(\text{rank}_a(R) - q_a)}_{\text{unnecessary char's}}$  (no match),

$L \leftarrow L + 1$  (match).

## Jumping algorithm: update rules

$q = (312)$

|          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |
|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|
|          | 0        | 1        | 2        | 3        | 4        | 5        | 6        | 7        | 8        | 9        | 10       | 11       | 12       | 13       | 14       | 15       | 16       | 17       | 18       | 19       | 20       |
|          | <b>b</b> | <b>b</b> | <b>a</b> | <b>c</b> | <b>a</b> | <b>L</b> | <b>c</b> | <b>c</b> | <b>a</b> | <b>b</b> | <b>a</b> | <b>b</b> | <b>b</b> | <b>a</b> | <b>R</b> | <b>b</b> | <b>c</b> | <b>c</b> | <b>a</b> | <b>a</b> | <b>a</b> |
| <b>a</b> | 0        | 0        | 1        | 1        | 2        | 2        | 2        | 2        | 3        | 3        | 4        | 4        | 4        | 5        | 5        | 5        | 5        | 6        | 7        | 8        | 8        |
| <b>b</b> | 1        | 2        | 2        | 2        | 2        | 2        | 2        | 2        | 3        | 3        | 4        | 5        | 5        | 6        | 6        | 6        | 6        | 6        | 6        | 6        | 6        |
| <b>c</b> | 0        | 0        | 0        | 1        | 1        | 2        | 3        | 3        | 3        | 3        | 3        | 3        | 3        | 3        | 3        | 4        | 5        | 5        | 5        | 5        | 6        |

update  $R$ :  $R \leftarrow \max_{a \in \Sigma} \underbrace{\text{select}_a(\text{rank}_a(L) + q_a)}_{\text{necessary char's}}$

update  $L$ :  $L \leftarrow \max_{a \in \Sigma} \underbrace{\text{select}_a(\text{rank}_a(R) - q_a)}_{\text{unnecessary char's}}$  (no match),

$L \leftarrow L + 1$  (match).

## Jumping algorithm: update rules

$q = (312)$

|          |          |          |   |   |   |   |   |   |   |   |          |          |          |    |    |    |    |    |    |    |    |
|----------|----------|----------|---|---|---|---|---|---|---|---|----------|----------|----------|----|----|----|----|----|----|----|----|
|          | 0        | 1        | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10       | 11       | 12       | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |
|          | <b>b</b> | <b>b</b> | a | c | a | c | c | a | b | a | <b>b</b> | <b>a</b> | <b>b</b> | c  | c  | a  | a  | a  | c  |    |    |
| <b>a</b> | 0        | 0        | 1 | 1 | 2 | 2 | 2 | 3 | 3 | 4 | 4        | 4        | 5        | 5  | 5  | 5  | 6  | 7  | 8  | 8  |    |
| <b>b</b> | 1        | 2        | 2 | 2 | 2 | 2 | 2 | 3 | 3 | 3 | 4        | 5        | 5        | 6  | 6  | 6  | 6  | 6  | 6  | 6  |    |
| <b>c</b> | 0        | 0        | 0 | 1 | 1 | 2 | 3 | 3 | 3 | 3 | 3        | 3        | 3        | 3  | 4  | 5  | 5  | 5  | 5  | 6  |    |

update  $R$ :  $R \leftarrow \max_{a \in \Sigma} \underbrace{\text{select}_a(\text{rank}_a(L) + q_a)}_{\text{necessary char's}}$

update  $L$ :  $L \leftarrow \max_{a \in \Sigma} \underbrace{\text{select}_a(\text{rank}_a(R) - q_a)}_{\text{unnecessary char's}}$  (no match),

$L \leftarrow L + 1$  (match).

## Jumping algorithm: update rules

$q = (312)$

|          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |    |   |
|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----|---|
|          | 0        | 1        | 2        | 3        | 4        | 5        | 6        | 7        | 8        | 9        | 10       | 11       | 12       | 13       | 14       | 15       | 16       | 17       | 18       | 19       | 20 |   |
|          | <b>b</b> | <b>b</b> | <b>a</b> | <b>c</b> | <b>a</b> | <b>c</b> | <b>c</b> | <b>a</b> | <b>b</b> | <b>a</b> | <b>b</b> | <b>a</b> | <b>b</b> | <b>c</b> | <b>c</b> | <b>a</b> | <b>a</b> | <b>a</b> | <b>a</b> | <b>c</b> |    |   |
| <b>a</b> | 0        | 0        | 1        | 1        | 2        | 2        | 2        | 3        | 3        | 3        | 4        | 4        | 4        | 4        | 5        | 5        | 5        | 5        | 6        | 7        | 8  | 8 |
| <b>b</b> | 1        | 2        | 2        | 2        | 2        | 2        | 2        | 2        | 3        | 3        | 3        | 4        | 5        | 5        | 6        | 6        | 6        | 6        | 6        | 6        | 6  | 6 |
| <b>c</b> | 0        | 0        | 0        | 1        | 1        | 2        | 3        | 3        | 3        | 3        | 3        | 3        | 3        | 3        | 3        | 4        | 5        | 5        | 5        | 5        | 6  |   |

update  $R$ :  $R \leftarrow \max_{a \in \Sigma} \underbrace{\text{select}_a(\text{rank}_a(L) + q_a)}_{\text{necessary char's}}$

update  $L$ :  $L \leftarrow \max_{a \in \Sigma} \underbrace{\text{select}_a(\text{rank}_a(R) - q_a)}_{\text{unnecessary char's}}$  (no match),

$L \leftarrow L + 1$  (match).

## Jumping algorithm: update rules

$q = (312)$

|          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |    |   |
|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----|---|
|          | 0        | 1        | 2        | 3        | 4        | 5        | 6        | 7        | 8        | 9        | 10       | 11       | 12       | 13       | 14       | 15       | 16       | 17       | 18       | 19       | 20 |   |
|          | <b>b</b> | <b>b</b> | <b>a</b> | <b>c</b> | <b>a</b> | <b>c</b> | <b>c</b> | <b>a</b> | <b>b</b> | <b>a</b> | <b>b</b> | <b>b</b> | <b>a</b> | <b>b</b> | <b>c</b> | <b>c</b> | <b>a</b> | <b>a</b> | <b>a</b> | <b>c</b> |    |   |
| <b>a</b> | 0        | 0        | 1        | 1        | 2        | 2        | 2        | 3        | 3        | 4        | 4        | 4        | 4        | 5        | 5        | 5        | 5        | 5        | 6        | 7        | 8  | 8 |
| <b>b</b> | 1        | 2        | 2        | 2        | 2        | 2        | 2        | 2        | 3        | 3        | 3        | 4        | 5        | 5        | 6        | 6        | 6        | 6        | 6        | 6        | 6  | 6 |
| <b>c</b> | 0        | 0        | 0        | 1        | 1        | 2        | 3        | 3        | 3        | 3        | 3        | 3        | 3        | 3        | 3        | 4        | 5        | 5        | 5        | 5        | 6  |   |

update  $R$ :  $R \leftarrow \max_{a \in \Sigma} \underbrace{\text{select}_a(\text{rank}_a(L) + q_a)}_{\text{necessary char's}}$

update  $L$ :  $L \leftarrow \max_{a \in \Sigma} \underbrace{\text{select}_a(\text{rank}_a(R) - q_a)}_{\text{unnecessary char's}}$  (no match),

$L \leftarrow L + 1$  (match).

## Jumping algorithm: update rules

$q = (312)$

|          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |    |
|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----|
|          | 0        | 1        | 2        | 3        | 4        | 5        | 6        | 7        | 8        | 9        | 10       | 11       | 12       | 13       | 14       | 15       | 16       | 17       | 18       | 19       | 20 |
|          | <b>b</b> | <b>b</b> | <b>a</b> | <b>c</b> | <b>a</b> | <b>c</b> | <b>c</b> | <b>a</b> | <b>b</b> | <b>a</b> | <b>b</b> | <b>b</b> | <b>a</b> | <b>b</b> | <b>c</b> | <b>c</b> | <b>a</b> | <b>a</b> | <b>a</b> | <b>c</b> |    |
| <b>a</b> | 0        | 0        | 1        | 1        | 2        | 2        | 2        | 3        | 3        | 4        | 4        | 4        | 5        | 5        | 5        | 5        | 5        | 6        | 7        | 8        | 8  |
| <b>b</b> | 1        | 2        | 2        | 2        | 2        | 2        | 2        | 2        | 3        | 3        | 4        | 5        | 5        | 6        | 6        | 6        | 6        | 6        | 6        | 6        | 6  |
| <b>c</b> | 0        | 0        | 0        | 1        | 1        | 2        | 3        | 3        | 3        | 3        | 3        | 3        | 3        | 3        | 4        | 5        | 5        | 5        | 5        | 5        | 6  |

update  $R$ :  $R \leftarrow \max_{a \in \Sigma} \underbrace{\text{select}_a(\text{rank}_a(L) + q_a)}_{\text{necessary char's}}$

update  $L$ :  $L \leftarrow \max_{a \in \Sigma} \underbrace{\text{select}_a(\text{rank}_a(R) - q_a)}_{\text{unnecessary char's}}$  (no match),

$L \leftarrow L + 1$  (match).

## Jumping algorithm: update rules

$q = (312)$

|          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |          |    |   |
|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----|---|
|          | 0        | 1        | 2        | 3        | 4        | 5        | 6        | 7        | 8        | 9        | 10       | 11       | 12       | 13       | 14       | 15       | 16       | 17       | 18       | 19       | 20 |   |
|          | <b>b</b> | <b>b</b> | <b>a</b> | <b>c</b> | <b>a</b> | <b>c</b> | <b>c</b> | <b>a</b> | <b>b</b> | <b>a</b> | <b>b</b> | <b>b</b> | <b>a</b> | <b>b</b> | <b>c</b> | <b>c</b> | <b>a</b> | <b>a</b> | <b>a</b> | <b>c</b> |    |   |
| <b>a</b> | 0        | 0        | 1        | 1        | 2        | 2        | 2        | 3        | 3        | 4        | 4        | 4        | 4        | 5        | 5        | 5        | 5        | 5        | 6        | 7        | 8  | 8 |
| <b>b</b> | 1        | 2        | 2        | 2        | 2        | 2        | 2        | 2        | 3        | 3        | 3        | 4        | 5        | 5        | 6        | 6        | 6        | 6        | 6        | 6        | 6  | 6 |
| <b>c</b> | 0        | 0        | 0        | 1        | 1        | 2        | 3        | 3        | 3        | 3        | 3        | 3        | 3        | 3        | 3        | 4        | 5        | 5        | 5        | 5        | 6  |   |

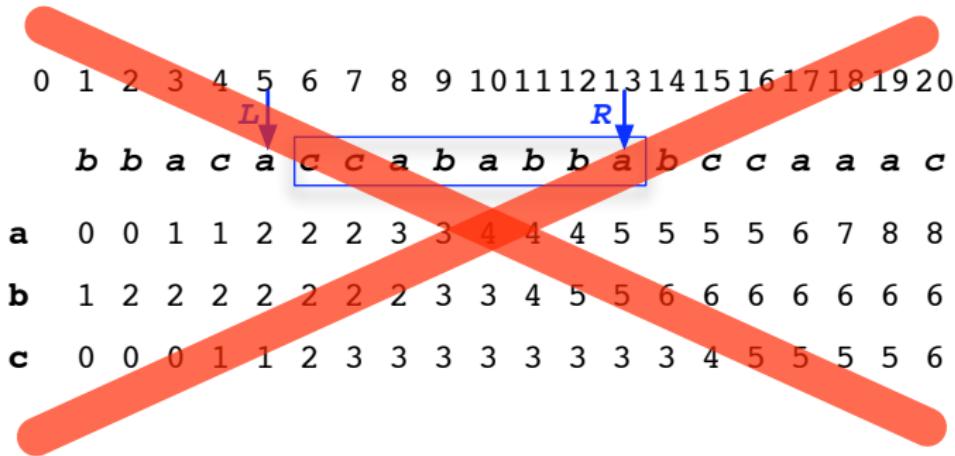
update  $R$ :  $R \leftarrow \max_{a \in \Sigma} \underbrace{\text{select}_a(\text{rank}_a(L) + q_a)}_{\text{necessary char's}}$

update  $L$ :  $L \leftarrow \max_{a \in \Sigma} \underbrace{\text{select}_a(\text{rank}_a(R) - q_a)}_{\text{unnecessary char's}}$  (no match),

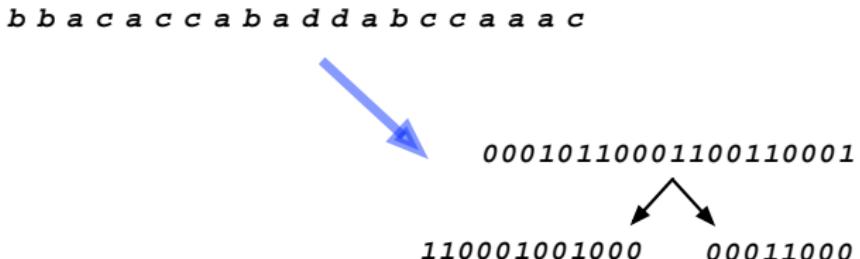
$L \leftarrow L + 1$  (match).

## Jumping algorithm: Analysis

Note that we do not need to store the prefix table, nor the string  $s$ .



## Jumping algorithm: Analysis



Using a wavelet tree [Grossi et al., SODA 2003] as index, we have

- space  $O(n)$  (for wavelet tree)
- construction time  $O(n \log \sigma)$
- every update (jump) in  $O(\sigma)$  time
- query time  $O(J\sigma)$ , where  $J = \text{number of jumps (updates)}$
- expected running time:  $O(n \sqrt{\frac{\sigma}{\log \sigma}} \frac{1}{\sqrt{m}})$ , where  $m = \sum_i q_i$ .

# Jumbled string matching on general alphabets

Current best result: Kociumaka, Radoszewski, Rytter (ESA 2013)

- for arbitrary constant size alphabet
- $o(n^2)$  index size
- $o(n)$  query time (worst-case)
- $O(n^2)$  construction time

(for any  $\delta \in (0, 1)$  construct index of size  $O(n^{2-\delta})$  with query time  $O(m^{\delta(2\sigma-1)})$ )

# Jumbled string matching on general alphabets

Current best result: Kociumaka, Radoszewski, Rytter (ESA 2013)

- for arbitrary constant size alphabet
- $o(n^2)$  index size
- $o(n)$  query time (worst-case)
- $O(n^2)$  construction time

(for any  $\delta \in (0, 1)$  construct index of size  $O(n^{2-\delta})$  with query time  $O(m^{\delta(2\sigma-1)})$ )

Compare to:

1. no index (online):  $O(n)$  space,  $O(n)$  query time
2. store all:  $O(n^2)$  index size,  $O(\log n)$  query time
3. Jumping algo:  $O(n)$  index size,  $o(n)$  query time in expectation

## Jumbled string matching on general alphabets

**Open problem was:** Find something better.

But:

Amir, Chan, Lewenstein, Lewenstein (ICALP 2014) showed that the indexing problem is 3-sum-hard for  $\sigma \geq 3$

More precisely, under different versions of 3-sum-hardness assumption,

- for **non-constant  $\sigma$** , either preprocessing time is  $\Omega(n^{2-\epsilon})$ , or query time is  $\Omega(n^{1-\delta})$ , for all  $\epsilon, \delta$
- for  $\sigma \geq 3$ , same except  $\epsilon, \delta$  depend on  $\sigma$

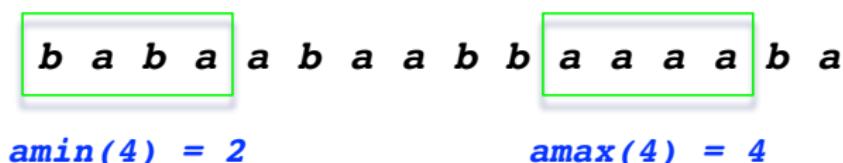
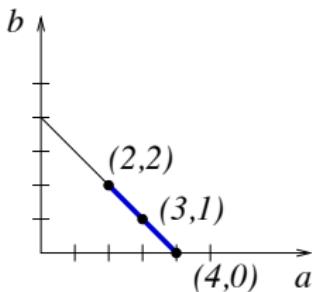
Authors also promised algorithms matching the bounds (forthcoming?)

### 3. Binary alphabets

## Binary alphabets: Interval property

### Lemma

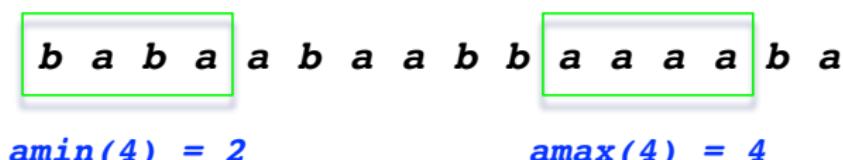
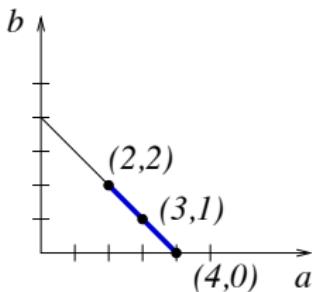
If  $(x, m - x)$  and  $(y, m - y)$  both occur in  $s$ , then so does  $(z, m - z)$  for any  $x \leq z \leq y$ .



## Binary alphabets: Interval property

### Lemma

If  $(x, m - x)$  and  $(y, m - y)$  both occur in  $s$ , then so does  $(z, m - z)$  for any  $x \leq z \leq y$ .



### Corollary

All sub-Parikh vectors of  $s$  of length  $m$  build a set  
 $\{(x, m - x) : \text{amin}(m) \leq x \leq \text{amax}(m)\}$ .

## Binary alphabets: Interval algorithm for **decision** queries

- Index: Table of  $\text{amin}(m)$  and  $\text{amax}(m)$ , for  $1 \leq m \leq n$   
size  $O(n)$
- Query  $(x, y)$  occurs in  $s$  iff  $\text{amin}(x + y) \leq x \leq \text{amax}(x + y)$ .
- Query time  $O(1)$ .

$s = \mathbf{ababbaabaabbbaaaabbab}$

| $m$ | amin | amax |
|-----|------|------|
| ... | ...  | ...  |
| 3   | 0    | 3    |
| 4   | 1    | 3    |
| 5   | 2    | 4    |
| 6   | 2    | 4    |
| ... | ...  | ...  |

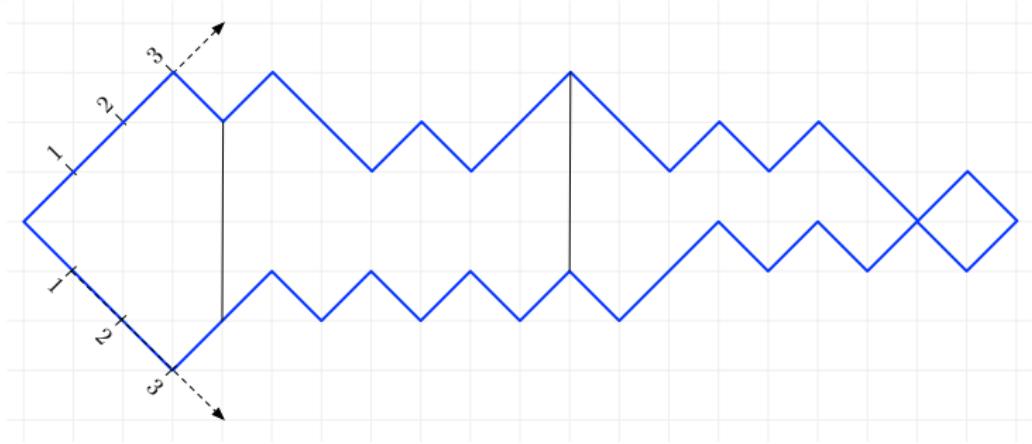
query (3, 2) — YES  
query (1, 4) — NO

# Construction of index

## Goal

Given a binary string of length  $n$ , find, for all  $1 \leq m \leq n$ , the minimum (maximum) number of a's in a window of size  $m$ .

- $O(n^2)$  time—Cicalese, Fici, L. (PSC 2009)
- $O(n^2 / \log n)$  time
  - Burcsi, Cicalese, Fici, L. (FUN 2010); Moosa, Rahman (IPL 2010)
- $O(n^2 / \log^2 n)$  time in word-RAM model
  - Moosa, Rahman (JDA 2012)
- approximate index with one-sided error in  $O(n^{1+\epsilon})$  time
  - Cicalese, Laber, Weimann, Yuster (CPM 2012)
- Corner Index: construction time and index size depend on  $r = \text{runlength enc. of } s$ 
  - Badkobeh, Fici, Kroon, L. (IPL 2013); Giaquinta, Grabowski (IPL 2013)
- $n^2 / 2^{\Omega(\log n / \log \log n)^{1/2}}$ —Hermelin, Landau, Rabinovich, Weimann (unpubl.)



| $m$ | amin | amax |
|-----|------|------|
| ... | ...  | ...  |
| 3   | 0    | 3    |
| 4   | 1    | 3    |
| 5   | 2    | 4    |
| 6   | 2    | 4    |
| ... | ...  | ...  |

The Parikh set of  
 $s = ababbaabaabbbaaababb.$   
 Verticals are fixed length sub-  
 Parikh vectors, for  $m = 4, 11$

# Construction of index

## Goal

Given a binary string of length  $n$ , find, for all  $1 \leq m \leq n$ , the minimum (maximum) number of a's in a window of size  $m$ .

- $O(n^2)$  time—Cicalese, Fici, L. (PSC 2009)
- $O(n^2 / \log n)$  time
  - Burcsi, Cicalese, Fici, L. (FUN 2010); Moosa, Rahman (IPL 2010)
- $O(n^2 / \log^2 n)$  time in word-RAM model
  - Moosa, Rahman (JDA 2012)
- approximate index with one-sided error in  $O(n^{1+\epsilon})$  time
  - Cicalese, Laber, Weimann, Yuster (CPM 2012)
- Corner Index: construction time and index size depend on  $r = \text{runlength enc. of } s$ 
  - Badkobeh, Fici, Kroon, L. (IPL 2013); Giaquinta, Grabowski (IPL 2013)
- $n^2 / 2^{\Omega(\log n / \log \log n)^{1/2}}$ —Hermelin, Landau, Rabinovich, Weimann (unpubl.)

## Construction of index

**Open problem:** Faster construction of binary index?

## 4. Prefix normal words

# Prefix normal words

Fici, L. (DLT 2011)

## Definition

A word  $s \in \{a, b\}^*$  is a **prefix normal word** (w.r.t.  $a$ ) if  $\forall 0 \leq m \leq |s|$  no substring of length  $m$  has more  $a$ 's than the prefix of  $s$  of length  $m$ .

## Example

$$s = ababbaabaabbbaaaabbab$$

$$s' = aaababbabaabbababbab$$

# Prefix normal words

Fici, L. (DLT 2011)

## Definition

A word  $s \in \{a, b\}^*$  is a **prefix normal word** (w.r.t.  $a$ ) if  $\forall 0 \leq m \leq |s|$  no substring of length  $m$  has more  $a$ 's than the prefix of  $s$  of length  $m$ .

## Example

$s = ababbaabaabbbaaaabbab$     NO

$s' = aaababbabaabbabababbab$     YES

## Prefix normal forms

Recall  $\text{amax}_a(s, m) = \text{maximum number of } a\text{'s in a substring of } s \text{ of length } m$ .

$$s = ababbaabaabbbaaaabbab$$

| $m$             | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |
|-----------------|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|
| $\text{amax}_a$ | 0 | 1 | 2 | 3 | 3 | 4 | 4 | 4 | 5 | 5 | 6  | 7  | 7  | 7  | 8  | 8  | 9  | 9  | 9  | 10 | 10 |

### Theorem

Let  $s \in \{a, b\}^*$ . Then there exists a unique prefix normal word  $s'$  s.t. for all  $0 \leq m \leq |s|$ ,  $\text{amax}_a(s, m) = \text{amax}_a(s', m)$ , called its **prefix normal form w.r.t. a**,  $\text{PNF}_a(s)$ .

# Prefix normal forms

$$s = ababbaabaabbbaaaabbab$$

| $m$         | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |
|-------------|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|
| $a_{max_a}$ | 0 | 1 | 2 | 3 | 3 | 4 | 4 | 4 | 5 | 5 | 6  | 7  | 7  | 7  | 8  | 8  | 9  | 9  | 9  | 10 | 10 |

## Prefix normal forms

$$s = ababbaabaabbbaaaabbab$$

| $m$       | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |
|-----------|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|
| $a\max_a$ | 0 | 1 | 2 | 3 | 3 | 4 | 4 | 4 | 5 | 5 | 6  | 7  | 7  | 7  | 8  | 8  | 9  | 9  | 9  | 10 | 10 |

Define the word  $s'$  by

$$s'_m = \begin{cases} a & \text{if } a\max_a(s, m) = 1 + a\max_a(s, m - 1) \\ b & \text{if } a\max_a(s, m) = a\max_a(s, m - 1) \end{cases}$$

## Prefix normal forms

$$s = ababbaabaabbbaaaabbab$$

| $m$       | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |
|-----------|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|
| $a\max_a$ | 0 | 1 | 2 | 3 | 3 | 4 | 4 | 4 | 5 | 5 | 6  | 7  | 7  | 7  | 8  | 8  | 9  | 9  | 9  | 10 | 10 |

Define the word  $s'$  by

$$s'_m = \begin{cases} a & \text{if } a\max_a(s, m) = 1 + a\max_a(s, m - 1) \\ b & \text{if } a\max_a(s, m) = a\max_a(s, m - 1) \end{cases}$$

## Prefix normal forms

$$s = ababbaabaabbbaaaabbab$$

| $m$       | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |
|-----------|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|
| $a\max_a$ | 0 | 1 | 2 | 3 | 3 | 4 | 4 | 4 | 5 | 5 | 6  | 7  | 7  | 7  | 8  | 8  | 9  | 9  | 9  | 10 | 10 |

Define the word  $s'$  by

$$s'_m = \begin{cases} a & \text{if } a\max_a(s, m) = 1 + a\max_a(s, m - 1) \\ b & \text{if } a\max_a(s, m) = a\max_a(s, m - 1) \end{cases}$$

$$s' = \textcolor{red}{a}$$

## Prefix normal forms

$$s = ababbaabaabbbaaaabbab$$

| $m$       | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |
|-----------|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|
| $a\max_a$ | 0 | 1 | 2 | 3 | 3 | 4 | 4 | 4 | 5 | 5 | 6  | 7  | 7  | 7  | 8  | 8  | 9  | 9  | 9  | 10 | 10 |

Define the word  $s'$  by

$$s'_m = \begin{cases} a & \text{if } a\max_a(s, m) = 1 + a\max_a(s, m - 1) \\ b & \text{if } a\max_a(s, m) = a\max_a(s, m - 1) \end{cases}$$

$$s' = a\textcolor{red}{a}$$

## Prefix normal forms

$$s = ababbaabaabbbaaaabbab$$

| $m$       | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |
|-----------|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|
| $a\max_a$ | 0 | 1 | 2 | 3 | 3 | 4 | 4 | 4 | 5 | 5 | 6  | 7  | 7  | 7  | 8  | 8  | 9  | 9  | 9  | 10 | 10 |

Define the word  $s'$  by

$$s'_m = \begin{cases} a & \text{if } a\max_a(s, m) = 1 + a\max_a(s, m - 1) \\ b & \text{if } a\max_a(s, m) = a\max_a(s, m - 1) \end{cases}$$

$$s' = aaa$$

## Prefix normal forms

$$s = ababbaabaabbbaaaabbab$$

| $m$       | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |
|-----------|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|
| $a\max_a$ | 0 | 1 | 2 | 3 | 3 | 4 | 4 | 4 | 5 | 5 | 6  | 7  | 7  | 7  | 8  | 8  | 9  | 9  | 9  | 10 | 10 |

Define the word  $s'$  by

$$s'_m = \begin{cases} a & \text{if } a\max_a(s, m) = 1 + a\max_a(s, m - 1) \\ b & \text{if } a\max_a(s, m) = a\max_a(s, m - 1) \end{cases}$$

$$s' = aaa\color{red}{b}$$

## Prefix normal forms

$$s = ababbaabaabbbaaaabbab$$

| $m$       | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |
|-----------|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|
| $a\max_a$ | 0 | 1 | 2 | 3 | 3 | 4 | 4 | 4 | 5 | 5 | 6  | 7  | 7  | 7  | 8  | 8  | 9  | 9  | 9  | 10 | 10 |

Define the word  $s'$  by

$$s'_m = \begin{cases} a & \text{if } a\max_a(s, m) = 1 + a\max_a(s, m - 1) \\ b & \text{if } a\max_a(s, m) = a\max_a(s, m - 1) \end{cases}$$

$$s' = aaababbabaabbababbab$$

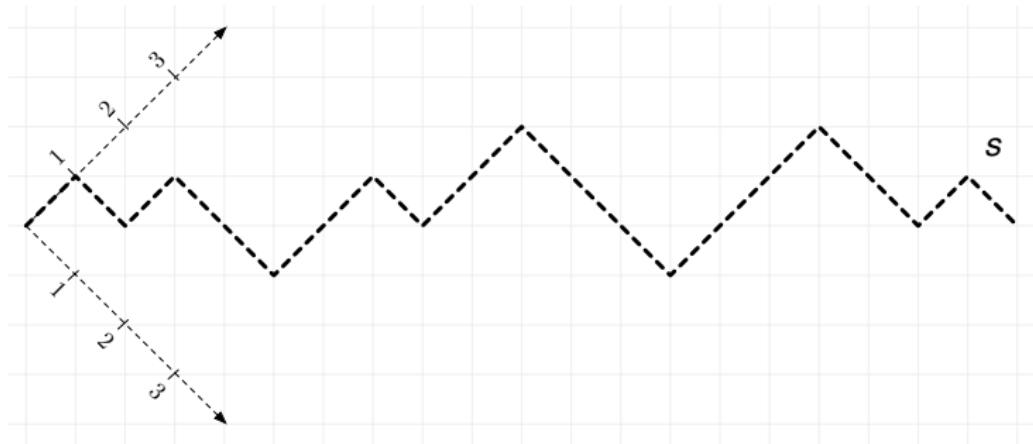
## BJPM with prefix normal forms

### Theorem

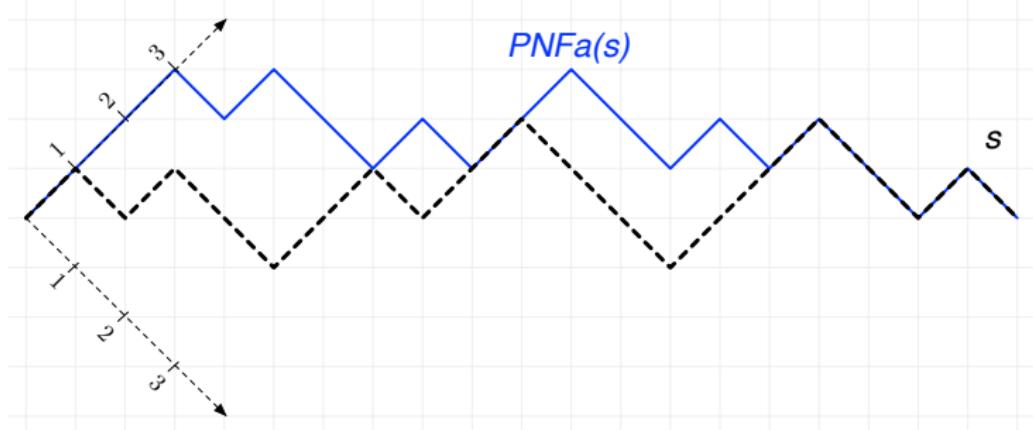
Two strings  $s, t \in \{a, b\}^*$  have the same Parikh set if and only if  $\text{PNF}_a(s) = \text{PNF}_a(t)$  and  $\text{PNF}_b(s) = \text{PNF}_b(t)$ .

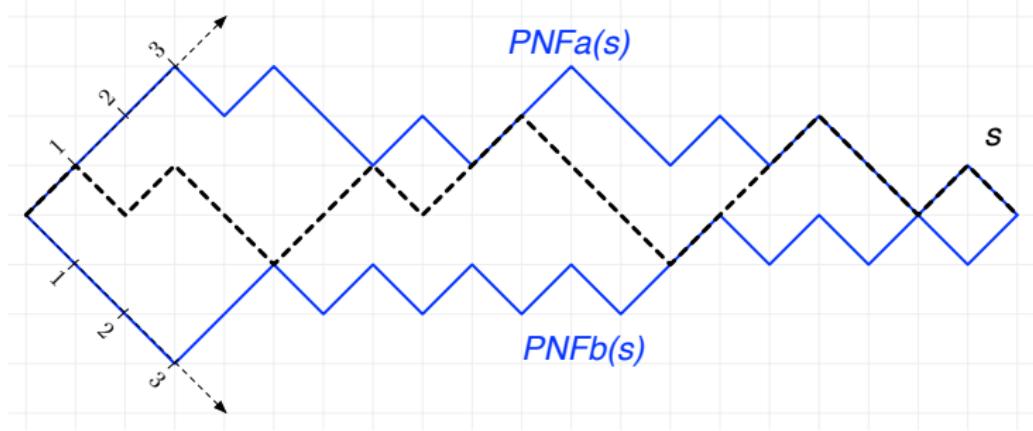
(Parikh set of a word  $s$  is the set of Parikh vectors of the substrings of  $s$ .  
BJPM = binary jumbled pattern matching)

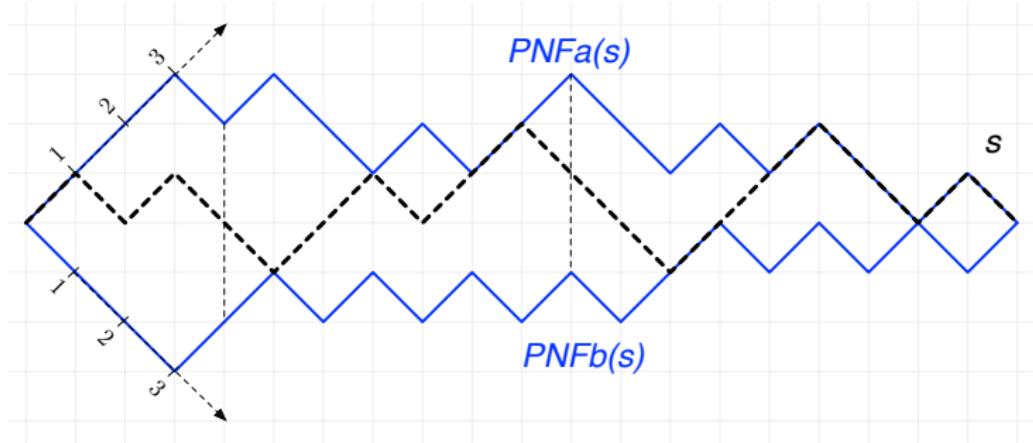
$$s = ababbaabaabbbaaaabbab$$



$$\swarrow = a, \searrow = b,$$

$s = ababbaabaabbbaaaabbab$  $\swarrow = a, \searrow = b,$

$s = ababbaabaabbbaaaabbab$  $\diagup = a, \diagdown = b,$

$s = ababbaabaabbbaaabbaab$  $\diagup = a, \diagdown = b,$

## BJPM with prefix normal forms

Does  $s = \mathbf{ababbaabaabbbaaabab}$  have a substring of length 11 containing exactly 5 a's?

# a's in  $\text{pref}(\text{PNF}_a(s), 11) = 7$

# a's in  $\text{pref}(\text{PNF}_b(s), 11) = 5$

$$7 \geq \mathbf{5} \geq 5 \rightsquigarrow \text{YES}$$

Thus, fast computation of PNFs yields fast solution to BJPM.

## BJPM with prefix normal forms

Does  $s = \mathbf{ababbaabaabbbaaaabbab}$  have a substring of length 11 containing exactly 5 a's?

# a's in  $\text{pref}(\text{PNF}_a(s), 11) = 7$

# a's in  $\text{pref}(\text{PNF}_b(s), 11) = 5$

$$7 \geq \underline{5} \geq 5 \rightsquigarrow \text{YES}$$

Thus, fast computation of PNFs yields fast solution to BJPM.  
Better understanding p.n. words might be useful for BJPM.

## BJPM with prefix normal forms

Does  $s = \mathbf{ababbaabaabbbaaabab}$  have a substring of length 11 containing exactly 5 a's?

# a's in  $\text{pref}(\text{PNF}_a(s), 11) = 7$

# a's in  $\text{pref}(\text{PNF}_b(s), 11) = 5$

$$7 \geq \underline{5} \geq 5 \rightsquigarrow \text{YES}$$

Thus, fast computation of PNFs yields fast solution to BJPM.  
Better understanding p.n. words might be useful for BJPM.  
Or might just be fun.

## Computation of PNFs

**Open problem:** Compute the PNFs fast. (Many other open problems on PNFs, e.g. testing.)

### Lit:

- Fici, L. (DLT 2011)
- Burcsi, Fici, L., Ruskey, Sawada (CPM 2014): combinatorial generation of all p.n. words, based on them being a **bubble language**.
- Burcsi, Fici, L., Ruskey, Sawada (FUN 2014): partial results on enumeration of p.n. words
- OEIS seq. no. A194850: no. of p.n. words of length  $n$
- OEIS seq. no. A238109: size of largest equivalence class for length  $n$

# Conclusion

We saw: JPM on

1. general alphabets (constant-size)
  - expected: jumping algorithm—sublinear expected time, linear space,
  - worst-case: Kociumaka et al—just sublinear time and just subquadratic space (tradeoff)
2. binary alphabets (decision queries): linear index, main question is how to compute it
3. prefix normal words / prefix normal forms

# Conclusion

We saw: JPM on

1. general alphabets (constant-size)
  - expected: jumping algorithm—sublinear expected time, linear space,
  - worst-case: Kociumaka et al—just sublinear time and just subquadratic space (tradeoff)
2. binary alphabets (decision queries): linear index, main question is how to compute it
3. prefix normal words / prefix normal forms

Take home message

1. Jumbled pattern matching is **interesting**.
2. It is definitely **hype**.
3. Prefix normal words are **cool!**

# AHKNOTUY!

zsuzsanna.liptak@univr.it

# APPENDIX

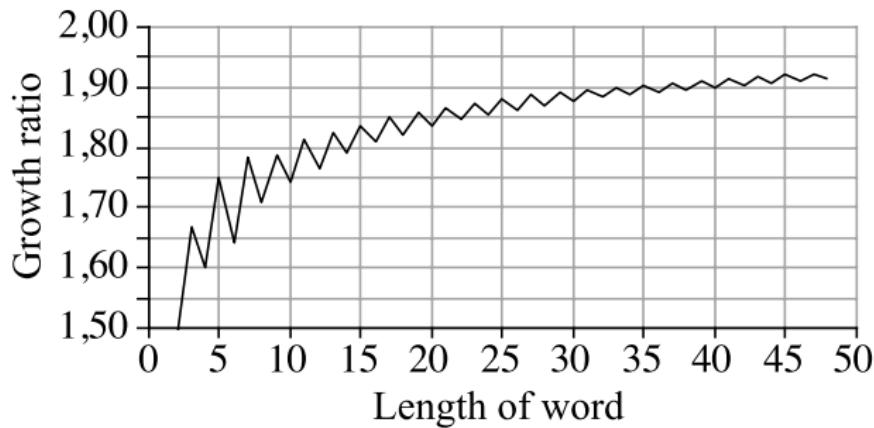
## Other variants

- approximate jumbled string matching: e.g. find all occurrences between (12, 5, 7) and (8, 2, 4): variant of jumping algorithm expected sublinear (Burcsi, Cicalese, Fici, L., FUN 2010, ToCS 2012)
- (mostly binary) on **trees** and **graphs with bounded treewidth** (Gagie, Hermelin, Landau, Weimann, ESA 2013)
- (mostly binary) locating one occurrence in **strings, trees, graphs** (Cicalese, Gagie, Giaquinta, Laber, L., Rizzi, Tomescu, SPIRE 2013)
- Jumbled matching on **color point sets** in the plane (Gagie et al: CPM 2014)
- binary **locating** queries: included in (Gagie, Hermelin, Landau, Weimann, ESA 2013)

## Enumerating p.n. words

OEIS sequence no. A194850

| $n$         | 1 | 2 | 3 | 4 | 5  | 6  | 7  | 8  | 9   | 10  | 11  | 12  | 13   | 14   | 15   | 16   |
|-------------|---|---|---|---|----|----|----|----|-----|-----|-----|-----|------|------|------|------|
| $p_{nw}(n)$ | 2 | 3 | 5 | 8 | 14 | 23 | 41 | 70 | 125 | 218 | 395 | 697 | 1273 | 2279 | 4185 | 7568 |



The prefix normal form  $\text{PNF}_a$  induces an equivalence relation on  $\Sigma^*$ , namely  $u \equiv_{\text{PNF}_a} v$  if and only if  $\text{PNF}_a(u) = \text{PNF}_a(v)$ .

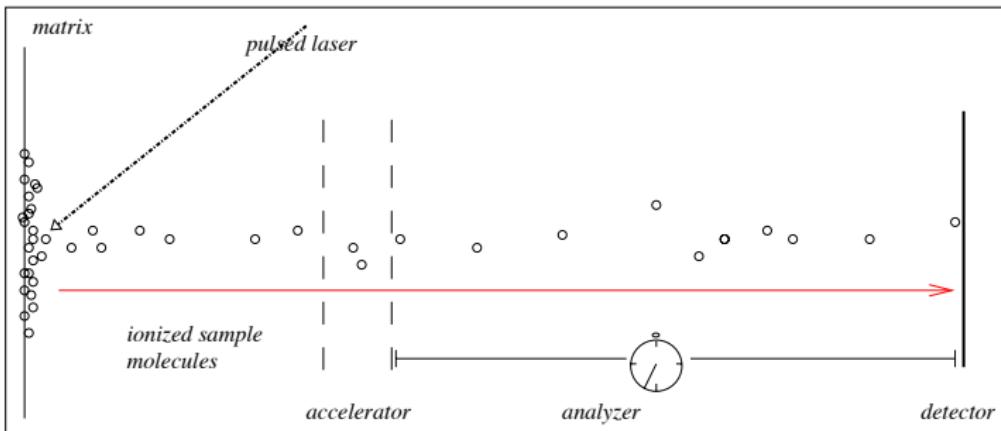
| $\text{PNF}_a$ | class                        | card. |
|----------------|------------------------------|-------|
| $aaaa$         | $\{aaaa\}$                   | 1     |
| $aaab$         | $\{aaab, baaa\}$             | 2     |
| $aaba$         | $\{aaba, abaa\}$             | 2     |
| $aabb$         | $\{aabb, baab, bbaa\}$       | 3     |
| $abab$         | $\{abab, baba\}$             | 2     |
| $abba$         | $\{abba\}$                   | 1     |
| $abbb$         | $\{abbb, babb, bbab, bbba\}$ | 4     |
| $bbbb$         | $\{bbbb\}$                   | 1     |

Table : The classes of words of length 4 having the same prefix normal form.

| PNF <sub>a</sub> | card. |
|------------------|-------|------------------|-------|------------------|-------|------------------|-------|
| aaaaaaaa         | 1     | aaabaabb         | 6     | aabababa         | 6     | abababba         | 2     |
| aaaaaaaab        | 2     | aaababaa         | 2     | aabababb         | 9     | abababbb         | 4     |
| aaaaaaaba        | 2     | aaababab         | 6     | aababbaa         | 2     | ababbaba         | 1     |
| aaaaaaabb        | 3     | aaababba         | 4     | aababbab         | 8     | ababbabb         | 6     |
| aaaaaabaa        | 2     | aaababbb         | 8     | aababbba         | 4     | ababbbab         | 4     |
| aaaaaabab        | 4     | aaabbaaa         | 1     | aababbbb         | 10    | ababbbba         | 2     |
| aaaaabba         | 2     | aaabbaab         | 4     | aabbaabb         | 3     | ababbbbb         | 6     |
| aaaaabbb         | 4     | aaabbaba         | 2     | aabbabab         | 4     | abbabbab         | 2     |
| aaaabaaa         | 2     | aaabbabb         | 6     | aabbabba         | 3     | abbabbba         | 2     |
| aaaabaab         | 4     | aaabbbaa         | 2     | aabbabbb         | 8     | abbabbbb         | 5     |
| aaaababa         | 3     | aaabbbab         | 4     | aabbbaab         | 2     | abbbabbb         | 4     |
| aaaababb         | 6     | aaabbbba         | 2     | aabbabba         | 2     | abbbbabb         | 3     |
| aaaabbaa         | 2     | aaabbbbb         | 6     | aabbabb          | 6     | abbbbabb         | 2     |
| aaaabbab         | 4     | aabaabaa         | 1     | aabbbaa          | 1     | abbbbbaa         | 1     |
| aaaabbba         | 2     | aabaabab         | 4     | aabbbbab         | 4     | abbbbbb          | 8     |
| aaaabbbb         | 5     | aabaabba         | 2     | aabbbbba         | 2     | bbbbbbb          | 1     |
| aaabaaaab        | 2     | aabaabbb         | 4     | aabbbbbb         | 7     |                  |       |
| aaabaaba         | 4     | aababaab         | 2     | abababab         | 2     |                  |       |

Table : The cardinalities of the 70 classes of words of length 8 having the same prefix normal form.

# What is mass spectrometry?



example: MALDI-TOF mass spectrometer

## Reduction to min-plus-convolution

### Min-plus-convolution

For real vectors  $x = (x_0, \dots, x_n)$  and  $y = (y_0, \dots, y_n)$ , define  $x \star y = z$  by

$$z_k = \min_i \{x_i + y_{k-i}\}, \quad k = 0, 1, \dots, 2n,$$

(replace *min* by *sum* and  $+$  by  $\cdot$  to get classic convolution).

# Reduction to min-plus-convolution

## Min-plus-convolution

For real vectors  $x = (x_0, \dots, x_n)$  and  $y = (y_0, \dots, y_n)$ , define  $x \star y = z$  by

$$z_k = \min_i \{x_i + y_{k-i}\}, \quad k = 0, 1, \dots, 2n,$$

## Reduction

Set  $x_i := |\text{pref}_i(s)|_a$ , the number of  $a$ 's in the prefix of  $s$  of length  $i$ .

Set  $y_i := -x_{n-i}$ ,  $i = 0, 1, \dots, n$ . Then

$$\text{amin}(m) = \min_{i=0}^{n-m} (x_{i+m} - x_i) = \min_{i=0}^{n-m} (x_{i+m} + y_{n-i}) = z_{n+m}.$$

## Jumping algo analysis: Estimating no. of jumps

We compute the expectation of  $J$  over a random string  $s$  (**uniform i.i.d.**) and **balanced** query  $q = (\frac{m}{\sigma}, \dots, \frac{m}{\sigma})$  as

$$\mathbb{E}(J) = \frac{n}{\mathbb{E}(\text{length of jump})}, \text{ and}$$

$$\mathbb{E}(\text{length of jump}) = \Pr(\text{match}) \cdot 1 + \Pr(\text{no match}) \cdot (\mathbb{E}(R_{\text{new}} - L) - m).$$

## Jumping algo analysis: Estimating no. of jumps

We compute the expectation of  $J$  over a random string  $s$  (**uniform i.i.d.**) and **balanced** query  $q = (\frac{m}{\sigma}, \dots, \frac{m}{\sigma})$  as

$$\mathbb{E}(J) = \frac{n}{\mathbb{E}(\text{length of jump})}, \text{ and}$$

$$\mathbb{E}(\text{length of jump}) = \Pr(\text{match}) \cdot 1 + \Pr(\text{no match}) \cdot (\mathbb{E}(R_{\text{new}} - L) - m).$$

- $\mathbb{E}(R_{\text{new}} - L) = \text{length of wait before we have seen } \frac{m}{\sigma} \text{ many of all } \sigma \text{ characters} \approx \sqrt{2m\sigma \ln \frac{\sigma}{\sqrt{2\pi}}} + m$   
(R. May, Coupon collecting with quotas, Electr. J Comb., 2008)
- $\Pr(\text{match}) = \Pr(\text{random string of size } m \text{ has Parikh vector } q) = \frac{\binom{m}{\frac{m}{\sigma}, \dots, \frac{m}{\sigma}}}{m^\sigma} \approx \sqrt{\frac{\sigma^\sigma}{m^{\sigma-1}}} \quad \text{very small for } m \gg \sigma$

## Jumping algo analysis: Estimating no. of jumps

We compute the expectation of  $J$  over a random string  $s$  (**uniform i.i.d.**) and **balanced** query  $q = (\frac{m}{\sigma}, \dots, \frac{m}{\sigma})$  as

$$\mathbb{E}(J) = \frac{n}{\mathbb{E}(\text{length of jump})}, \text{ and}$$

$$\begin{aligned}\mathbb{E}(\text{length of jump}) &= \Pr(\text{match}) \cdot 1 + \Pr(\text{no match}) \cdot (\mathbb{E}(R_{\text{new}} - L) - m) \\ &\approx (\mathbb{E}(R_{\text{new}} - L) - m) \approx \sqrt{2m\sigma \ln \frac{\sigma}{\sqrt{2\pi}}}\end{aligned}$$

So,

$$\mathbb{E}(J) = O\left(\frac{n}{\sqrt{m}\sqrt{\sigma \log \sigma}}\right).$$

## Jumping algo analysis: Average case

Altogether, we get

$$\mathbb{E}(\text{runtime}) = O(J\sigma) = O\left(\frac{n\sigma}{\sqrt{m}\sqrt{\sigma \log \sigma}}\right) = O(n\sqrt{\frac{\sigma}{\log \sigma}}\frac{1}{\sqrt{m}}).$$

Recall: We assume a **random string  $s$  (uniform i.i.d.)** and **balanced query  $q = (\frac{m}{\sigma}, \dots, \frac{m}{\sigma})$** .

These are worst-case, both acc. to our model and our simulations:

- on non-random strings (DNA)  $J$  decreases
- on non-balanced Parikh vectors  $J$  decreases