

Algorithms in bioinformatics

Anton Bankevich
Center for Algorithmic Biotechnology
SPbSU

Bioinformatics vs Algorithmic Biology

Algorithmic biology:

development and application of data-analytical and theoretical methods, mathematical modeling and computational simulation techniques to the study of biological, behavioral, and social systems.

Bioinformatics:

field that develops methods and software tools for understanding biological data

Bioinformatics vs Algorithmic Biology

Algorithmic biology:



Bioinformatics:



Bioinformatics vs Algorithmic Biology

Algorithmic biology:



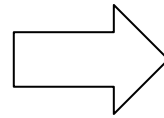
Bioinformatics:



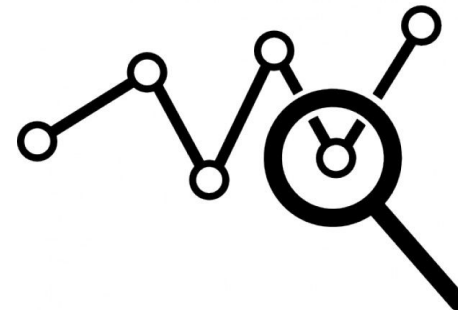
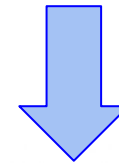
Discovery



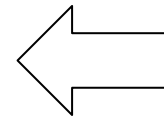
Technology



Data



Analysis



Discovery

Data

What is biological data?

- Sequencing reads
- Mass spectrometry
- Expression measurements

Questions about sequences

- Are they the same?
- Is one of them a part of another?
- What is common between sequences?

Questions about sequences

- Are they similar?
- Is one of them similar to a part of another?
- What is similar in the sequences?

Alignment

AACGCTAACGGTAA

AACCGCGAACTAA

Alignment

AACGCTAACGGTAA

AACCGCGAACTAA



AAC - GCTAACGGTAA

AACCGCGAAC - - TAA

Needleman–Wunsch algorithm

A C G T T A G
A C C T - A G

Score(x, x) = 1 (match)

Score(x, y) = -1 (mismatch)

Score(x, -) = Score(-, x) = -1 (Indel)

Score(ACGTTAG, ACCTAG) = 3

Needleman–Wunsch algorithm

	A	C	G	T	T	A	G
A							
C							
C							
T							
A							
G							

Needleman–Wunsch algorithm

	A	C	G	T	T	A	G	
	0	-1	-2	-3	-4	-5	-6	-7
A	-1							
C	-2							
C	-3							
T	-4							
A	-5							
G	-6							

Needleman–Wunsch algorithm

	A	C	G	T	T	A	G	
	0	-1	-2	-3	-4	-5	-6	-7
A	-1							
C	-2							
C	-3							
T	-4							
A	-5							
G	-6							

Needleman–Wunsch algorithm

	A	C	G	T	T	A	G	
	0	-1	-2	-3	-4	-5	-6	-7
A	-1							
C	-2							
C	-3							
T	-4							
A	-5							
G	-6							

Needleman–Wunsch algorithm

	A	C	G	T	T	A	G	
	0	-1	-2	-3	-4	-5	-6	-7
A	-1	1						
C	-2							
C	-3							
T	-4							
A	-5							
G	-6							

Needleman–Wunsch algorithm

	A	C	G	T	T	A	G	
	0	-1	-2	-3	-4	-5	-6	-7
A	-1	1						
C	-2							
C	-3							
T	-4							
A	-5							
G	-6							

Needleman–Wunsch algorithm

	A	C	G	T	T	A	G	
	0	-1	-2	-3	-4	-5	-6	-7
A	-1	1/-2						
C	-2							
C	-3							
T	-4							
A	-5							
G	-6							

Needleman–Wunsch algorithm

	A	C	G	T	T	A	G	
	0	-1	-2	-3	-4	-5	-6	-7
A	-1	1/-2						
C	-2							
C	-3							
T	-4							
A	-5							
G	-6							

Needleman–Wunsch algorithm

	A	C	G	T	T	A	G	
	0	-1	-2	-3	-4	-5	-6	-7
A	-1	1/-2/-2						
C	-2							
C	-3							
T	-4							
A	-5							
G	-6							

Needleman–Wunsch algorithm

	A	C	G	T	T	A	G	
	0	-1	-2	-3	-4	-5	-6	-7
A	-1	1						
C	-2							
C	-3							
T	-4							
A	-5							
G	-6							

Needleman–Wunsch algorithm

	A	C	G	T	T	A	G	
	0	-1	-2	-3	-4	-5	-6	-7
A	-1	1						
C	-2							
C	-3							
T	-4							
A	-5							
G	-6							

Needleman–Wunsch algorithm

	A	C	G	T	T	A	G	
	0	-1	-2	-3	-4	-5	-6	-7
A	-1	1	-2					
C	-2							
C	-3							
T	-4							
A	-5							
G	-6							

Needleman–Wunsch algorithm

	A	C	G	T	T	A	G	
	0	-1	-2	-3	-4	-5	-6	-7
A	-1	1	-2/0					
C	-2							
C	-3							
T	-4							
A	-5							
G	-6							

Needleman–Wunsch algorithm

	A	C	G	T	T	A	G	
	0	-1	-2	-3	-4	-5	-6	-7
A	-1	1	-2/0/-3					
C	-2							
C	-3							
T	-4							
A	-5							
G	-6							

Needleman–Wunsch algorithm

	A	C	G	T	T	A	G	
	0	-1	-2	-3	-4	-5	-6	-7
A	-1	1	0					
C	-2							
C	-3							
T	-4							
A	-5							
G	-6							

Needleman–Wunsch algorithm

	A	C	G	T	T	A	G	
	0	-1	-2	-3	-4	-5	-6	-7
A	-1	1	0	-1	-2	-3	-4	-5
C	-2	0						
C	-3							
T	-4							
A	-5							
G	-6							

The image shows a Needleman–Wunsch algorithm matrix. The columns are labeled A, C, G, T, T, A, G. The rows are labeled A, C, C, T, A, G. The matrix contains numerical values representing the alignment score. Red arrows indicate the path of the optimal alignment, starting from the bottom-left cell (row G, column A) and moving up and right to the top-right cell (row A, column G). The cells (row A, column T) and (row A, column A) are highlighted in green, and the cell (row A, column A) is highlighted in orange.

Needleman–Wunsch algorithm

	A	C	G	T	T	A	G	
	0	-1	-2	-3	-4	-5	-6	-7
A	-1	1	0	-1	-2	-3	-4	-5
C	-2	0						
C	-3							
T	-4							
A	-5							
G	-6							

Needleman–Wunsch algorithm

	A	C	G	T	T	A	G	
	0	-1	-2	-3	-4	-5	-6	-7
A	-1	1	0	-1	-2	-3	-4	-5
C	-2	0	2					
C	-3							
T	-4							
A	-5							
G	-6							

Needleman–Wunsch algorithm

	A	C	G	T	T	A	G	
	0	-1	-2	-3	-4	-5	-6	-7
A	-1	1	0	-1	-2	-3	-4	-5
C	-2	0	2/-1					
C	-3							
T	-4							
A	-5							
G	-6							

Needleman–Wunsch algorithm

	A	C	G	T	T	A	G	
	0	-1	-2	-3	-4	-5	-6	-7
A	-1	1	0	-1	-2	-3	-4	-5
C	-2	0	2/-1/-1					
C	-3							
T	-4							
A	-5							
G	-6							

The image shows a dynamic programming table for the Needleman–Wunsch algorithm. The table is a grid with rows and columns labeled with nucleotide bases (A, C, G, T). The top row and left column contain the sequence indices (0 to -7). The cell at row A, column C (value 0) is highlighted in green. The cell at row C, column C (value 0) is highlighted in orange and contains the text "2/-1/-1". Red arrows indicate the backpointers for the sequence alignment, showing the path from the top-right cell (0, -7) to the bottom-left cell (-7, 0).

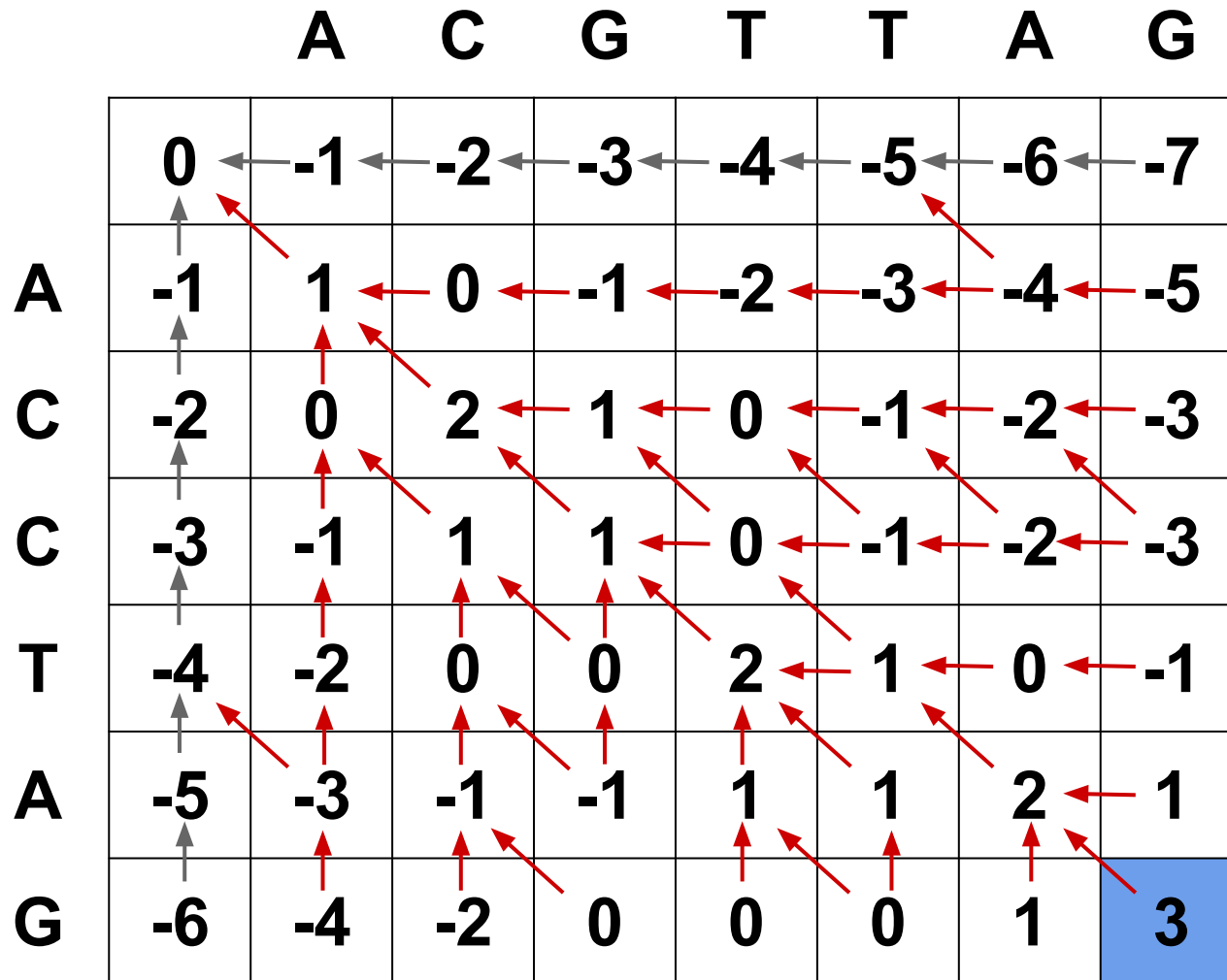
Needleman–Wunsch algorithm

	A	C	G	T	T	A	G	
	0	-1	-2	-3	-4	-5	-6	-7
A	-1	1	0	-1	-2	-3	-4	-5
C	-2	0	2					
C	-3							
T	-4							
A	-5							
G	-6							

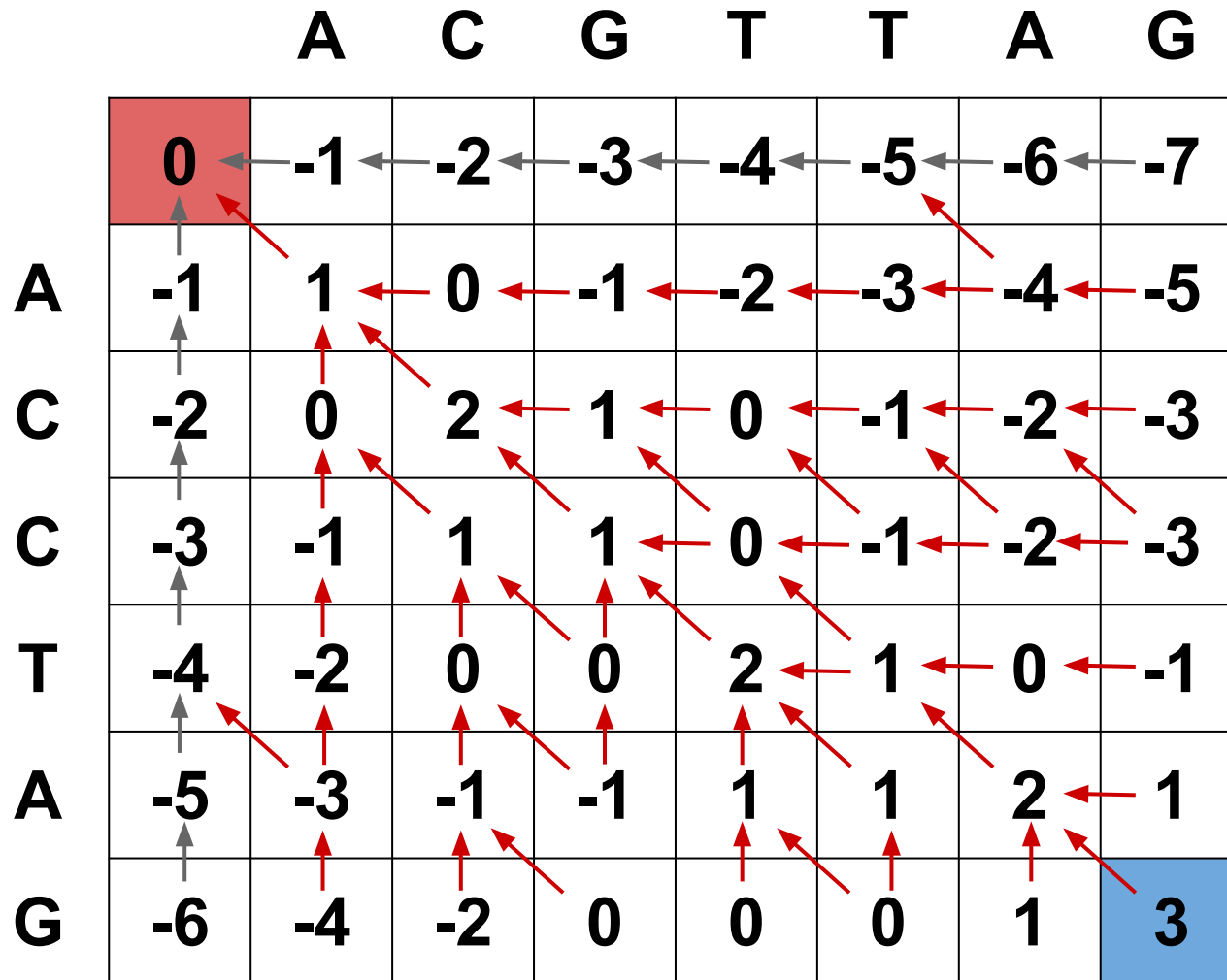
Needleman–Wunsch algorithm

$$M(i, j) = \text{MAX} \begin{cases} M(i - 1, j - 1) + \text{Score}(X(i), Y(j)) \\ M(i - 1, j) + \text{Score}(X(i), '-') \\ M(i, j - 1) + \text{Score}('-', Y(j)) \end{cases}$$

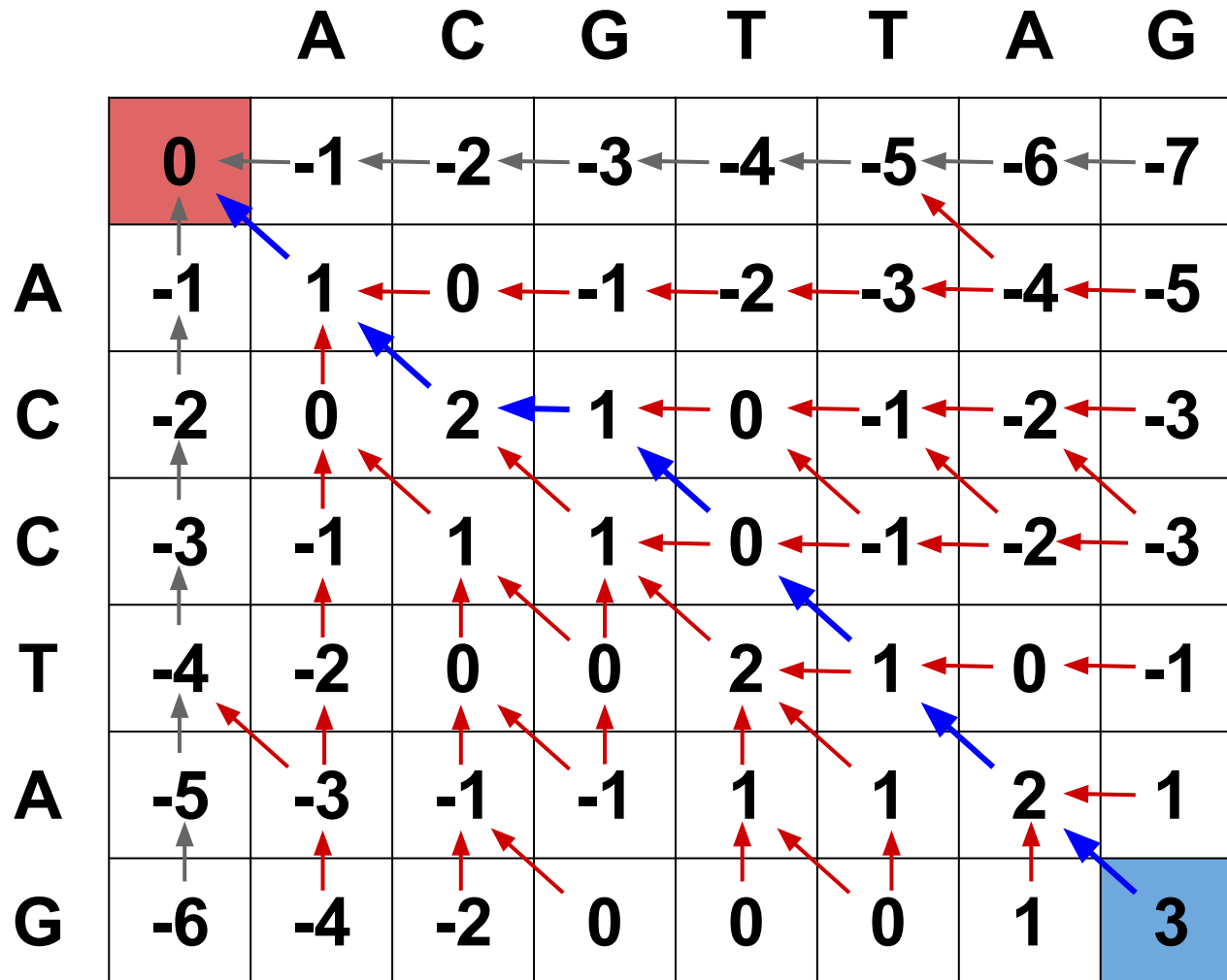
Needleman–Wunsch algorithm



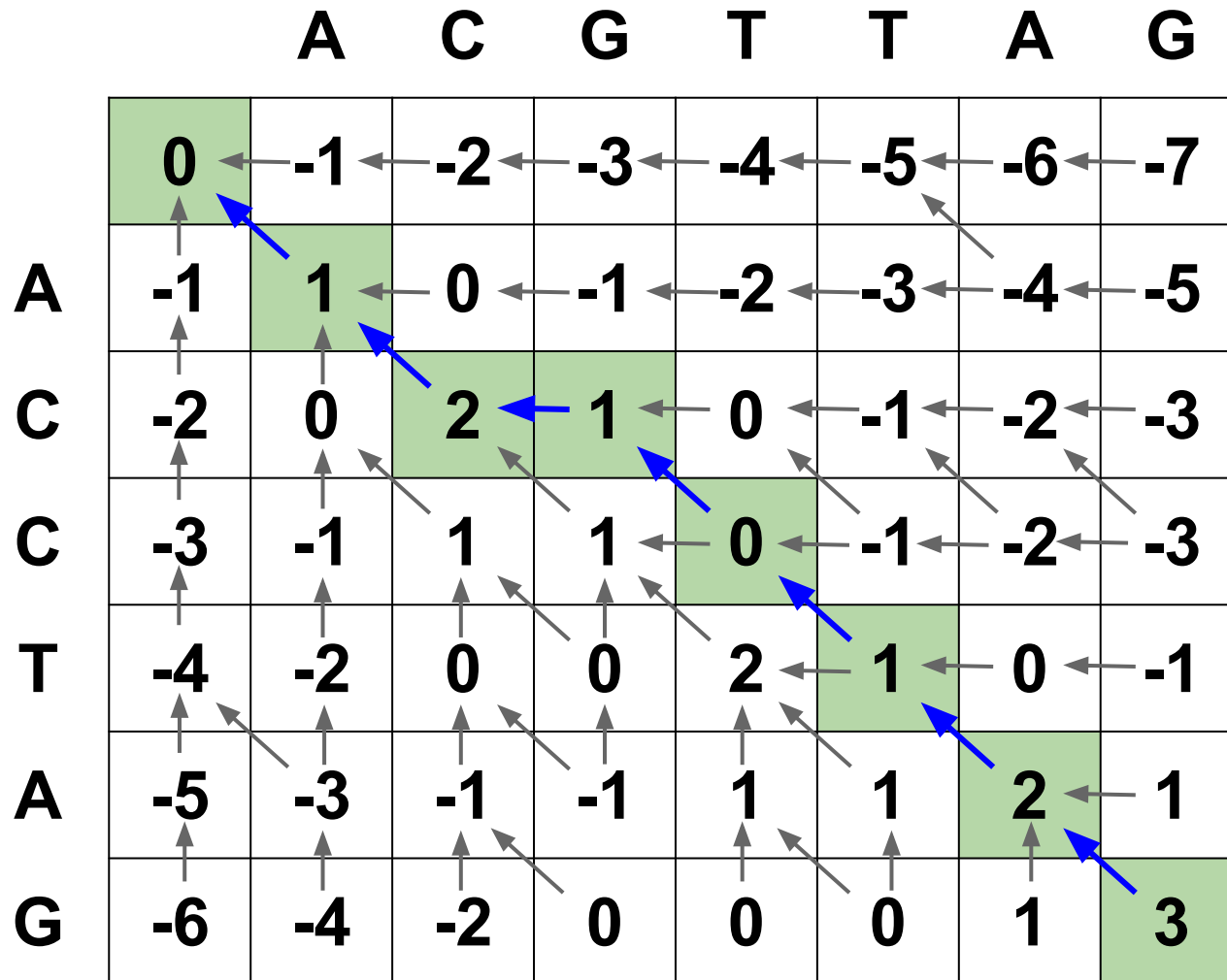
Needleman–Wunsch algorithm



Needleman–Wunsch algorithm



Needleman–Wunsch algorithm



Needleman–Wunsch algorithm

G	G

	A	C	G	T	T	A	G	
	0	-1	-2	-3	-4	-5	-6	-7
A	-1	1	0	-1	-2	-3	-4	-5
C	-2	0	2	1	0	-1	-2	-3
C	-3	-1	1	1	0	-1	-2	-3
T	-4	-2	0	0	2	1	0	-1
A	-5	-3	-1	-1	1	1	2	1
G	-6	-4	-2	0	0	0	1	3

Needleman–Wunsch algorithm

A	A
G	G

	A	C	G	T	T	A	G	
	0	-1	-2	-3	-4	-5	-6	-7
A	-1	1	0	-1	-2	-3	-4	-5
C	-2	0	2	1	0	-1	-2	-3
C	-3	-1	1	1	0	-1	-2	-3
T	-4	-2	0	0	2	1	0	-1
A	-5	-3	-1	-1	1	1	2	1
G	-6	-4	-2	0	0	0	1	3

Needleman–Wunsch algorithm

T	T
A	A
G	G

	A	C	G	T	T	A	G	
	0	-1	-2	-3	-4	-5	-6	-7
A	-1	1	0	-1	-2	-3	-4	-5
C	-2	0	2	1	0	-1	-2	-3
C	-3	-1	1	1	0	-1	-2	-3
T	-4	-2	0	0	2	1	0	-1
A	-5	-3	-1	-1	1	1	2	1
G	-6	-4	-2	0	0	0	1	3

Needleman–Wunsch algorithm

C	T
T	T
A	A
G	G

	A	C	G	T	T	A	G	
	0	-1	-2	-3	-4	-5	-6	-7
A	-1	1	0	-1	-2	-3	-4	-5
C	-2	0	2	1	0	-1	-2	-3
C	-3	-1	1	1	0	-1	-2	-3
T	-4	-2	0	0	2	1	0	-1
A	-5	-3	-1	-1	1	1	2	1
G	-6	-4	-2	0	0	0	1	3

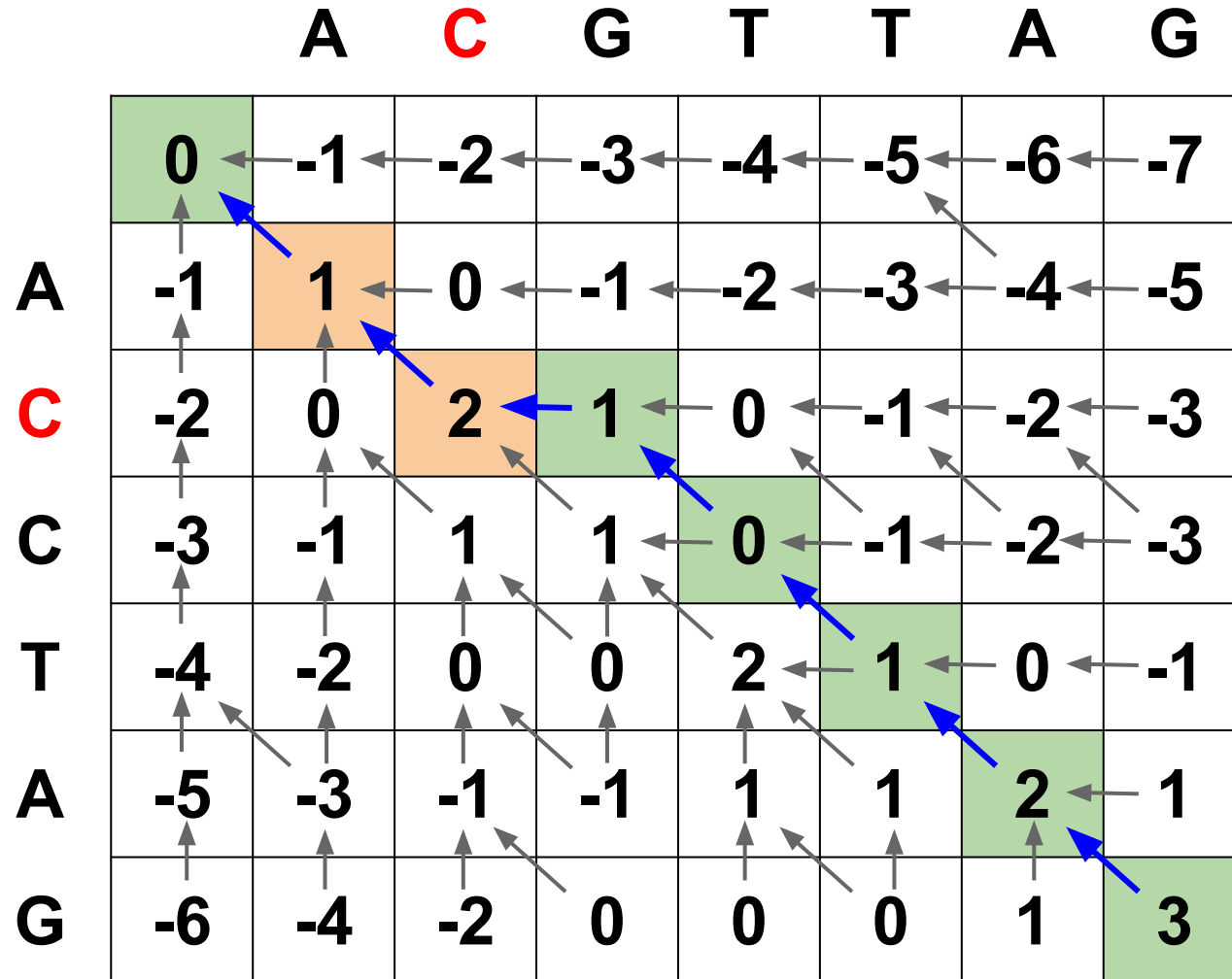
Needleman–Wunsch algorithm

-	G
C	T
T	T
A	A
G	G

	A	C	G	T	T	A	G	
	0	-1	-2	-3	-4	-5	-6	-7
A	-1	1	0	-1	-2	-3	-4	-5
C	-2	0	2	1	0	-1	-2	-3
C	-3	-1	1	1	0	-1	-2	-3
T	-4	-2	0	0	2	1	0	-1
A	-5	-3	-1	-1	1	1	2	1
G	-6	-4	-2	0	0	0	1	3

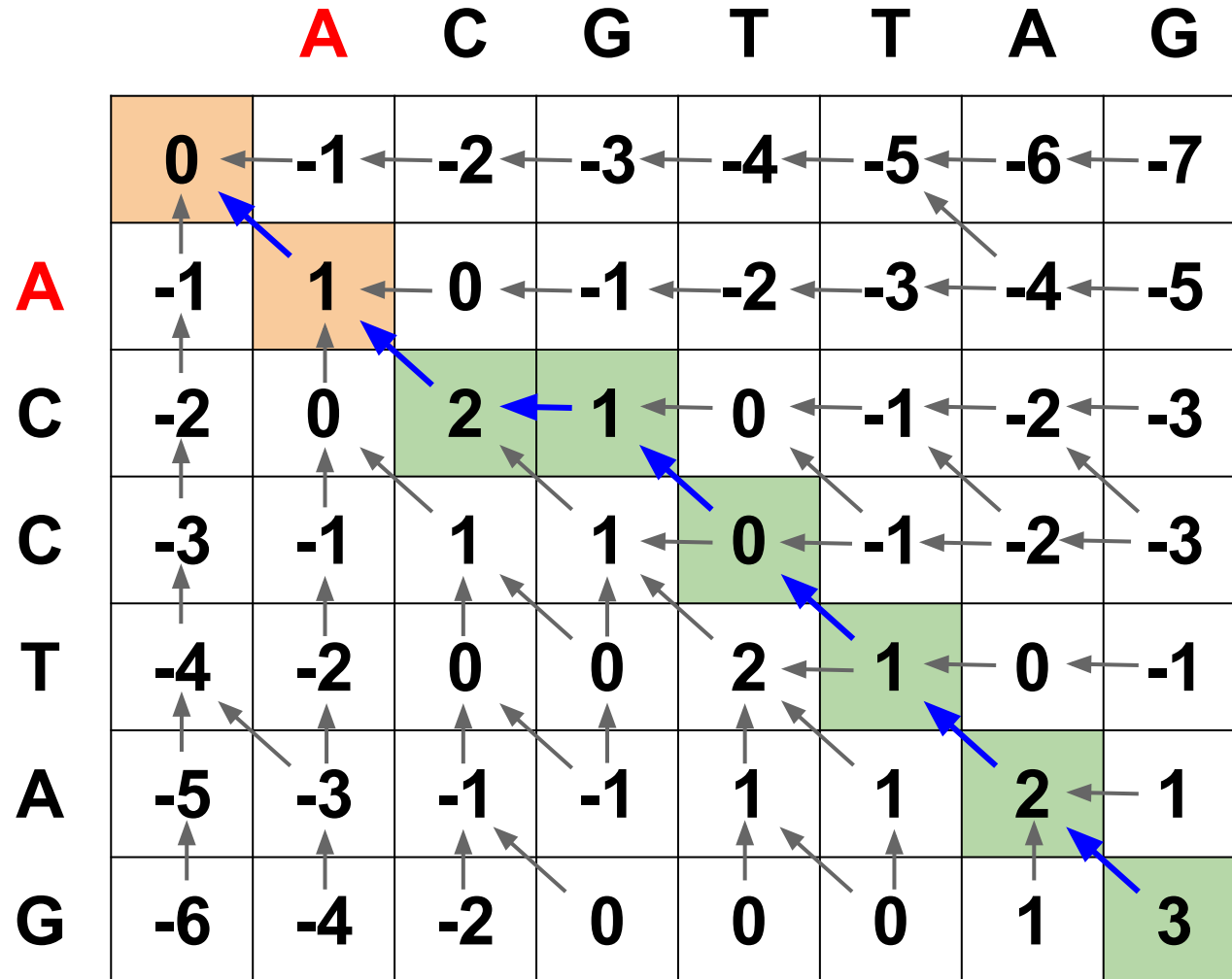
Needleman–Wunsch algorithm

C	C
-	G
C	T
T	T
A	A
G	G



Needleman–Wunsch algorithm

A	A
C	C
-	G
C	T
T	T
A	A
G	G

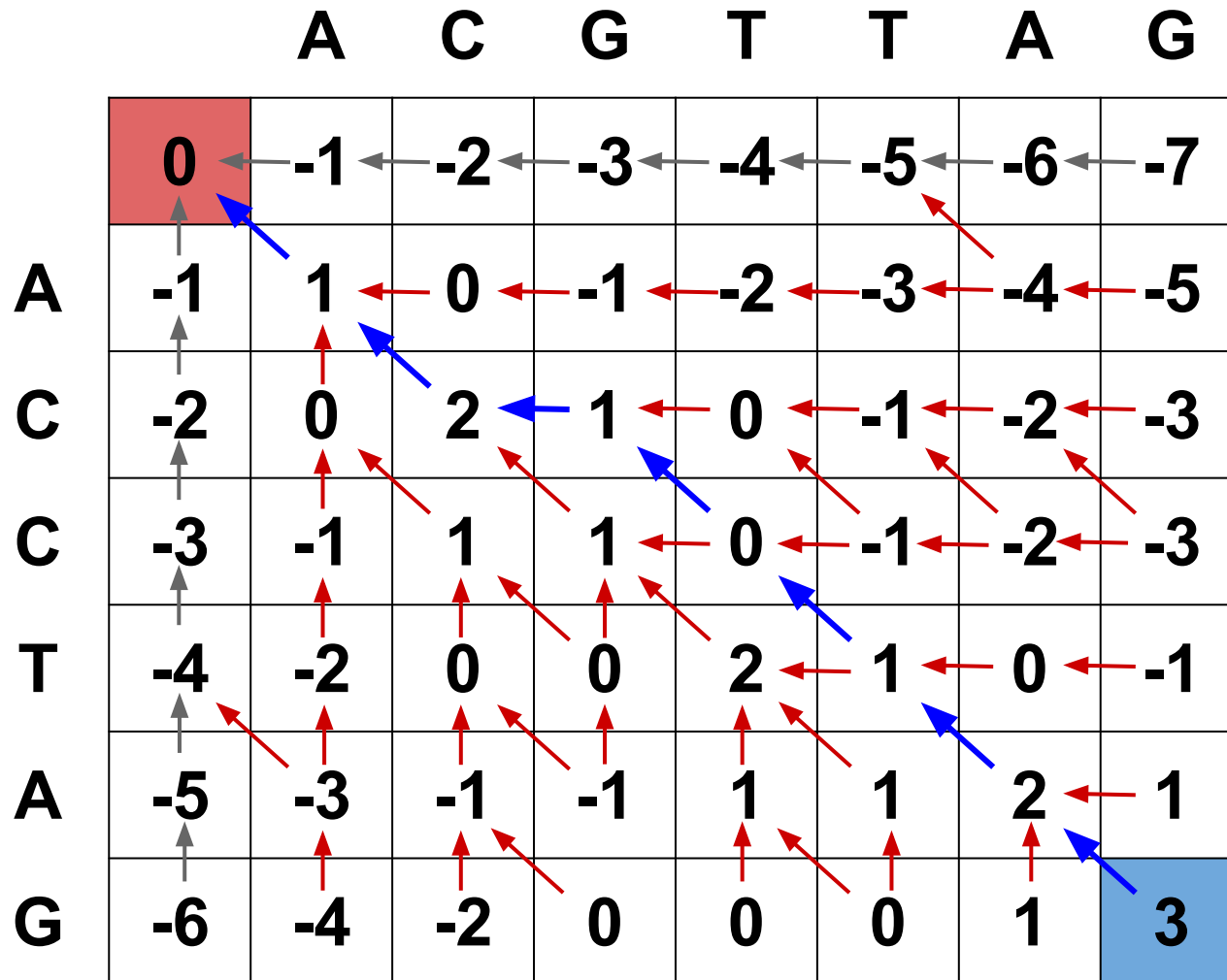


Needleman–Wunsch algorithm

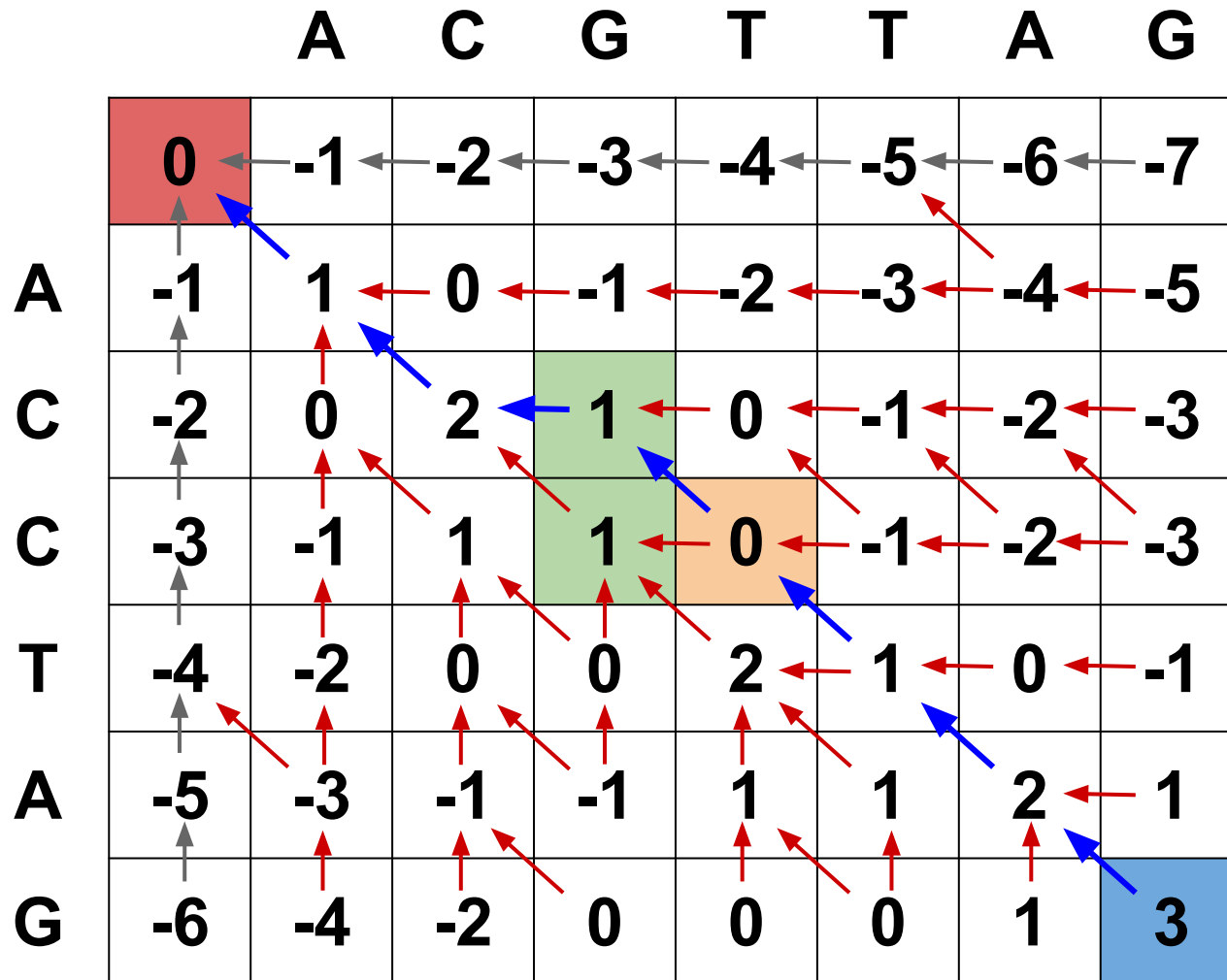
A	C	G	T	T	A	G
A	C	C	T	-	A	G
A	C	G	T	T	A	G
A	C	-	C	T	A	G

Score(ACGTTAG, ACCTAG) = 3

Needleman–Wunsch algorithm



Needleman–Wunsch algorithm



Needleman–Wunsch algorithm

	A	C	G	T	T	A	G	
	0	-1	-2	-3	-4	-5	-6	-7
A	-1	1	0	-1	-2	-3	-4	-5
C	-2	0	2	1	0	-1	-2	-3
C	-3	-1	1	1	0	-1	-2	-3
T	-4	-2	0	0	2	1	0	-1
A	-5	-3	-1	-1	1	1	2	1
G	-6	-4	-2	0	0	0	1	3

Needleman–Wunsch algorithm

	A	C	G	T	T	A	G	
	0	-1	-2	-3	-4	-5	-6	-7
A	-1	1	0	-1	-2	-3	-4	-5
C	-2	0	2	1	0	-1	-2	-3
C	-3	-1	1	1	0	-1	-2	-3
T	-4	-2	0	0	2	1	0	-1
A	-5	-3	-1	-1	1	1	2	1
G	-6	-4	-2	0	0	0	1	3

Needleman–Wunsch algorithm

Time =

Memory =

Needleman–Wunsch algorithm

$$\text{Time} = O(n^2)$$

$$\text{Memory} = O(n^2)$$

How about the remaining questions?

- Are they similar?
- Is one of them similar to a part of another?
- What is similar in the sequences?

Needleman–Wunsch algorithm

	A	C	G	C	G	A	G
	0	0	0	0	0	0	0
G	-1	-1	1	0	1	0	1
C	-2	-2	0	2	1	0	-1
G	-3	-3	-1	-1	1	3	2

- Is one of them similar to a part of another?

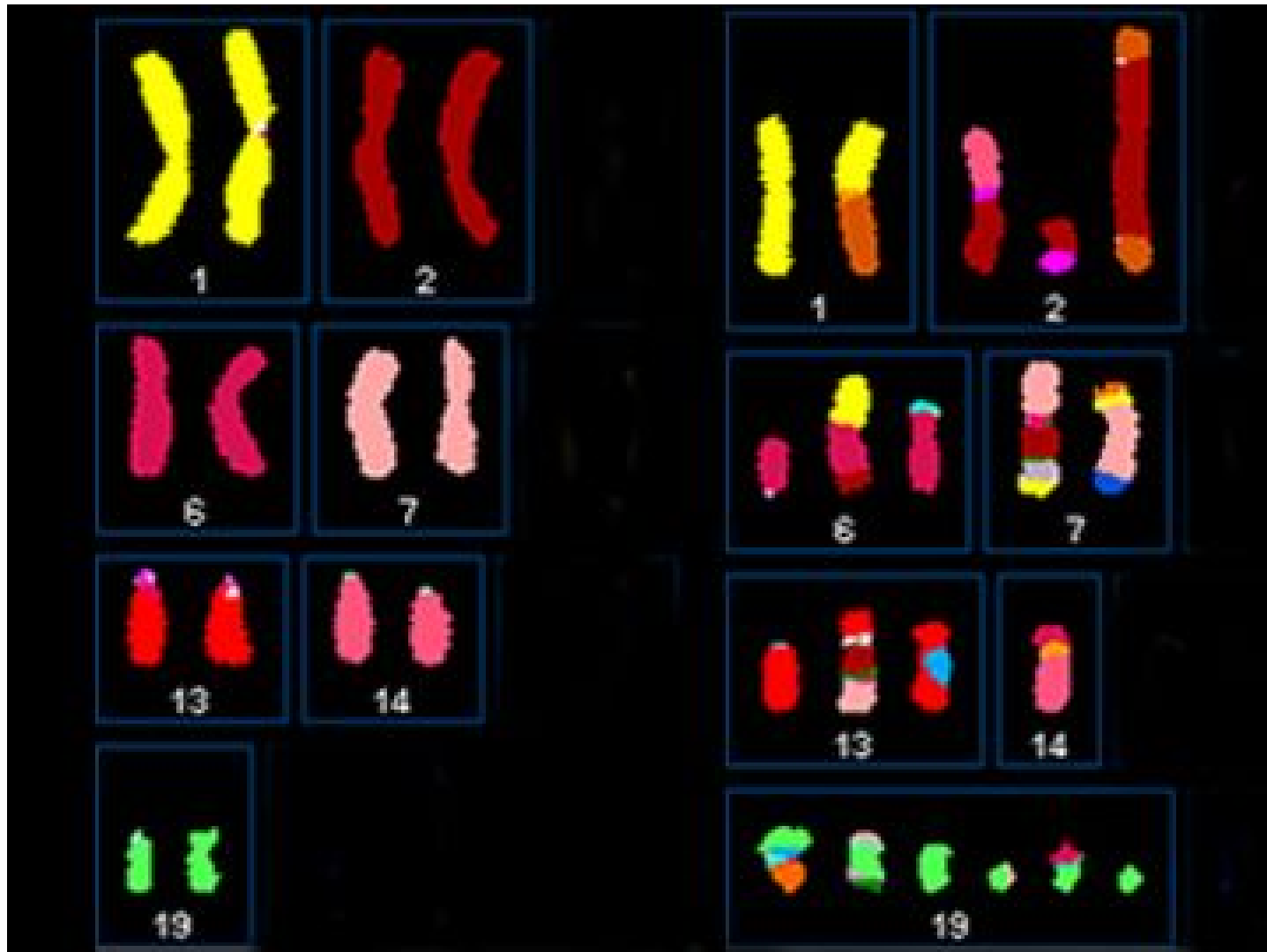
How about the remaining questions?

- Are they similar?
- Is one of them similar to a part of another?
- What is similar in the sequences?

Why is this not enough?

- We deal with genomes of extreme size
- We have extreme number of sequencing reads
- We need more careful alignment that can handle structural rearrangements

Rearrangements



Burrows-Wheeler transform

a c a a c g

Burrows-Wheeler transform

a c a a c g \$

Burrows-Wheeler transform

a c a a c g \$
\$ a c a a c g

Burrows-Wheeler transform

a c a a c g \$

\$ a c a a c g

g \$ a c a a c

Burrows-Wheeler transform

a c a a c g \$
\$ a c a a c g
g \$ a c a a c
c g \$ a c a a
a c g \$ a c a
a a c g \$ a c
c a a c g \$ a

Burrows-Wheeler transform

\$ a c a a c g
a a c g \$ a c
a c a a c g \$
a c g \$ a c a
c a a c g \$ a
c g \$ a c a a
g \$ a c a a c

Burrows-Wheeler transform

\$ a c a a c g
a a c g \$ a c
a c a a c g \$
a c g \$ a c a
c a a c g \$ a
c g \$ a c a a
g \$ a c a a c

Burrows-Wheeler transform

a c a a c g \$



g c \$ a a a c

Burrows-Wheeler transform

a c a a c g \$
↑ ?
g c \$ a a a c

Burrows-Wheeler transform

\$ a c a a c g
a a c g \$ a c
a c a a c g \$
a c g \$ a c a
c a a c g \$ a
c g \$ a c a a
g \$ a c a a c

Burrows-Wheeler transform

\$ a c a a c g

a a c g \$ a c

a c a a c g \$

a c g \$ a c a

c a a c g \$ a

c g \$ a c a a

g \$ a c a a c

Burrows-Wheeler transform

g \$ a c a a c
c a a c g \$ a
\$ a c a a c g
a a c g \$ a c
a c a a c g \$
a c g \$ a c a
c g \$ a c a a

Burrows-Wheeler transform

\$ a c a a c g
a a c g \$ a c
a c a a c g \$
a c g \$ a c a
c a a c g \$ a
c g \$ a c a a
g \$ a c a a c

Burrows-Wheeler transform

g	\$ a c a a c g
c	a a c g \$ a c
\$	a c a a c g \$
a	a c g \$ a c a
a	c a a c g \$ a
a	c g \$ a c a a
c	g \$ a c a a c

Burrows-Wheeler transform

\$	\$ a c a a c g
a	a a c g \$ a c
a	a c a a c g \$
a	a c g \$ a c a
c	c a a c g \$ a
c	c g \$ a c a a
g	g \$ a c a a c

Burrows-Wheeler transform

\$
a
a
a
c
c
g

\$	a	c	a	a	c	g
a	a	c	g	\$	a	c
a	c	a	a	c	g	\$
a	c	g	\$	a	c	a
c	a	a	c	g	\$	a
c	g	\$	a	c	a	a
g	\$	a	c	a	a	c

Burrows-Wheeler transform

g \$

c a

\$ a

a a

a c

a c

c g

\$ a c a a c **g**

a a c g \$ a **c**

a c a a c g **\$**

a c g \$ a c **a**

c a a c g \$ **a**

c g \$ a c a **a**

g \$ a c a a **c**

Burrows-Wheeler transform

\$ a

a a

a c

a c

c a

c g

g \$

\$ a c a a c g

a a c g \$ a c

a c a a c g \$

a c g \$ a c a

c a a c g \$ a

c g \$ a c a a

g \$ a c a a c

Burrows-Wheeler transform

\$ a

a a

a c

a c

c a

c g

g \$

\$ a c a a c g

a a c g \$ a c

a c a a c g \$

a c g \$ a c a

c a a c g \$ a

c g \$ a c a a

g \$ a c a a c

Burrows-Wheeler transform

g \$ a

c a a

\$ a c

a a c

a c a

a c g

c g \$

\$ a c a a c **g**

a a c g \$ a **c**

a c a a c g **\$**

a c g \$ a c **a**

c a a c g \$ **a**

c g \$ a c a **a**

g \$ a c a a **c**

Burrows-Wheeler transform

\$ a c

a a c

a c a

a c g

c a a

c g \$

g \$ a

\$ a c a a c **g**

a a c g \$ a **c**

a c a a c g **\$**

a c g \$ a c **a**

c a a c g \$ **a**

c g \$ a c a **a**

g \$ a c a a **c**

Burrows-Wheeler transform

\$ a c

a a c

a c a

a c g

c a a

c g \$

g \$ a

\$ a c a a c g

a a c g \$ a c

a c a a c g \$

a c g \$ a c a

c a a c g \$ a

c g \$ a c a a

g \$ a c a a c

Burrows-Wheeler transform

g \$ a c

c a a c

\$ a c a

a a c g

a c a a

a c g \$

c g \$ a

\$ a c a a c **g**

a a c g \$ a **c**

a c a a c g **\$**

a c g \$ a c **a**

c a a c g \$ **a**

c g \$ a c a **a**

g \$ a c a a **c**

Burrows-Wheeler transform

\$ a c a
a a c g
a c a a
a c g \$
c a a c
c g \$ a
g \$ a c

\$ a c a a c g
a a c g \$ a c
a c a a c g \$
a c g \$ a c a
c a a c g \$ a
c g \$ a c a a
g \$ a c a a c

Burrows-Wheeler transform

g \$ a c a

c a a c g

\$ a c a a

a a c g \$

a c a a c

a c g \$ a

c g \$ a c

\$ a c a a c **g**

a a c g \$ a **c**

a c a a c g **\$**

a c g \$ a c **a**

c a a c g \$ **a**

c g \$ a c a **a**

g \$ a c a a **c**

Burrows-Wheeler transform

\$ a c a a
a a c g \$
a c a a c
a c g \$ a
c a a c g
c g \$ a c
g \$ a c a

\$ a c a a c g
a a c g \$ a c
a c a a c g \$
a c g \$ a c a
c a a c g \$ a
c g \$ a c a a
g \$ a c a a c

Burrows-Wheeler transform

g \$ a c a a

c a a c g \$

\$ a c a a c

a a c g \$ a

a c a a c g

a c g \$ a c

c g \$ a c a

\$ a c a a c **g**

a a c g \$ a **c**

a c a a c g **\$**

a c g \$ a c **a**

c a a c g \$ **a**

c g \$ a c a **a**

g \$ a c a a **c**

Burrows-Wheeler transform

\$ a c a a c
a a c g \$ a
a c a a c g
a c g \$ a c
c a a c g \$
c g \$ a c a
g \$ a c a a

\$ a c a a c g
a a c g \$ a c
a c a a c g \$
a c g \$ a c a
c a a c g \$ a
c g \$ a c a a
g \$ a c a a c

Burrows-Wheeler transform

g \$ a c a a c
c a a c g \$ a
\$ a c a a c g
a a c g \$ a c
a c a a c g \$
a c g \$ a c a
c g \$ a c a a

\$ a c a a c **g**
a a c g \$ a **c**
a c a a c g **\$**
a c g \$ a c **a**
c a a c g \$ **a**
c g \$ a c a **a**
g \$ a c a a **c**

Burrows-Wheeler transform

\$ a c a a c g

a a c g \$ a c

a c a a c g \$

a c g \$ a c a

c a a c g \$ a

c g \$ a c a a

g \$ a c a a c

\$ a c a a c g

a a c g \$ a c

a c a a c g \$

a c g \$ a c a

c a a c g \$ a

c g \$ a c a a

g \$ a c a a c

First-last property

$\$$ ₁	a	c	a	a	c	g
a ₁	a	c	g	$\$$	a	c
a ₂	c	a	a	c	g	$\$$
a ₃	c	g	$\$$	a	c	a
c ₁	a	a	c	g	$\$$	a
c ₂	g	$\$$	a	c	a	a
g ₁	$\$$	a	c	a	a	c

First-last property

$\$$ ₁ a c a a c **g**₁
a₁ a c g \$ a **c**
a₂ c a a c g **\$**₁
a₃ c g \$ a c **a**
c₁ a a c g \$ **a**
c₂ g \$ a c a **a**
g₁ \$ a c a a **c**

First-last property

$\$$ ₁ a c a a c **g**₁

a₁ a c g \$ a **c**

a₂ c a a c g **\$**₁

a₃ c g \$ a c **a**

c₁ a a c g \$ **a**

c₂ g \$ a c a **a**

g₁ \$ a c a a **c**

First-last property

a_3 c g \$ a c **a**
 c_1 a a c g \$ **a**
 c_2 g \$ a c a **a**

$\$$ ₁ a c a a c **g**₁
 a ₁ a c g \$ a **c**
 a ₂ c a a c g **\$**₁
a₃ c g \$ a c **a**
c₁ a a c g \$ **a**
c₂ g \$ a c a **a**
g₁ \$ a c a a **c**

First-last property

a a₃ c g \$ a c
a c₁ a a c g \$
a c₂ g \$ a c a

\$₁ a c a a c **g**₁
a₁ a c g \$ a **c**
a₂ c a a c g \$₁
a₃ c g \$ a c **a**
c₁ a a c g \$ **a**
c₂ g \$ a c a **a**
g₁ \$ a c a a **c**

First-last property

a a₃ c g \$ a c
a c₁ a a c g \$
a c₂ g \$ a c a

\$₁ a c a a c **g**₁
a₁ a c g \$ a **c**
a₂ c a a c g \$₁
a₃ c g \$ a c **a**
c₁ a a c g \$ **a**
c₂ g \$ a c a **a**
g₁ \$ a c a a **c**

First-last property

a₁ **a₃** c g \$ a c
a₂ **c₁** a a c g \$
a₃ **c₂** g \$ a c a

\$₁ a c a a c **g₁**
a₁ a c g \$ a **c**
a₂ c a a c g **\$₁**
a₃ c g \$ a c **a**
c₁ a a c g \$ **a**
c₂ g \$ a c a **a**
g₁ \$ a c a a **c**

First-last property

a_3 c g \$ a c a_1
 c_1 a a c g \$ a_2
 c_2 g \$ a c a a_3

$\$1$ a c a a c g_1
 a_1 a c g \$ a c
 a_2 c a a c g $\$1$
 a_3 c g \$ a c a
 c_1 a a c g \$ a
 c_2 g \$ a c a a
 g_1 \$ a c a a c

First-last property

a_3 c g \$ a c a_1
 c_1 a a c g \$ a_2
 c_2 g \$ a c a a_3

$\$1$ a c a a c g_1
 a_1 a c g \$ a c
 a_2 c a a c g $\$1$
 a_3 c g \$ a c a
 c_1 a a c g \$ a
 c_2 g \$ a c a a
 g_1 \$ a c a a c

First-last property

a_3 c g \$ a c a_1
 c_1 a a c g \$ a_2
 c_2 g \$ a c a a_3

$\$1$ a c a a c g_1
 a_1 a c g \$ a c
 a_2 c a a c g $\$1$
 a_3 c g \$ a c a_1
 c_1 a a c g \$ a_2
 c_2 g \$ a c a a_3
 g_1 \$ a c a a c

First-last property

$\$$ ₁	a	c	a	a	c	g	g ₁
a ₁	a	c	g	$\$$	a	c ₁	
a ₂	c	a	a	c	g	$\$$ ₁	
a ₃	c	g	$\$$	a	c	a ₁	
c ₁	a	a	c	g	$\$$	a ₂	
c ₂	g	$\$$	a	c	a	a ₃	
g ₁	$\$$	a	c	a	a	c ₂	

Pattern search

$\$$ ₁	a	c	a	a	c	g	₁		
a	₁	a	c	g	$\$$	a	c	₁	
a	₂	c	a	a	c	g	$\$$	₁	
a	₃	c	g	$\$$	a	c	a	₁	
c	₁	a	a	c	g	$\$$	a	₂	
c	₂	g	$\$$	a	c	a	a	₃	
g	₁	$\$$	a	c	a	a	a	c	₂

a a c

Pattern search

$\$$ ₁	a	c	a	a	c	g	₁	a	a	c
a	a	c	g	\$	a	c	₁			
a	c	a	a	c	g	\$	₁			
a	c	g	\$	a	c	a	₁			
c	a	a	c	g	\$	a	₂			
c	g	\$	a	c	a	a	₃			
g	\$	a	c	a	a	c	₂			

a a c

Pattern search

$\$$ ₁ a c a a c **g**₁

a a c

a₁ a c g \$ a **c**₁

a₂ c a a c g **\$**₁

a₃ c g \$ a c **a**₁

c₁ a a c g \$ **a**₂

c₂ g \$ a c a **a**₃

g₁ \$ a c a a **c**₂

Pattern search

$\$$ ₁ a c a a c **g**₁

a a c

a₁ a c g \$ a **c**₁

a₂ c a a c g **\$**₁

a₃ c g \$ a c **a**₁

c₁ a a c g \$ **a**₂

c₂ g \$ a c a **a**₃

g₁ \$ a c a a **c**₂

Pattern search

$\$$ ₁	a	c	a	a	c	g	g ₁
a ₁	a	c	g	$\$$	a	c ₁	
a ₂	c	a	a	c	g	$\$$ ₁	
a ₃	c	g	$\$$	a	c	a ₁	
c ₁	a	a	c	g	$\$$	a ₂	
c ₂	g	$\$$	a	c	a	a ₃	
g ₁	$\$$	a	c	a	a	c ₂	

a a c

Pattern search

$\$$ ₁	a	c	a	a	c	g	g ₁
a	a	c	g	$\$$	a	c	c ₁
a	c	a	a	c	g	$\$$	$\$$ ₁
a	c	g	$\$$	a	c	a	a ₁
c	a	a	c	g	$\$$	a	a ₂
c	g	$\$$	a	c	a	a	a ₃
g	$\$$	a	c	a	a	a	c ₂

a a c

Pattern search

$\$$ ₁	a	c	a	a	c	g	g ₁
a	a	c	g	$\$$	a	c	c ₁
a	c	a	a	c	g	$\$$	$\$$ ₁
a	c	g	$\$$	a	c	a	a ₁
c	a	a	c	g	$\$$	a	a ₂
c	g	$\$$	a	c	a	a	a ₃
g	$\$$	a	c	a	a	c	c ₂

a a c

Pattern search

$\$$ ₁	a	c	a	a	c	g	g ₁
a ₁	a	c	g	$\$$	a	c ₁	
a ₂	c	a	a	c	g	$\$$ ₁	
a ₃	c	g	$\$$	a	c	a ₁	
c ₁	a	a	c	g	$\$$	a ₂	
c ₂	g	$\$$	a	c	a	a ₃	
g ₁	$\$$	a	c	a	a	c ₂	

a a c

Pattern search

$\$$ ₁	a	c	a	a	c	g	g ₁
a ₁	a	c	g	$\$$	a	c ₁	
a ₂	c	a	a	c	g	$\$$ ₁	
a ₃	c	g	$\$$	a	c	a ₁	
c ₁	a	a	c	g	$\$$	a ₂	
c ₂	g	$\$$	a	c	a	a ₃	
g ₁	$\$$	a	c	a	a	c ₂	

a a c



Pattern search

For how long does it work?

Pattern search

In first column find index of a letter with given rank

$\$$ ₁ a c a a c **g**₁

a₁ a c g \$ a **c**₁

a₂ c a a c g **\$**₁

a₃ c g \$ a c **a**₁

c₁ a a c g \$ **a**₂

c₂ g \$ a c a **a**₃

g₁ \$ a c a a **c**₂

Pattern search

In first column find index of a letter with given rank

→	\$₁	a	c	a	a	c	g₁
→	a₁	a	c	g	\$	a	c₁
	a₂	c	a	a	c	g	\$₁
	a₃	c	g	\$	a	c	a₁
→	c₁	a	a	c	g	\$	a₂
	c₂	g	\$	a	c	a	a₃
→	g₁	\$	a	c	a	a	c₂

Pattern search

In first column find index of a letter with given rank — $O(1)$

$\$$ ₁ a c a a c **g**₁

a₁ a c g \$ a **c**₁

a₂ c a a c g **\$**₁

a₃ c g \$ a c **a**₁

c₁ a a c g \$ **a**₂

c₂ g \$ a c a **a**₃

g₁ \$ a c a a **c**₂

Pattern search

In BWT find all given letters in specified range

$\$$ ₁ a c a a c **g**₁

a₁ a c g \$ a **c**₁

a₂ c a a c g **\$**₁

a₃ c g \$ a c **a**₁

c₁ a a c g \$ **a**₂

c₂ g \$ a c a **a**₃

g₁ \$ a c a a **c**₂

Pattern search

In BWT find a number of given letter above a certain row

$\$$ ₁ a c a a c **g**₁

a₁ a c g \$ a **c**₁

a₂ c a a c g **\$**₁

a₃ c g \$ a c **a**₁

c₁ a a c g \$ **a**₂

c₂ g \$ a c a **a**₃

g₁ \$ a c a a **c**₂

Pattern search

In BWT find a number of given letter above a certain row

$\$$ ₁	a	c	a	a	c	g	₁	0
a	a	c	g	\$	a	c	₁	0
a	c	a	a	c	g	\$	₁	0
a	c	g	\$	a	c	a	₁	1
c	a	a	c	g	\$	a	₂	2
c	g	\$	a	c	a	a	₃	3
g	\$	a	c	a	a	c	₂	3

Pattern search

In BWT find a number of given letter above a certain row

$\$$ ₁	a	c	a	a	c	g	₁	0
a	a	c	g	\$	a	c	₁	0
a	c	a	a	c	g	\$	₁	0
a	c	g	\$	a	c	a	₁	1
c	a	a	c	g	\$	a	₂	2
c	g	\$	a	c	a	a	₃	3
g	\$	a	c	a	a	c	₂	3

Pattern search

In BWT find a number of given letter above a certain row

$\$$ ₁	a	c	a	a	c	g	₁	0
a	a	c	g	\$	a	c	₁	0
a	c	a	a	c	g	\$	₁	0
a	c	g	\$	a	c	a	₁	1
c	a	a	c	g	\$	a	₂	2
c	g	\$	a	c	a	a	₃	3
g	\$	a	c	a	a	c	₂	3

Pattern search

In BWT find a number of given letter above a certain row

$\$$ ₁	a	c	a	a	c	g	₁	0
a	a	c	g	\$	a	c	₁	1
a	c	a	a	c	g	\$	₁	1
a	c	g	\$	a	c	a	₁	1
c	a	a	c	g	\$	a	₂	1
c	g	\$	a	c	a	a	₃	1
g	\$	a	c	a	a	c	₂	2

Pattern search

In BWT find a number of given letter above a certain row

$\$$ ₁	a	c	a	a	c	g	g ₁	0
a ₁	a	c	g	$\$$	a	c ₁		1
a ₂	c	a	a	c	g	$\\$ ₁		1
a ₃	c	g	$\$$	a	c	a ₁		1
c ₁	a	a	c	g	$\$$	a ₂		1
c ₂	g	$\$$	a	c	a	a ₃		1
g ₁	$\$$	a	c	a	a	c ₂		2

Pattern search

In BWT find a number of given letter above a certain row

$\$$ ₁	a	c	a	a	c	g	₁	0
a	a	c	g	$\$$	a	c	₁	1
a	c	a	a	c	g	$\$$	₁	1
a	c	g	$\$$	a	c	a	₁	1
c	a	a	c	g	$\$$	a	₂	1
c	g	$\$$	a	c	a	a	₃	1
g	$\$$	a	c	a	a	c	₂	2

Pattern search

In BWT find a number of given letter above a certain row — $O(1)$

$\$$ ₁ a c a a c **g**₁

a₁ a c g \$ a **c**₁

a₂ c a a c g **\$**₁

a₃ c g \$ a c **a**₁

c₁ a a c g \$ **a**₂

c₂ g \$ a c a **a**₃

g₁ \$ a c a a **c**₂

Pattern search

In BWT find all given letters in specified range — $O(1)$

$\$$ ₁ a c a a c **g**₁

a₁ a c g \$ a **c**₁

a₂ c a a c g **\$**₁

a₃ c g \$ a c **a**₁

c₁ a a c g \$ **a**₂

c₂ g \$ a c a **a**₃

g₁ \$ a c a a **c**₂

Pattern search

Pattern search works in linear time

Thank you!

Questions?