

Wiktor Zuba

MIMUW Colloquium

13.03.2025

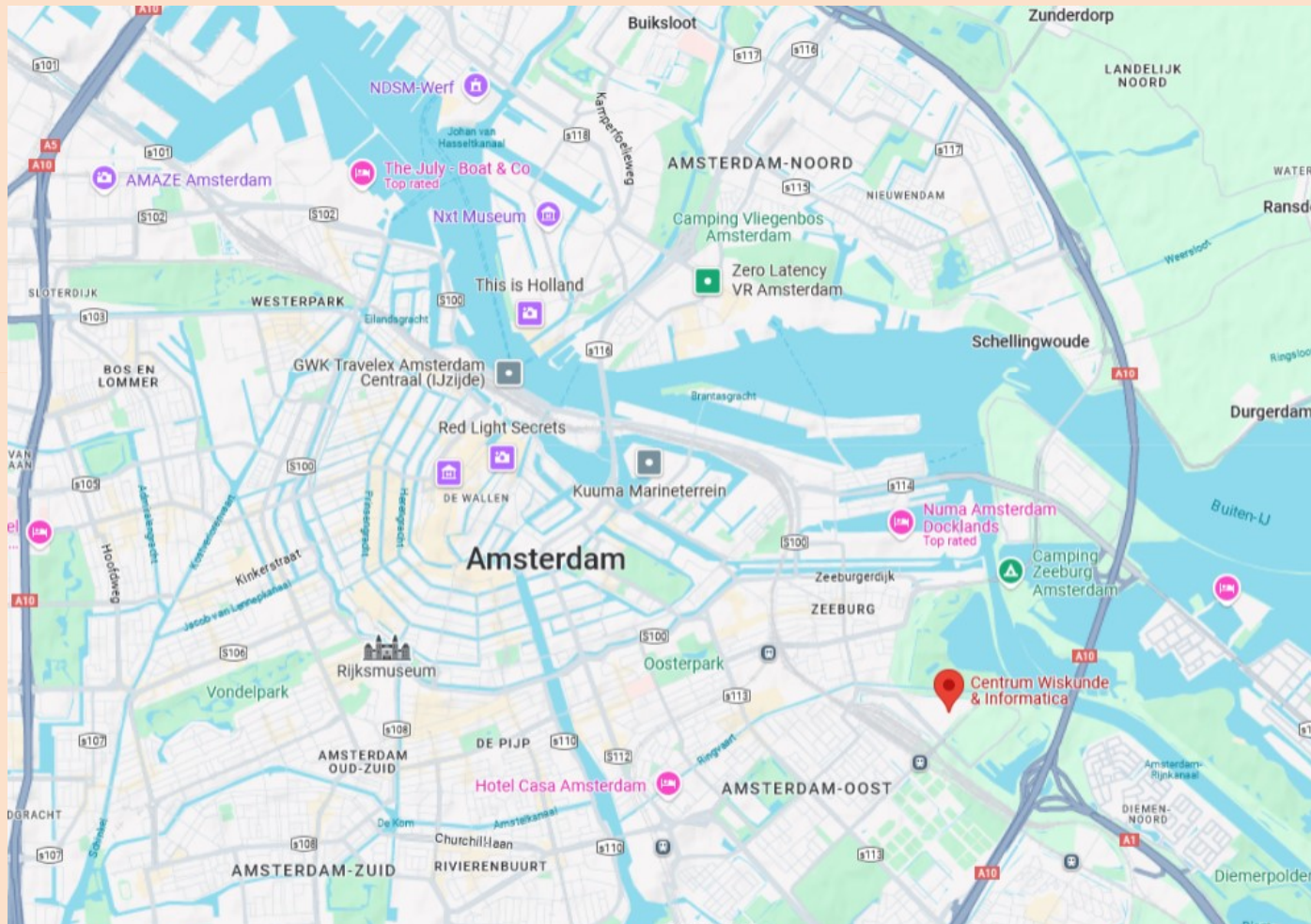
# Background

- 2011-2021: Studies at MIMUW
- 2022-2024: PostDoc at CWI
- 2024-present: Assistant Professor at MIMUW

# CWI



# CWI





# CWI

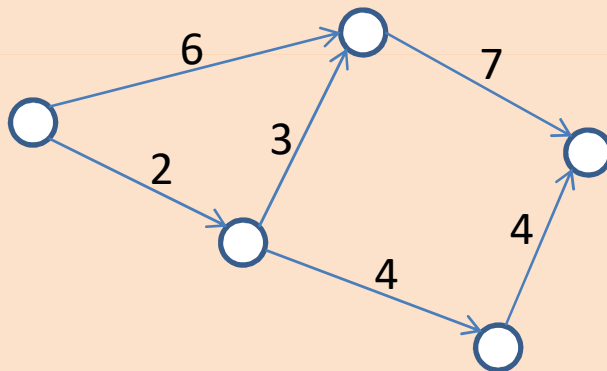




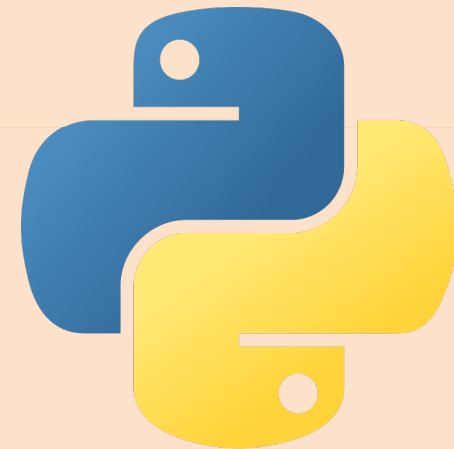
CWI



Dijkstra Algorithm 1956



Python 1980s



Now -



Text Algorithms Team



Now -



## Text Algorithms Team



Wojciech Rytter



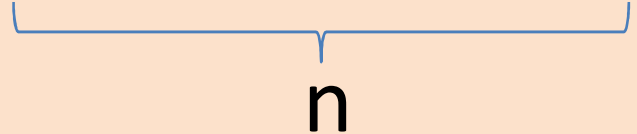
Jakub Radoszewski



Tomasz Waleń

# Text Algorithms

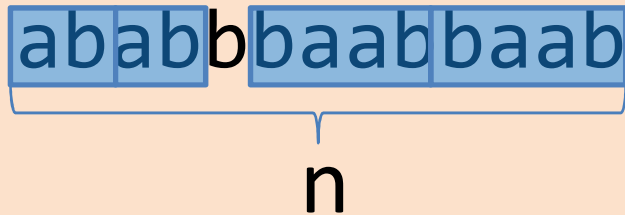
ababbbaabbaab



n

(Check / find all / count all) regular parts

# Text Algorithms

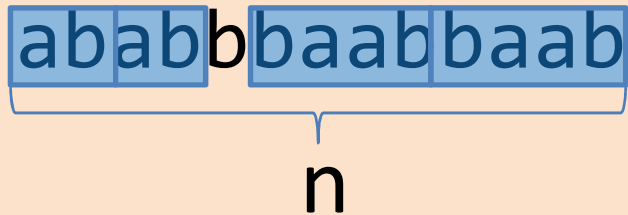


ababbabaabbaab

n

(Check / find all / count all) regular parts

# Text Algorithms



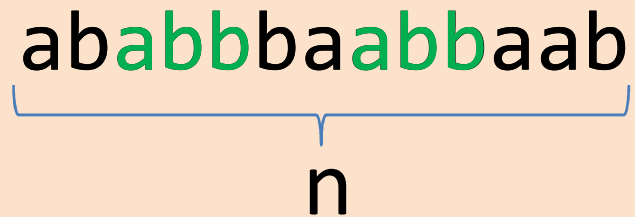
The diagram shows the string "ababbabaabbaab" where the first 10 characters "ababbabaab" are grouped by a blue bracket underneath. Below the bracket is the letter "n", indicating that the substring of length n is "ababbabaab".

(Check / find all / count all) regular parts  
(in a substring)



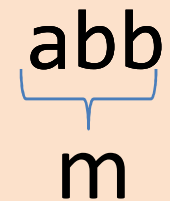
# Text Algorithms

ababbbaabbaab



n

abb



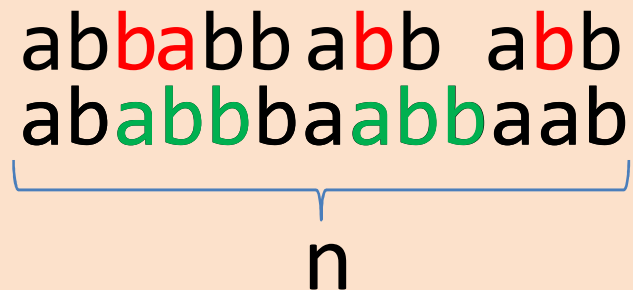
m

(Check / find all / count all) regular parts  
(in a substring)

Pattern matching

# Text Algorithms

ab**ba**bbab**ba**bb  
ab**abb**ba**abb**aab

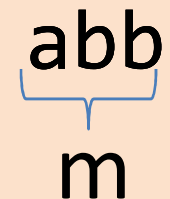


A horizontal blue bracket is positioned below the two lines of text, spanning the entire length of the string. Below the bracket is the letter 'n'.

n

(Check / find all / count all) regular parts  
(in a substring)

abb



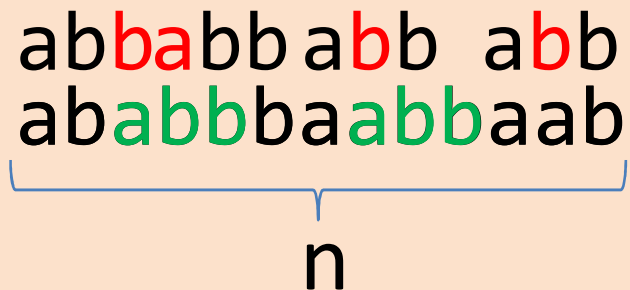
A horizontal blue bracket is positioned below the string 'abb'. Below the bracket is the letter 'm'.

m

(Approximate = up to  
k errors)  
Pattern matching

# Text Algorithms

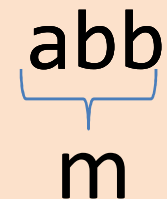
ab**ba**bbab**ba**bb  
ab**abb**ba**abb**aab



n

(Check / find all / count all) regular parts  
(in a substring)

abb



m

(Approximate = up to  
k errors)

Pattern matching

String similarities  
(common substrings, ...)

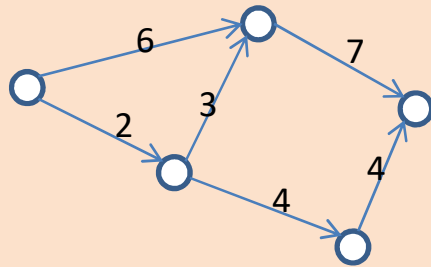
# Applications of regularities and similarities of strings

- **Bioinformatics** – tandem repeats are associated to genetic diseases,  
similarities between DNA sequences ~ close relationship between organisms
- **Compression** – high regularity = better compression rate
- **Data analysis** – plagiarism detection



# Aims

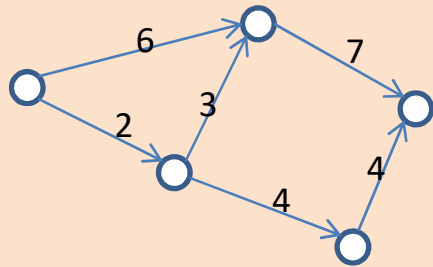
## Graph algorithms



brute force:  
usually exponential

# Aims

## Graph algorithms



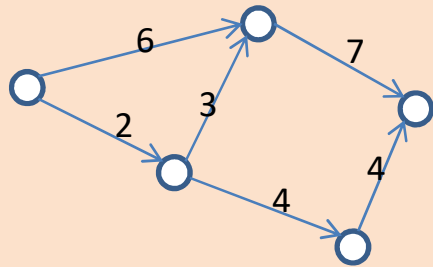
brute force:  
usually exponential

Is there an  $n^{O(1)}$ -time algorithm?

NP-complete

# Aims

## Graph algorithms



brute force:  
usually exponential

Is there an  $n^{O(1)}$ -time algorithm?

NP-complete

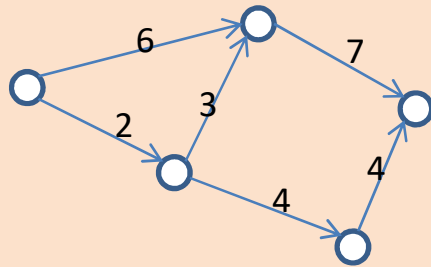
Smaller exponent?

Approximation?

FPT?

# Aims

## Graph algorithms



brute force:  
usually exponential

Is there an  $n^{O(1)}$ -time algorithm?

NP-complete

Smaller exponent?

Approximation?

FPT?

## Text algorithms

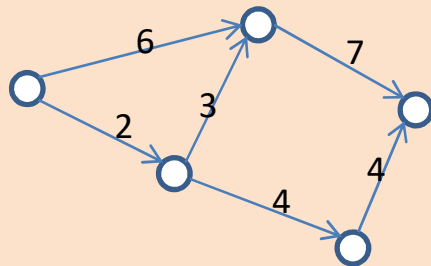
ababbbaabbaab

brute force:  
usually  $O(n^3)$  or  $O(n^2)$  time.



# Aims

## Graph algorithms



brute force:  
usually exponential

Is there an  $n^{O(1)}$ -time algorithm?

NP-complete

Smaller exponent?

Approximation?

FPT?

## Text algorithms

ababbbaabbaab

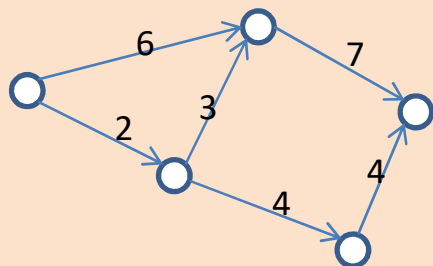
brute force:  
usually  $O(n^3)$  or  $O(n^2)$  time.

Try to reach  $O(n)$  or at least  $O(n \text{ polylog } n)$ .

No  $O(n^{2-\epsilon})$  time algorithm  
(by e.g. SETH)

# Aims

## Graph algorithms



brute force:  
usually exponential

Is there an  $n^{O(1)}$ -time algorithm?

NP-complete

Smaller exponent?

Approximation?

FPT?

## Text algorithms


ababbbaabbaab

brute force:  
usually  $O(n^3)$  or  $O(n^2)$  time. Sometimes even  $O(n/\log n)$ .

Try to reach  $O(n)$  or at least  $O(n \text{ polylog } n)$ .

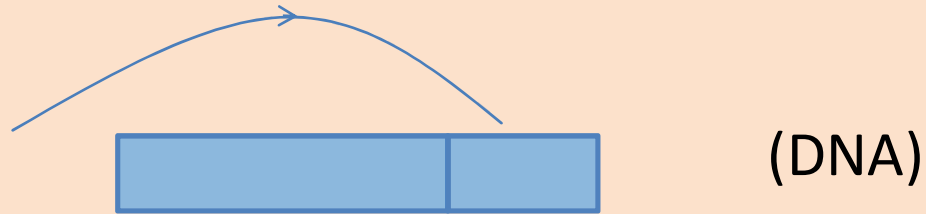
No  $O(n^{2-\epsilon})$  time algorithm  
(by e.g. SETH)

# Different models

- Circular 
- Approximate
- Abelian
- Order preserving
- 2D
- Labelled graphs

# Different models

- Circular
- Approximate
- Abelian
- Order preserving
- 2D
- Labelled graphs





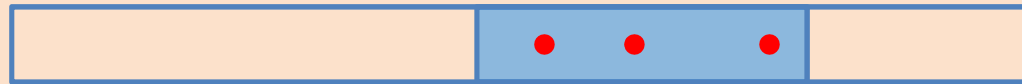
# Different models

- Circular



(DNA)

- Approximate




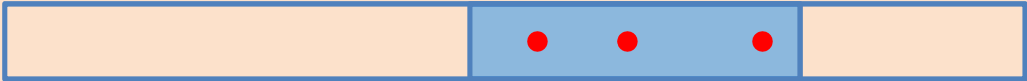
- Abelian

- Order preserving

- 2D

- Labelled graphs

# Different models

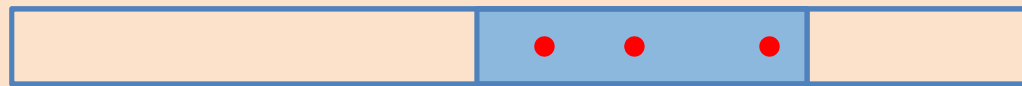
- Circular  (DNA)
- Approximate 
- Abelian     ababbbaabbaab~ aaaaaabbbbbbbb
- Order preserving
- 2D
- Labelled graphs

# Different models

- Circular

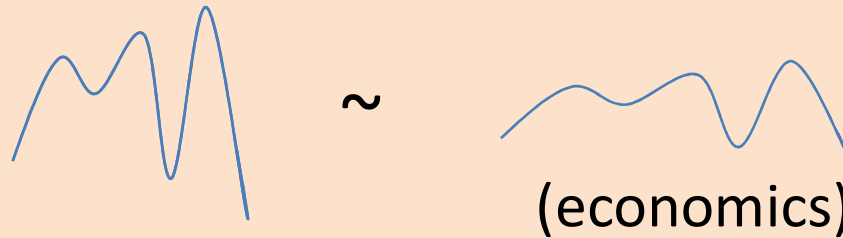


- Approximate



- Abelian     ababbbaabbaab~ aaaaaabbbbbbbb

- Order preserving



- 2D

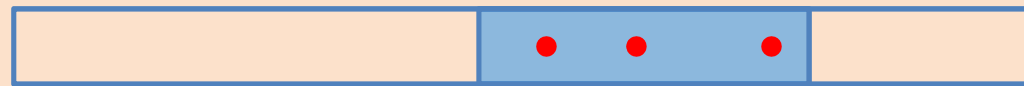
- Labelled graphs

# Different models

- Circular



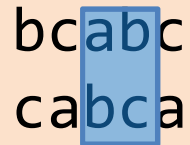
- Approximate



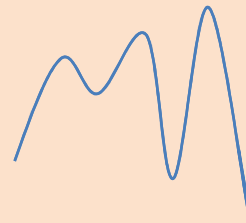
- Abelian     ababbbaabbaab~ aaaaaabbbbbbbb

- Order preserving

- 2D     abcab  
         bcabc  
         cabca



(image processing)



~



(economics)

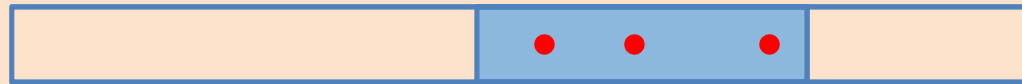
- Labelled graphs

# Different models

- Circular



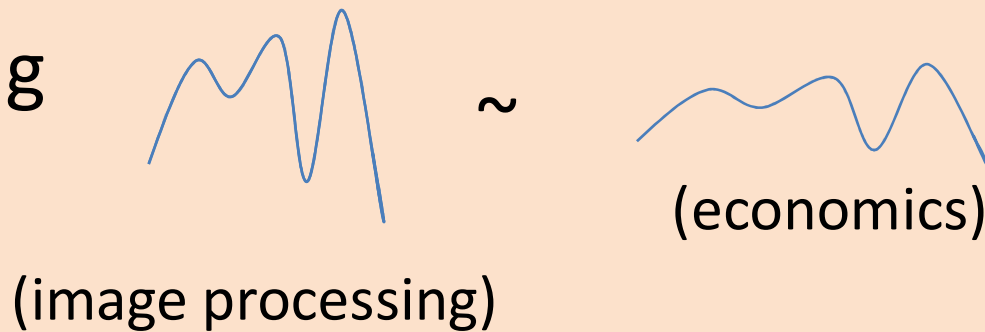
- Approximate



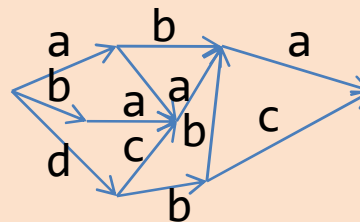
- Abelian  $ababbbaabbaab \sim aaaaaabbbbbbb$

- Order preserving

- 2D
  - ab cab
  - bc **abc**
  - ca **bca**



- Labelled graphs



(representation of pangenomes)

# RECENT EXAMPLES

(IN MY PUBLICATIONS)

## Space-Efficient Indexes for Uncertain Strings (ICDE 2024)

$X$	1	2	3	4	5	6
A	1	$1/2$	$3/4$	$4/5$	$1/2$	$1/4$
B	0	$1/2$	$1/4$	$1/5$	$1/2$	$3/4$

## Space-Efficient Indexes for Uncertain Strings (ICDE 2024)

$X$	1	2	3	4	5	6
A	1	1/2	3/4	4/5	1/2	1/4
B	0	1/2	1/4	1/5	1/2	3/4

$$P(X[3..5] = ABA) = 3/40$$



# Space-Efficient Indexes for Uncertain Strings (ICDE 2024)

$X$	1	2	3	4	5	6
A	1	1/2	3/4	4/5	1/2	1/4
B	0	1/2	1/4	1/5	1/2	3/4

$P(X[3..5] = ABA) = 3/40$

Occurrence if probability  $\geq 1/z$

$O(nz)$  size indexes are too large to be useful.

Our result:  $O(nz/l)$  size indexes for patterns of length  $\geq l$ .

# Space-Efficient Indexes for Uncertain Strings (ICDE 2024)

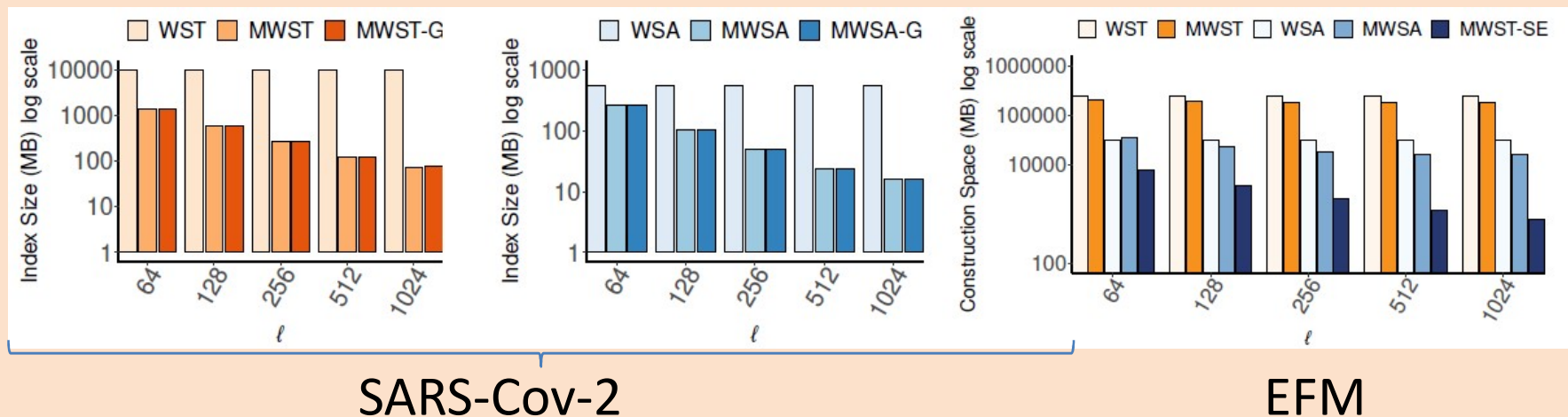
$X$	1	2	3	4	5	6
A	1	1/2	3/4	4/5	1/2	1/4
B	0	1/2	1/4	1/5	1/2	3/4

$$P(X[3..5] = ABA) = 3/40$$

Occurrence if probability  $\geq 1/z$

$O(nz)$  size indexes are too large to be useful.

Our result:  $O(nz/l)$  size indexes for patterns of length  $\geq l$ .



# Scalable Order-Preserving Pattern Mining (ICDM 2024)

Pattern mining in time series (IEEE Trans. Cybern. 2023)  
Their algorithm finds frequent patterns, but it takes  $\Omega(n^3)$  time.

# Scalable Order-Preserving Pattern Mining (ICDM 2024)

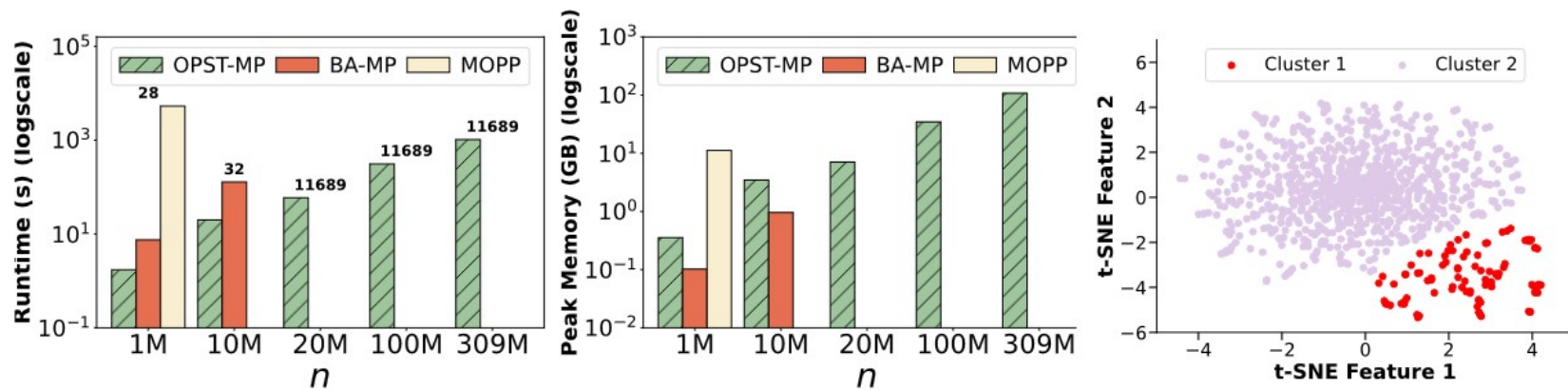
Pattern mining in time series (IEEE Trans. Cybern. 2023)  
Their algorithm finds frequent patterns, but it takes  $\Omega(n^3)$  time.

Our result: Structure that mines frequent patterns in  $O(n)$  time after  $O(n \log \sigma)$ -time construction, where  $\sigma$  is the alphabet size.

# Scalable Order-Preserving Pattern Mining (ICDM 2024)

Pattern mining in time series (IEEE Trans. Cybern. 2023)  
Their algorithm finds frequent patterns, but it takes  $\Omega(n^3)$  time.

Our result: Structure that mines frequent patterns in  $O(n)$  time after  $O(n \log \sigma)$ -time construction, where  $\sigma$  is the alphabet size.



Whale Signals

EEG of 2 groups of patients

Thank you for your attention!