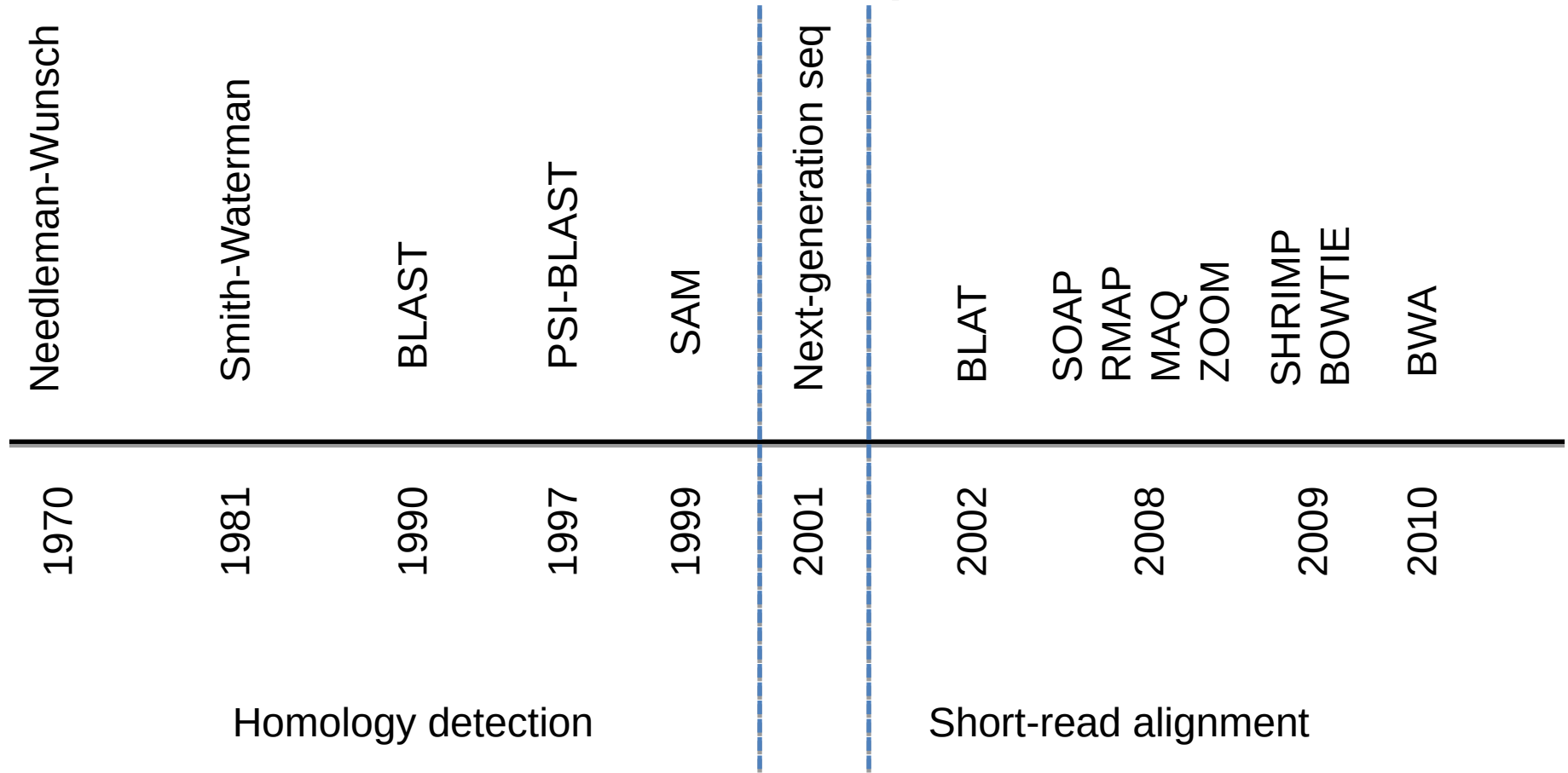


# Lecture 5

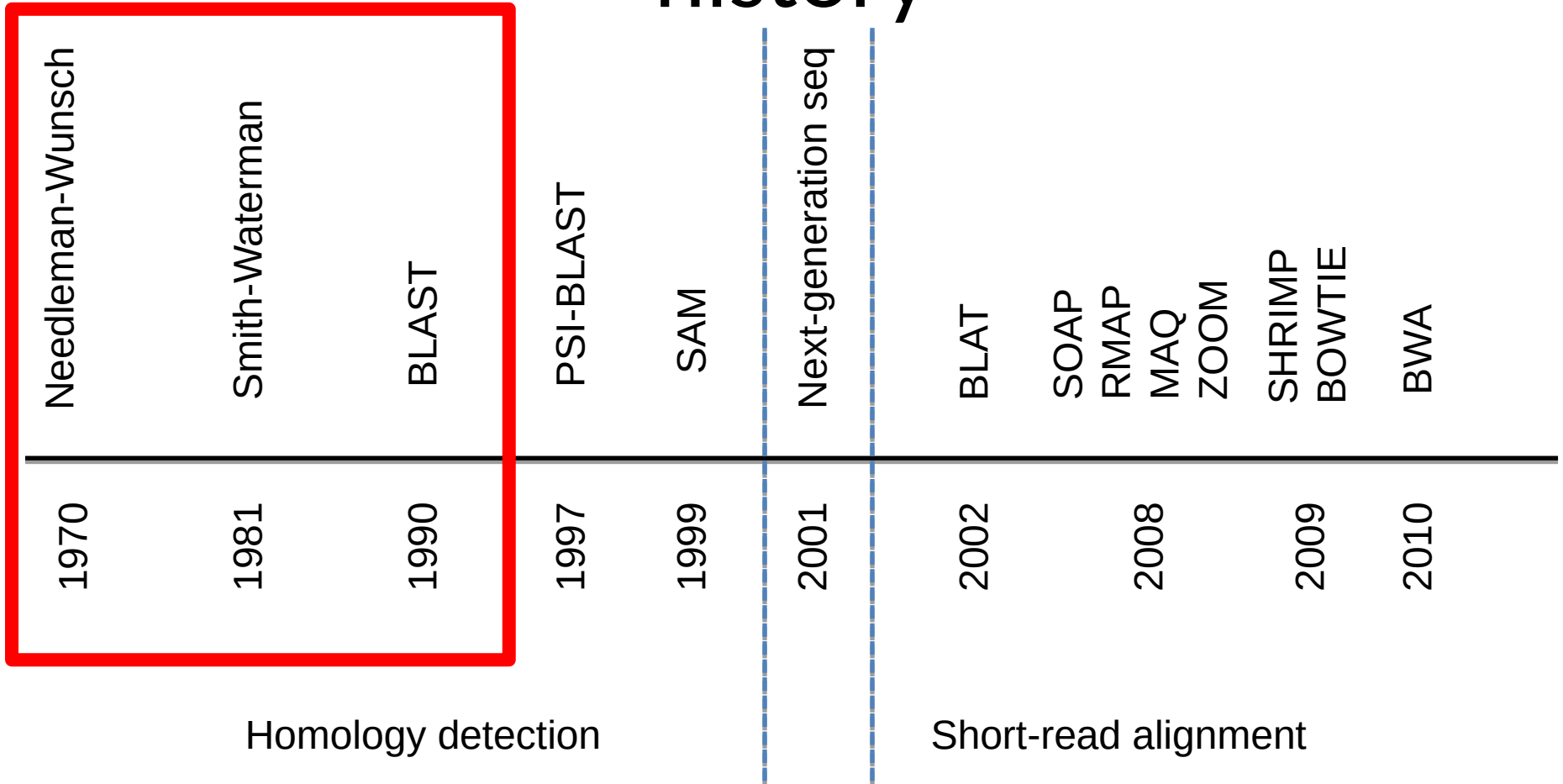
## Pairwise sequence alignment algorithms

Rachel Karchin  
BME 580.488, 580.688  
Spring 2019

# Some pairwise sequence alignment history



# Some pairwise sequence alignment history



# Alignment for homology detection

- When we compare sequences, we are looking for evidence that they have diverged from a common ancestor by a process of mutation and selection.

# Fundamentals

- Assume a simplified model of evolution
  - Substitutions
  - Deletions
  - Insertions

# Fundamentals

- Null hypothesis
  - Similarity between sequences is due to chance
- Alternative hypothesis
  - Similarity between sequences is due to a common ancestor

# Conceptualizing pairwise alignment

- Tabular representation
- Scoring system
- Search strategy

# Pairwise alignment basics

Two sequences

GACCC

GACAA

What is the best way to match up their shared equivalent positions?



# Pairwise alignment basics

Two sequences      GACCC  
                         GACAA

What is the best way to match up their shared equivalent positions?

GACCC	GACCC - -	GACC - C
GACAA	- - GACAA	GACAA -

G - ACCC	GACCC -	GACCC -
GACAA -	GA - CAA	GAC - AA

# Pairwise alignment basics

Two sequences

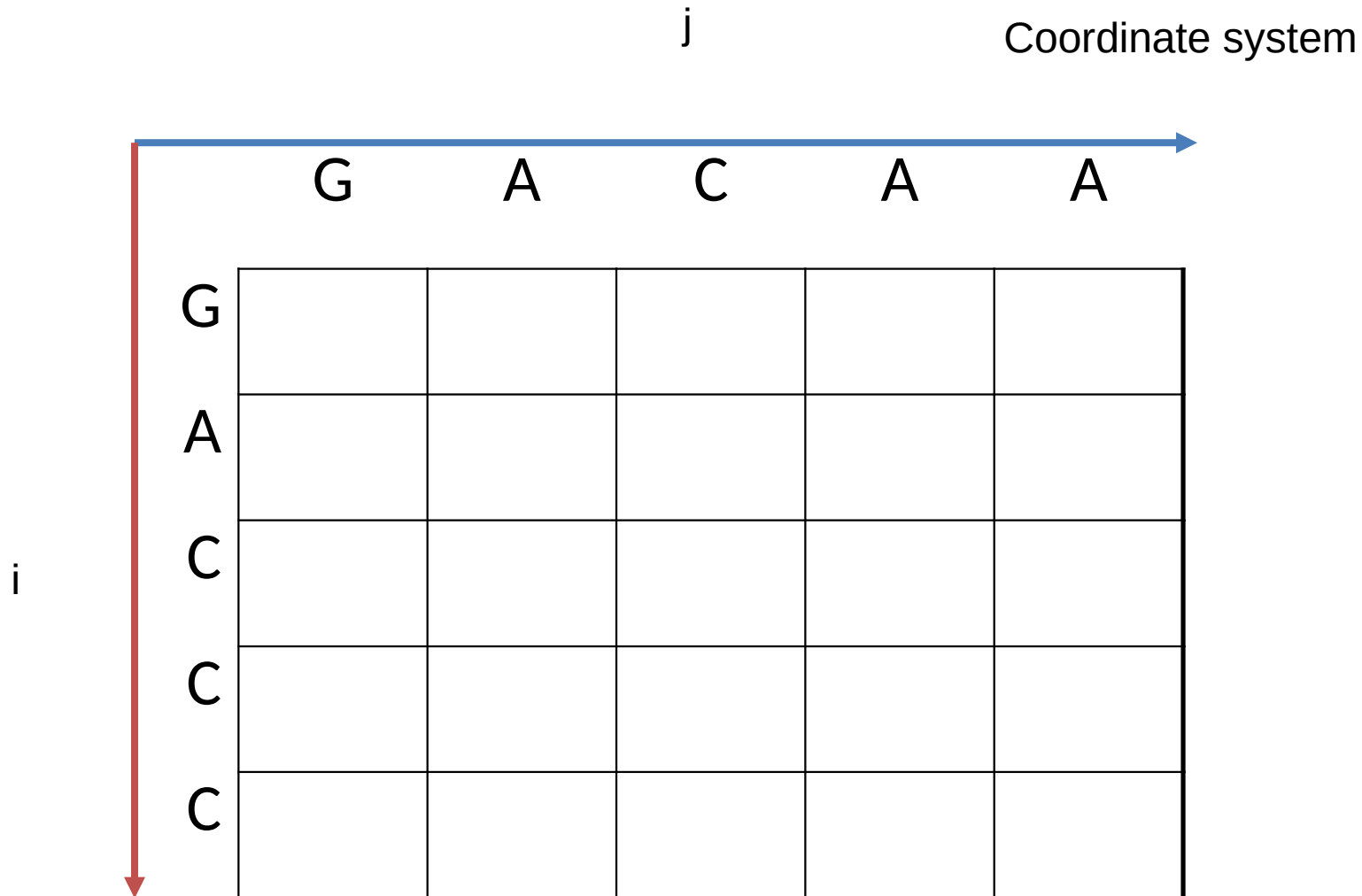
GACCC  
GACAA

Tabular representation

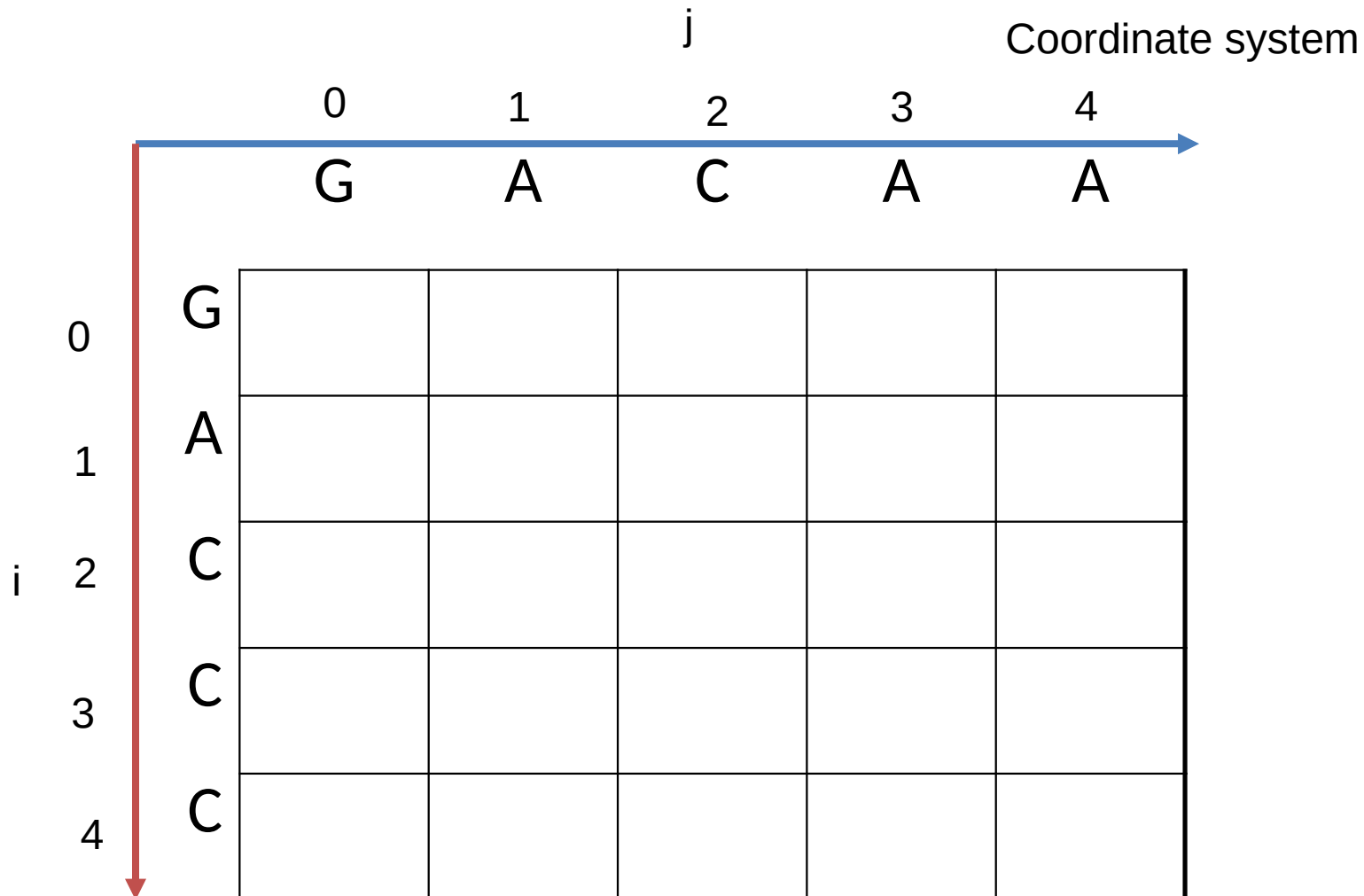
G      A      C      A      A

G					
A					
C					
C					
C					

# Pairwise alignment basics



# Pairwise alignment basics

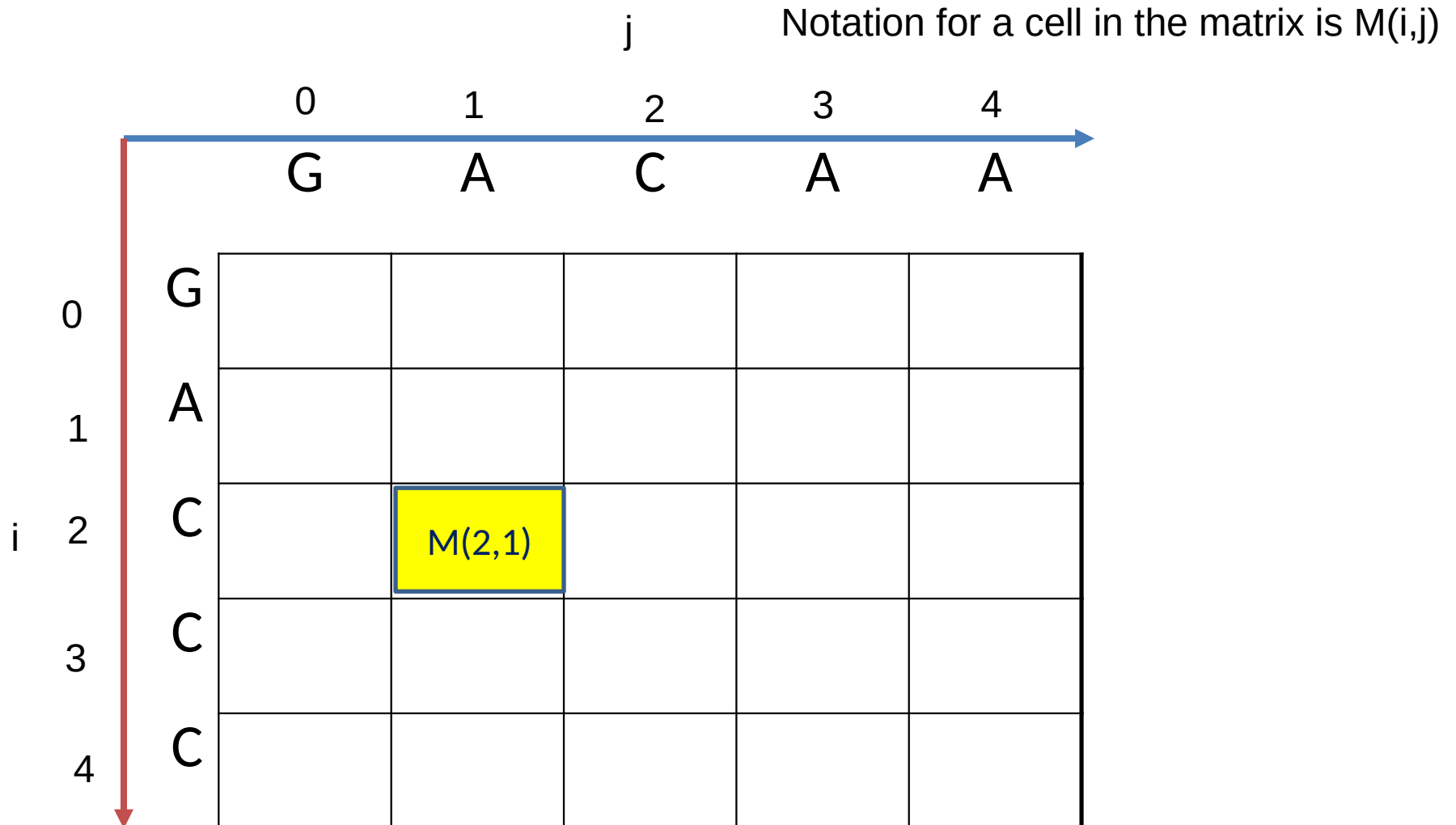


# Pairwise alignment basics

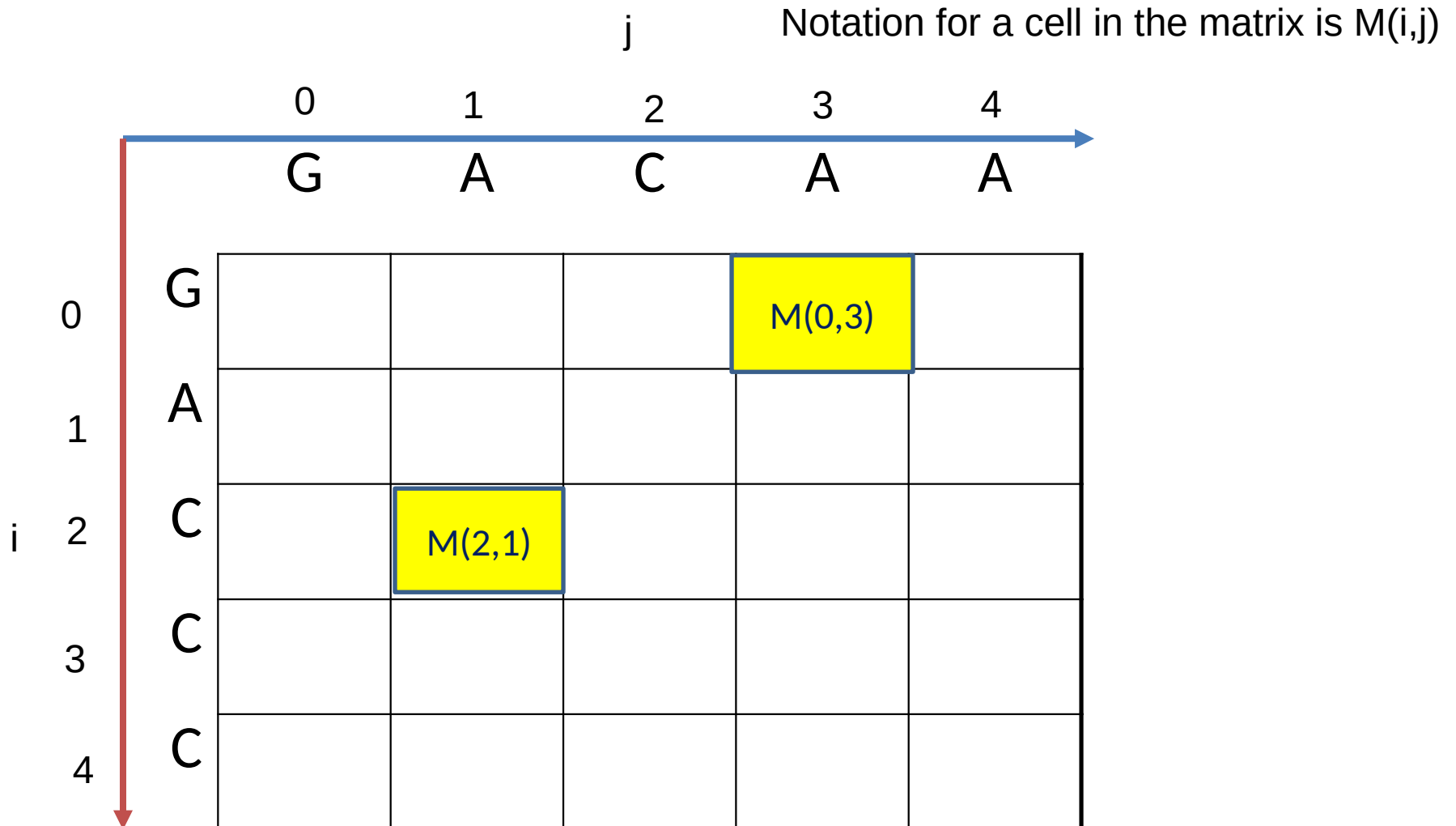
Notation for a cell in the matrix is  $M(i,j)$

		$j$				
		0	1	2	3	4
		G	A	C	A	A
$i$	0	G				
	1	A				
	2	C				
	3	C				
	4	C				

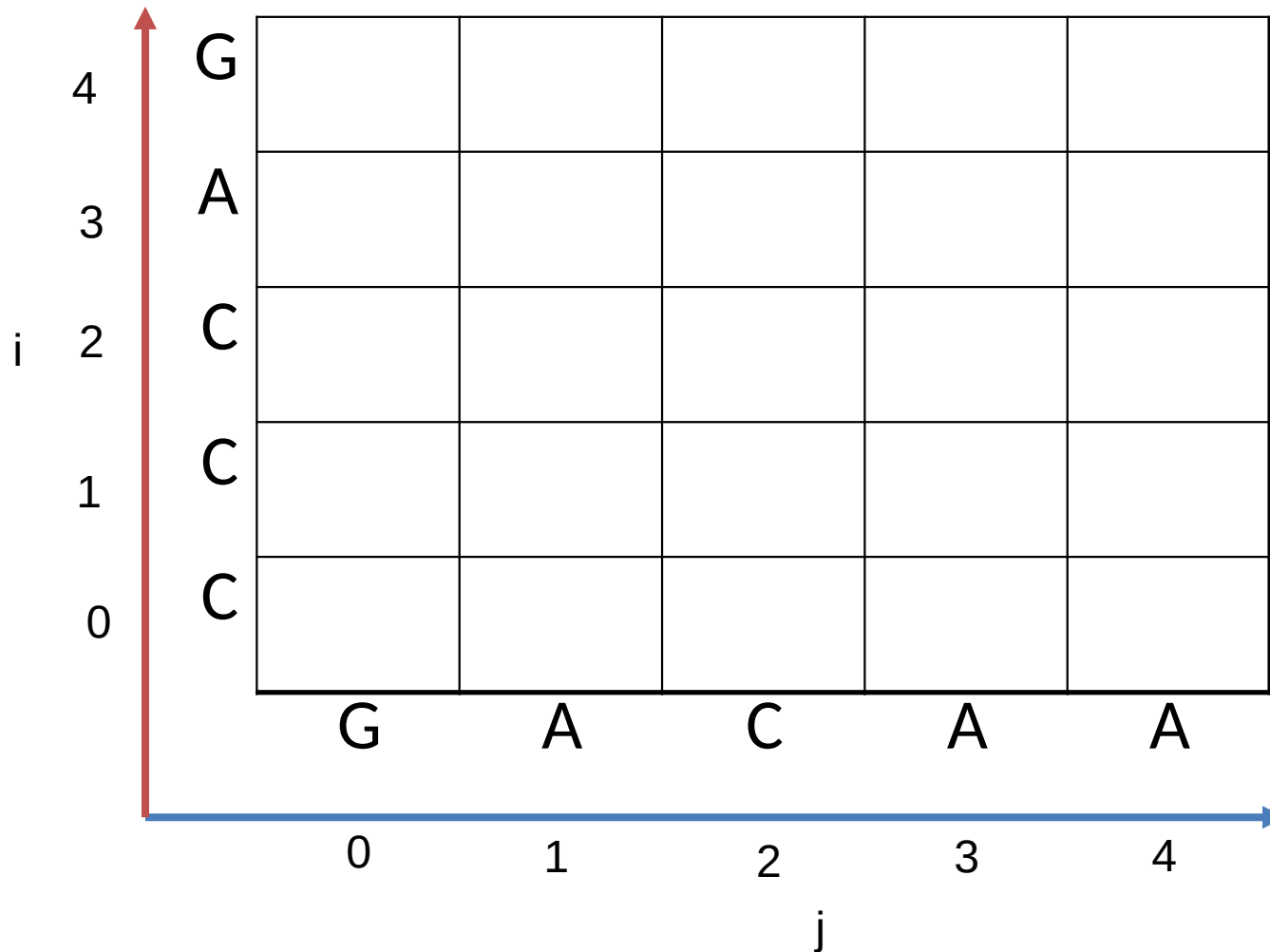
# Pairwise alignment basics



# Pairwise alignment basics



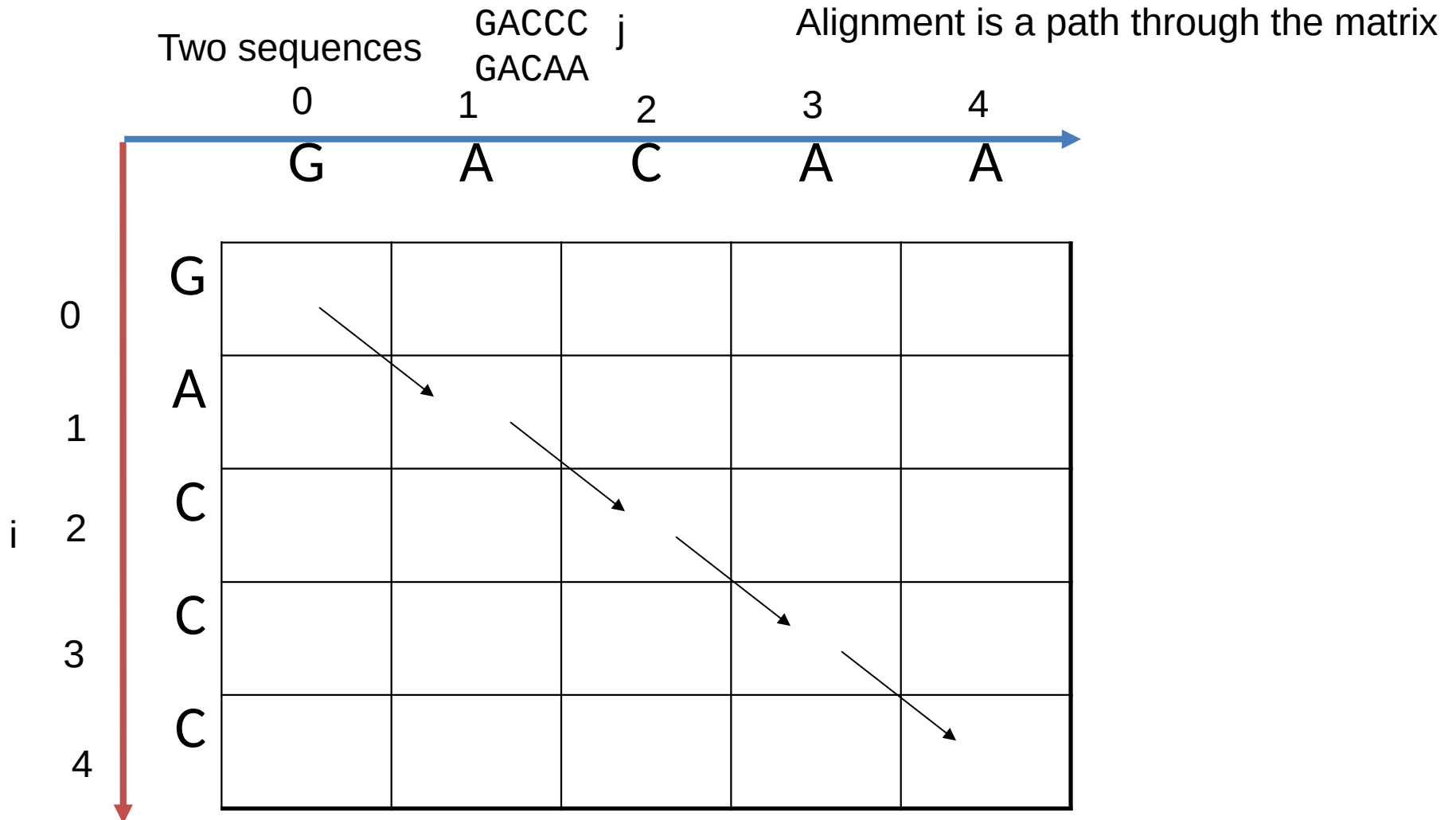
# Pairwise alignment basics



Coordinate system  
sometimes  
defined this way.



# Pairwise alignment basics



# Pairwise alignment basics

GACCC -  
G-ACAA

What path is this alignment?

G      A      C      A      A

G					
A					
C					
C					
C					

# Pairwise alignment basics

GACCC -  
G-ACAA

Let GACCC be the “reference”

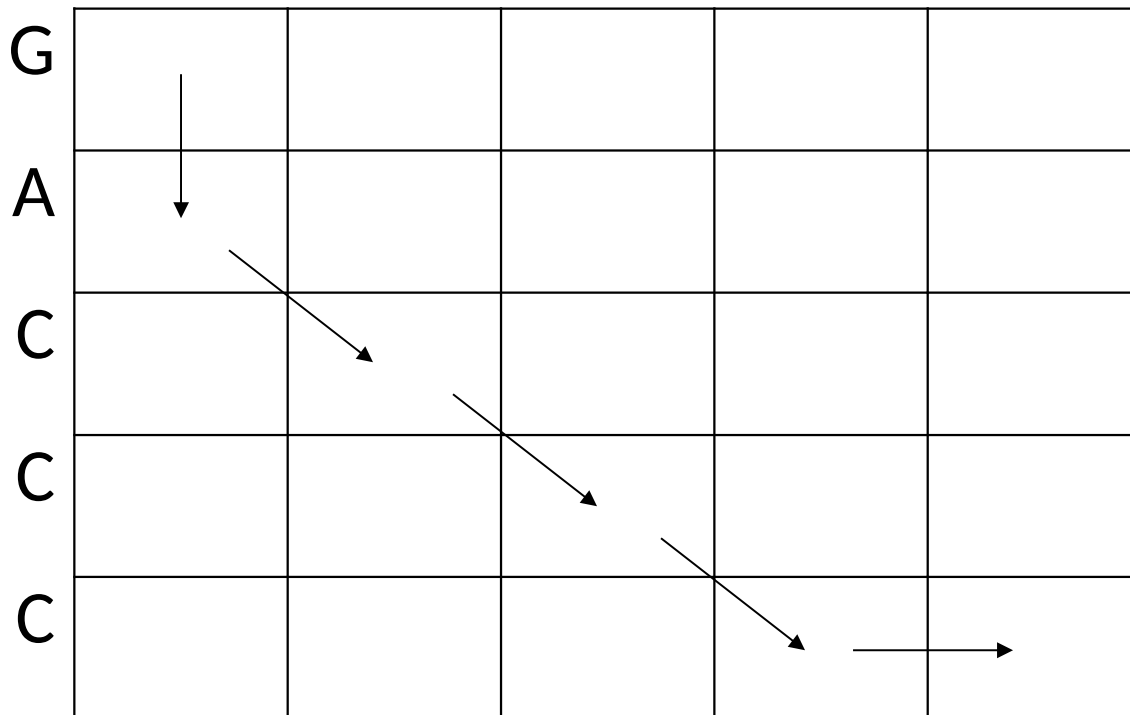
G      A      C      A      A

G					
A					
C					
C					
C					

# Pairwise alignment basics

GACCC -  
G-ACAA

G      A      C      A      A



# Pairwise alignment

- Giving the very large number of possible pairwise alignments for any two sequences, which one is best?

# Pairwise alignment

- Requires a **scoring system** that gives points or penalties for matches and gaps

# Pairwise alignment

- Also requires efficient **search strategy** to explore the very large alignment space.

# Scoring system

- “Match” points
- Gap costs (penalties) for insertions and deletions (indels)
- Any alignment can be scored by its matches (which are rewarded) and gap costs (which are penalized)

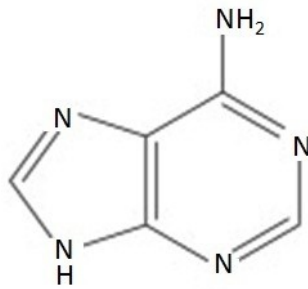


# What do we want from a scoring system?

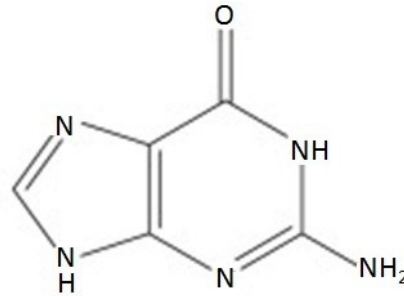
- For homology detection
  - Better scores for pairs of sequences that have diverged from a common ancestor
  - Worse scores for pairs that have not
- For short read alignment
  - Better scores for true matches
  - Worse scores for false matches

# DNA/RNA sequences are composed of nucleotides

purines



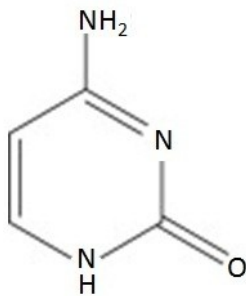
Adenine



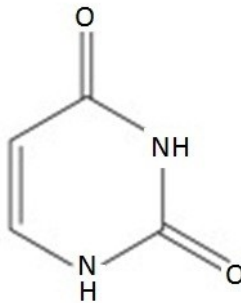
Guanine

DNA  
A  
G  
C  
T

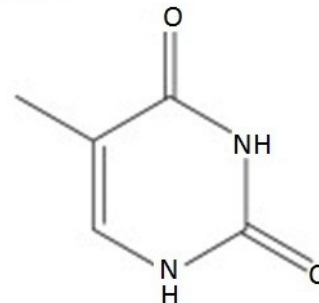
pyrimidines



Cytosine



Uracil



Thymine

RNA  
A  
G  
C  
U

4 nucleotide bases

# Scoring matrices

- “Score” associated with each possible type of “match” or substitution

	A	C	G	T
A	6	1	2	1
C	1	6	1	2
G	2	1	6	1
T	1	2	1	6

Kimura scoring matrix

# Gap penalties

- Constant gap penalty
- Linear gap penalty
- Affine gap penalty

# Constant Gap penalties

- Simplest approach
  - Cost for opening a gap  $\delta$
  - Every gap gets same penalty, regardless of length

# Linear Gap penalties

- Cost depends on the length of the gap
- Gap of length  $k$  gets cost of  $kc$
- Where  $c$  is the “unit” gap cost

# Affine Gap penalties

- A linear transformation followed by a translation
- Gap opening penalty  $\delta$  and linear gap extension penalty  $\varepsilon$
- Gap of length  $k$  gets penalty  $\delta + (k-1) \varepsilon$

# Gap penalties

- Constant gap penalty
- Linear gap penalty
- Affine gap penalty

Which works best?

How do we set parameters  $\delta$  and  $\epsilon$ ?



# Linear Gap penalties

- Each gap of one base gets penalized the same, regardless of whether or not it is contiguous with other gaps
- Alignment with fewer gaps favored
- Penalty for large gap is same as many small gaps

What do you think about this from a biological point of view?

$$\text{Gap cost} = \delta + (k-1)\varepsilon$$

# Affine Gap penalties

- Biologically, more likely to see a single gap of 10 than 10 single gaps
- We can model this using a gap opening penalty  $\delta$  and a gap extension penalty  $\varepsilon$
- Gap of length  $k$  gets penalty  $\delta + (k-1)\varepsilon$

Should  $\varepsilon$  be smaller than  $\delta$  ?  
Why?

# How can we efficiently search through alignment space?

# Dynamic programming for sequence alignment

- Key insight is that we can align prefixes and then extend them one position at a time
- Efficient search through the huge alignment space is done by caching intermediate results

# Global pairwise sequence alignment

- Global pairwise sequence alignment by dynamic programming introduced by Needleman and Wunsch in 1970.
- Known as the Needleman-Wunsch algorithm.

# Dynamic programming for sequence alignment

- General idea

A C A A

Matrix  
Initialization

	0	-2	-2	-2	-2
G	-2				
A	-2				
C	-2				
C	-2				
C	-2				

Two sequences

GACCC  
ACAA

# Dynamic programming for sequence alignment

- General idea

Aligning each nucleotide with the empty string is equivalent to putting a gap at the beginning of the alignment.

What kind of gap penalty is being used here?

Matrix  
Initialization

Why?

Two sequences

GACCC  
ACAA

	A	C	A	A
G	0	-2	-2	-2
A	-2			
C	-2			
C	-2			
C	-2			

# Dynamic programming for sequence alignment

- General idea

What is this?

		A	C	A	A
	0	-2	-2	-2	-2
G	-2				
A	-2				
C	-2				
C	-2				
C	-2				

Matrix  
Initialization

Two sequences

GACCC  
ACAA



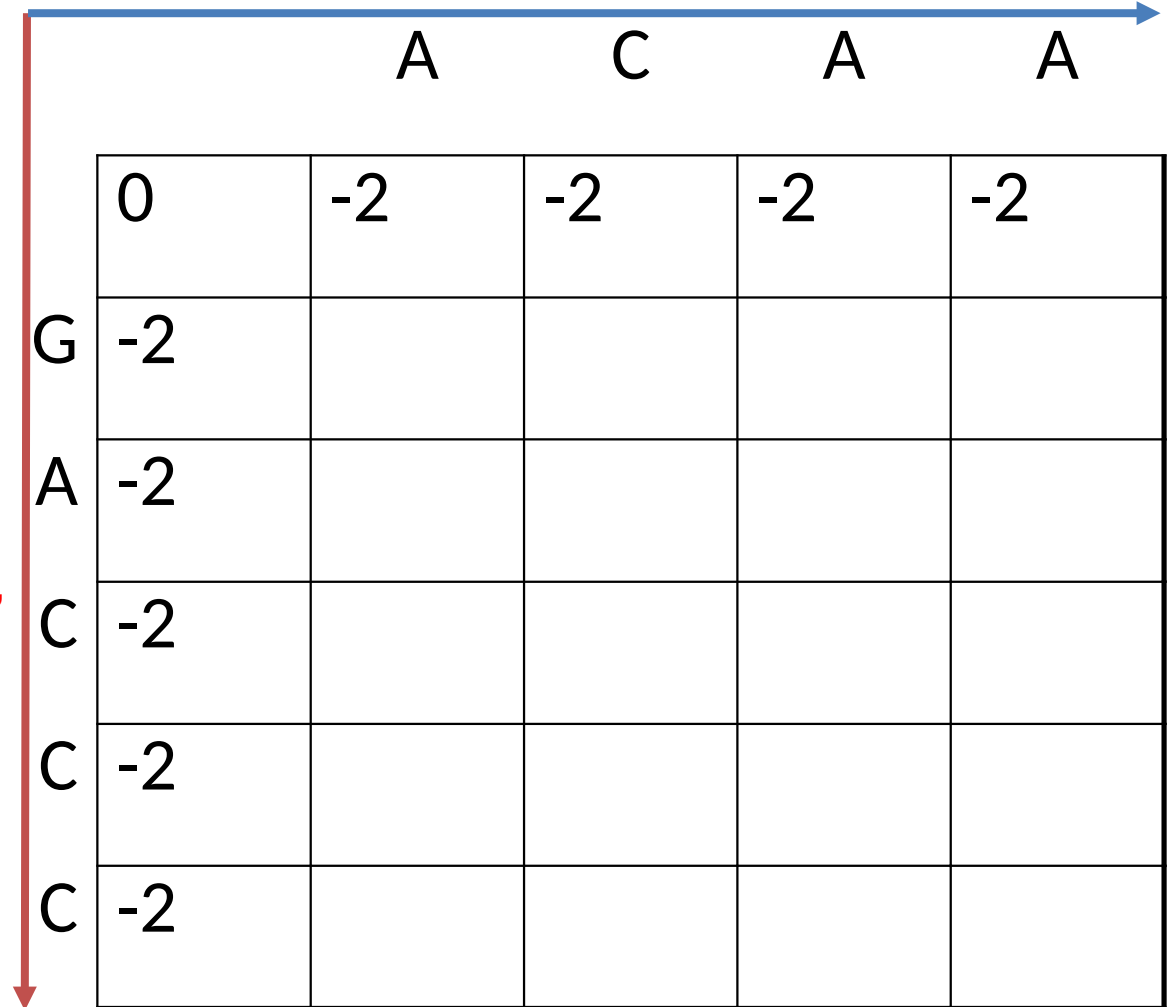
# Dynamic programming for sequence alignment

- General idea

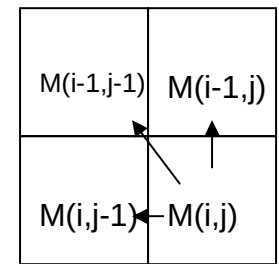
Matrix  
Fill

Recurrence relation

$$M(i,j) = \max( \\ M(i-1,j-1) + S(i,j), \\ M(i,j-1) + \omega(\text{gap in seq 1}), \\ M(i-1,j) + \omega(\text{gap in seq 2}) \\ )$$



		A	C	A	A
G	0	-2	-2	-2	-2
A	-2				
C	-2				
C	-2				
C	-2				



To fill in each  $M(i,j)$ , need to look where?

- General idea

Matrix  
Fill

Recurrence relation

$$M(i,j) = \max( \\ M(i-1,j-1) + S(i,j), \\ M(i,j-1) + \omega(\text{gap in seq 1}), \\ M(i-1,j) + \omega(\text{gap in seq 2}) \\ )$$

		A	C	A	A	
		0	-2	-2	-2	-2
G	-2					
A	-2					
C	-2					
C	-2					
C	-2					

# Dynamic programming for sequence alignment

- General idea

Matrix  
Fill

Recurrence relation

$$M(i,j) = \max( \\ M(i-1,j-1) + S(i,j), \\ M(i,j-1) + \omega(\text{gap in seq 1}), \\ M(i-1,j) + \omega(\text{gap in seq 2}) \\ )$$

$$M(1,1) = \max( \\ M(0,0) + S(G,A), \\ M(1,0) + \omega(\text{gap in seq 1}), \\ M(0,1) + \omega(\text{gap in seq 2}) \\ )$$

		A	C	A	A
	0	-2	-2	-2	-2
G	-2	M(1,1)			
A	-2				
C	-2				
C	-2				
C	-2				

# Dynamic programming for sequence alignment

- General idea

Matrix  
Fill

Recurrence relation

$$M(i,j) = \max( \\ M(i-1,j-1) + S(i,j), \\ M(i,j-1) + \omega(\text{gap in seq 1}), \\ M(i-1,j) + \omega(\text{gap in seq 2}) \\ )$$

$\omega(.) = -2$   
constant

$$M(1,1) = \max( \\ M(0,0) + S(G,A), \\ M(1,0) - 2, \\ M(0,1) - 2$$

		A	C	A	A
	0	-2	-2	-2	-2
G	-2	M(1,1)			
A	-2				
C	-2				
C	-2				
C	-2				

# Dynamic programming for sequence alignment

- General idea

Matrix  
Fill

Recurrence relation

Scoring system:

$S(i,j)$		Pur	Pyr
	Pur	2	-1
	Pyr	-1	2

$$M(1,1) = \max(\begin{aligned} &M(0,0) + 2, \\ &M(1,0) - 2, \\ &M(0,1) - 2 \end{aligned})$$

		A	C	A	A
	0	-2	-2	-2	-2
G	-2	M(1,1)			
A	-2				
C	-2				
C	-2				
C	-2				

# Dynamic programming for sequence alignment

Matrix  
Fill

Recurrence relation

$$M(i,j) = \max( \\ M(i-1,j-1) + S(i,j), \\ M(i,j-1) + \omega(\text{gap in seq 1}), \\ M(i-1,j) + \omega(\text{gap in seq 2}) \\ )$$

Scoring system:

	Pur	Pyr
S(i,j)		
Pur	2	-1
Pyr	-1	2

$$M(1,1) = \max( \\ 0 + 2, \\ -2 - 2, \\ -2 - 2$$

A C A A

	0	-2	-2	-2	-2
G	-2	2			
A	-2				
C	-2				
C	-2				
C	-2				

# Dynamic programming for sequence alignment

Matrix  
Fill

Recurrence relation

$$M(i,j) = \max( \\ M(i-1,j-1) + S(i,j), \\ M(i,j-1) + \omega(\text{gap in seq 1}), \\ M(i-1,j) + \omega(\text{gap in seq 2}) \\ )$$

Scoring system:

	Pur	Pyr
S(i,j)		
Pur	2	-1
Pyr	-1	2

$$M(1,1) = \max( \\ 0 + 2, \\ -2 - 2, \\ -2 - 2$$

A C A A

	0	-2	-2	-2	-2
G	-2	2			
A	-2				
C	-2				
C	-2				
C	-2				

Add a pointer to trace the step along the path we chose

# Dynamic programming for sequence alignment

Matrix  
Fill

Recurrence relation

$$M(i,j) = \max( \\ M(i-1,j-1) + S(i,j), \\ M(i,j-1) + \omega(\text{gap in seq 1}), \\ M(i-1,j) + \omega(\text{gap in seq 2}) \\ )$$

$$M(1,2) = \max( \\ -2 - 1, \\ 2 - 2, \\ -2 - 2$$

		A	C	A	A
	0	-2	-2	-2	-2
G	-2	2	0		
A	-2				
C	-2				
C	-2				
C	-2				



# Dynamic programming for sequence alignment

Matrix  
Fill

Recurrence relation

$$M(i,j) = \max( \\ M(i-1,j-1) + S(i,j), \\ M(i,j-1) + \omega(\text{gap in seq 1}), \\ M(i-1,j) + \omega(\text{gap in seq 2}) \\ )$$

$$M(1,2) = \max( \\ -2 - 1, \\ 2 - 2, \\ -2 - 2$$

		A	C	A	A
	0	-2	-2	-2	-2
G	-2	2	0		
A	-2				
C	-2				
C	-2				
C	-2				

# Dynamic programming for sequence alignment

Matrix  
Fill

Recurrence relation

$$M(i,j) = \max( \\ M(i-1,j-1) + S(i,j), \\ M(i,j-1) + \omega(\text{gap in seq 1}), \\ M(i-1,j) + \omega(\text{gap in seq 2}) \\ )$$

$$M(1,3) = \max( \\ -2 + 2, \\ 0 - 2, \\ -2 - 2$$

		A	C	A	A
	0	-2	-2	-2	-2
G	-2	2	0	0	
A	-2				
C	-2				
C	-2				
C	-2				

# Dynamic programming for sequence alignment

Matrix  
Fill

Recurrence relation

$$M(i,j) = \max( \\ M(i-1,j-1) + S(i,j), \\ M(i,j-1) + \omega(\text{gap in seq 1}), \\ M(i-1,j) + \omega(\text{gap in seq 2}) \\ )$$

$$M(1,3) = \max( \\ -2 + 2, \\ 0 - 2, \\ -2 - 2$$

		A	C	A	A
	0	-2	-2	-2	-2
G	-2	2	0	0	
A	-2				
C	-2				
C	-2				
C	-2				

# Dynamic programming for sequence alignment

Matrix  
Fill

Recurrence relation

$$M(i,j) = \max( \\ M(i-1,j-1) + S(i,j), \\ M(i,j-1) + \omega(\text{gap in seq 1}), \\ M(i-1,j) + \omega(\text{gap in seq 2}) \\ )$$

		A	C	A	A
	0	-2	-2	-2	-2
G	-2	2	0	0	0
A	-2				
C	-2				
C	-2				
C	-2				

# Dynamic programming for sequence alignment

Matrix  
Fill

Recurrence relation

$$M(i,j) = \max( \\ M(i-1,j-1) + S(i,j), \\ M(i,j-1) + \omega(\text{gap in seq 1}), \\ M(i-1,j) + \omega(\text{gap in seq 2}) \\ )$$

		A	C	A	A
	0	-2	-2	-2	-2
G	-2	2	0	0	0
A	-2				
C	-2				
C	-2				
C	-2				

# Dynamic programming for sequence alignment

Matrix  
Fill

Recurrence relation

$$M(i,j) = \max( \\ M(i-1,j-1) + S(i,j), \\ M(i,j-1) + \omega(\text{gap in seq 1}), \\ M(i-1,j) + \omega(\text{gap in seq 2}) \\ )$$

		A	C	A	A
	0	-2	-2	-2	-2
G	-2	2	0	0	0
A	-2	0			
C	-2				
C	-2				
C	-2				

# Dynamic programming for sequence alignment

Matrix  
Fill

Recurrence relation

$$M(i,j) = \max( \\ M(i-1,j-1) + S(i,j), \\ M(i,j-1) + \omega(\text{gap in seq 1}), \\ M(i-1,j) + \omega(\text{gap in seq 2}) \\ )$$

		A	C	A	A
	0	-2	-2	-2	-2
G	-2	2	0	0	0
A	-2	0			
C	-2				
C	-2				
C	-2				

The diagram illustrates the sequence alignment process. The matrix shows the scores for aligning the sequence 'GACAC' (rows) with the sequence 'AACA' (columns). The path of maximum alignment is highlighted with arrows: from (0,0) to (1,1) (diagonal), from (1,1) to (1,2) (horizontal), from (1,2) to (2,2) (vertical), from (2,2) to (2,3) (diagonal), and from (2,3) to (2,4) (diagonal). This path corresponds to the alignment: G aligned with A, A aligned with A, C aligned with C, and the remaining C and A in the first sequence aligned with gaps.

# Dynamic programming for sequence alignment

Matrix  
Fill

Recurrence relation

$$M(i,j) = \max( \\ M(i-1,j-1) + S(i,j), \\ M(i,j-1) + \omega(\text{gap in seq 1}), \\ M(i-1,j) + \omega(\text{gap in seq 2}) \\ )$$

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



# Dynamic programming for sequence alignment

Matrix  
Fill

Recurrence relation

$$M(i,j) = \max( \\ M(i-1,j-1) + S(i,j), \\ M(i,j-1) + \omega(\text{gap in seq 1}), \\ M(i-1,j) + \omega(\text{gap in seq 2}) \\ )$$

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

# alignment

$$M(i,j) = \max( \\ M(i-1,j-1) + S(i,j), \\ M(i,j-1) + \omega(\text{gap in seq 1}), \\ M(i-1,j) + \omega(\text{gap in seq 2}) \\ )$$

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

# Dynamic programming for sequence alignment

Complexity?

Matrix  
Fill

Recurrence relation

$$M(i,j) = \max( \\ M(i-1,j-1) + S(i,j), \\ M(i,j-1) + \omega(\text{gap in seq 1}), \\ M(i-1,j) + \omega(\text{gap in seq 2}) \\ )$$

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

# Dynamic programming for sequence alignment

Traceback

Where do we start?

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

# Dynamic programming for sequence alignment

Traceback

For global alignment

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

# Dynamic programming for sequence alignment

Traceback

For global alignment

In this coordinate system,  
bottom right cell gives  
us the optimal global  
alignment score

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

# Dynamic programming for sequence alignment

Traceback

Want to identify the alignment over the length of both sequences

Backtrack over the winning path

How?

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

# Dynamic programming for sequence alignment

Traceback

Want to recover the (an) optimal alignment over the length of both sequences

Backtrack over the winning path

How?

Have to identify the direction into which each cell was reached.

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



# Dynamic programming for sequence alignment

Traceback

Want to recover the (an)  
optimal alignment over the  
length of both sequences

Backtrack over the winning  
path

How?

Have to identify the direction  
into which each cell was  
reached.

Forward pointers tell us

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

# Dynamic programming for sequence alignment

Traceback

Follow the path backwards

Identify the (a) winning path

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

# Dynamic programming for sequence alignment

Traceback

Follow the path backwards

Identify the (a) winning path

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

# Dynamic programming for sequence alignment

Traceback

Follow the path backwards

Identify the (a) winning path

What does our alignment look like at this point?

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

# Dynamic programming for sequence alignment

Traceback

Follow the path backwards

Identify the (a) winning path

What does our alignment look like at this point?

CCC  
CAA

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

# Dynamic programming for sequence alignment

Traceback

Follow the path backwards

Identify the (a) winning path

ACCC  
-CAA

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

# Dynamic programming for sequence alignment

Traceback

Follow the path backwards

Identify the (a) winning path

GACCC  
A - CAA

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

- Global sequence alignment is not always desirable.
  - Useful when two sequences are very similar
- When there are short regions of high similarity and longer regions of low similarity, we want to use **local sequence alignment**



# Local pairwise sequence alignment

- Ten years after Needleman-Wunsch
- Smith-Waterman (1981)

# Local pairwise sequence alignment

- General idea

Recurrence relation

$$M(i,j) = \max(0, M(i-1,j-1) + S(i,j), M(i,j-1) + \omega(\text{gap in seq 1}), M(i-1,j) + \omega(\text{gap in seq 2}))$$

Scoring system:

	Pur	Pyr
Pur	2	-1
Pyr	-1	2

$\omega(.) = -2$   
constant

		C	C	A	A
	0	0	0	0	0
G	0				
A	0				
C	0				
C	0				
C	0				

# Local pairwise sequence alignment

Recurrence relation

$$M(i,j) = \max(0, M(i-1,j-1) + S(i,j), M(i,j-1) + \omega(\text{gap in seq 1}), M(i-1,j) + \omega(\text{gap in seq 2}))$$

Zero “win” => start of a new alignment (reset)

Alignment can begin and end anywhere in the matrix

		C	C	A	A
		0	0	0	0
G		0			
A		0			
C		0			
C		0			
C		0			

# Local pairwise sequence alignment

Recurrence relation

$$M(i,j) = \max(0, M(i-1,j-1) + S(i,j), M(i,j-1) + \omega(\text{gap in seq 1}), M(i-1,j) + \omega(\text{gap in seq 2}))$$

Scoring system:

S(i,j)		Pur	Pyr
	Pur	2	-1
	Pyr	-1	2

$\omega(.) = -2$   
constant

		C	C	A	A
		0	0	0	0
G	0	0			
A	0				
C	0				
C	0				
C	0				

# Local pairwise sequence alignment

Do we set a pointer from the parent to this cell?

Recurrence relation

$$M(i,j) = \max(0, M(i-1,j-1) + S(i,j), M(i,j-1) + \omega(\text{gap in seq 1}), M(i-1,j) + \omega(\text{gap in seq 2}))$$

Scoring system:

	Pur	Pyr
Pur	2	-1
Pyr	-1	2

$\omega(.) = -2$   
constant

		C	C	A	A
	0	0	0	0	0
G	0	0			
A	0				
C	0				
C	0				
C	0				

# Local pairwise sequence alignment

Recurrence relation

$$M(i,j) = \max(0, M(i-1,j-1) + S(i,j), M(i,j-1) + \omega(\text{gap in seq 1}), M(i-1,j) + \omega(\text{gap in seq 2}))$$

Scoring system:

S(i,j)		Pur	Pyr
	Pur	2	-1
	Pyr	-1	2

$\omega(.) = -2$   
constant

		C	C	A	A
		0	0	0	0
G	0	0	0		
A	0				
C	0				
C	0				
C	0				

# Local pairwise sequence alignment

Recurrence relation

$$M(i,j) = \max(0, M(i-1,j-1) + S(i,j), M(i,j-1) + \omega(\text{gap in seq 1}), M(i-1,j) + \omega(\text{gap in seq 2}))$$

Scoring system:

S(i,j)		Pur	Pyr
	Pur	2	-1
	Pyr	-1	2

$\omega(.) = -2$   
constant

		C	C	A	A
		0	0	0	0
G	0	0	0	2	
A	0				
C	0				
C	0				
C	0				

# Local pairwise sequence alignment

Do we set a pointer from the parent to this cell?

Recurrence relation

$$M(i,j) = \max(0, M(i-1,j-1) + S(i,j), M(i,j-1) + \omega(\text{gap in seq 1}), M(i-1,j) + \omega(\text{gap in seq 2}))$$

Scoring system:

	Pur	Pyr
Pur	2	-1
Pyr	-1	2

$\omega(.) = -2$   
constant

		C	C	A	A
		0	0	0	0
G	0	0	0	2	
A	0				
C	0				
C	0				
C	0				



# Local pairwise sequence alignment

Recurrence relation

$$M(i,j) = \max(0, M(i-1,j-1) + S(i,j), M(i,j-1) + \omega(\text{gap in seq 1}), M(i-1,j) + \omega(\text{gap in seq 2}))$$

Scoring system:

	Pur	Pyr
Pur	2	-1
Pyr	-1	2

$\omega(.) = -2$   
constant

		C	C	A	A
		0	0	0	0
G	0	0	0	2	2
A	0	0	0	2	4
C	0	2	2	0	2
C	0	2	4	2	0
C	0	2	4	3	1

# Local pairwise sequence alignment

Recurrence relation

$$M(i,j) = \max(0, M(i-1,j-1) + S(i,j), M(i,j-1) + \omega(\text{gap in seq 1}), M(i-1,j) + \omega(\text{gap in seq 2}))$$

Scoring system:

	Pur	Pyr
Pur	2	-1
Pyr	-1	2

$\omega(.) = -2$   
constant

		C	C	A	A
		0	0	0	0
G	0	0	0	2	2
A	0	0	0	2	4
C	0	2	2	0	2
C	0	2	4	2	0
C	0	2	4	3	1

# Local pairwise sequence alignment

Traceback

Want to recover the (an)  
optimal alignment ~~over the~~  
~~length of both sequences~~

Backtrack over the winning  
paths

How?

Where do we start?

		C	C	A	A
	0	0	0	0	0
G	0	0	0	2	2
A	0	0	0	2	4
C	0	2	2	0	2
C	0	2	4	2	0
C	0	2	4	3	1

# Local pairwise sequence alignment

Traceback

Want to recover the (an)  
optimal alignment ~~over the~~  
~~length of both sequences~~

Backtrack over the winning  
paths

How?

Where do we start?

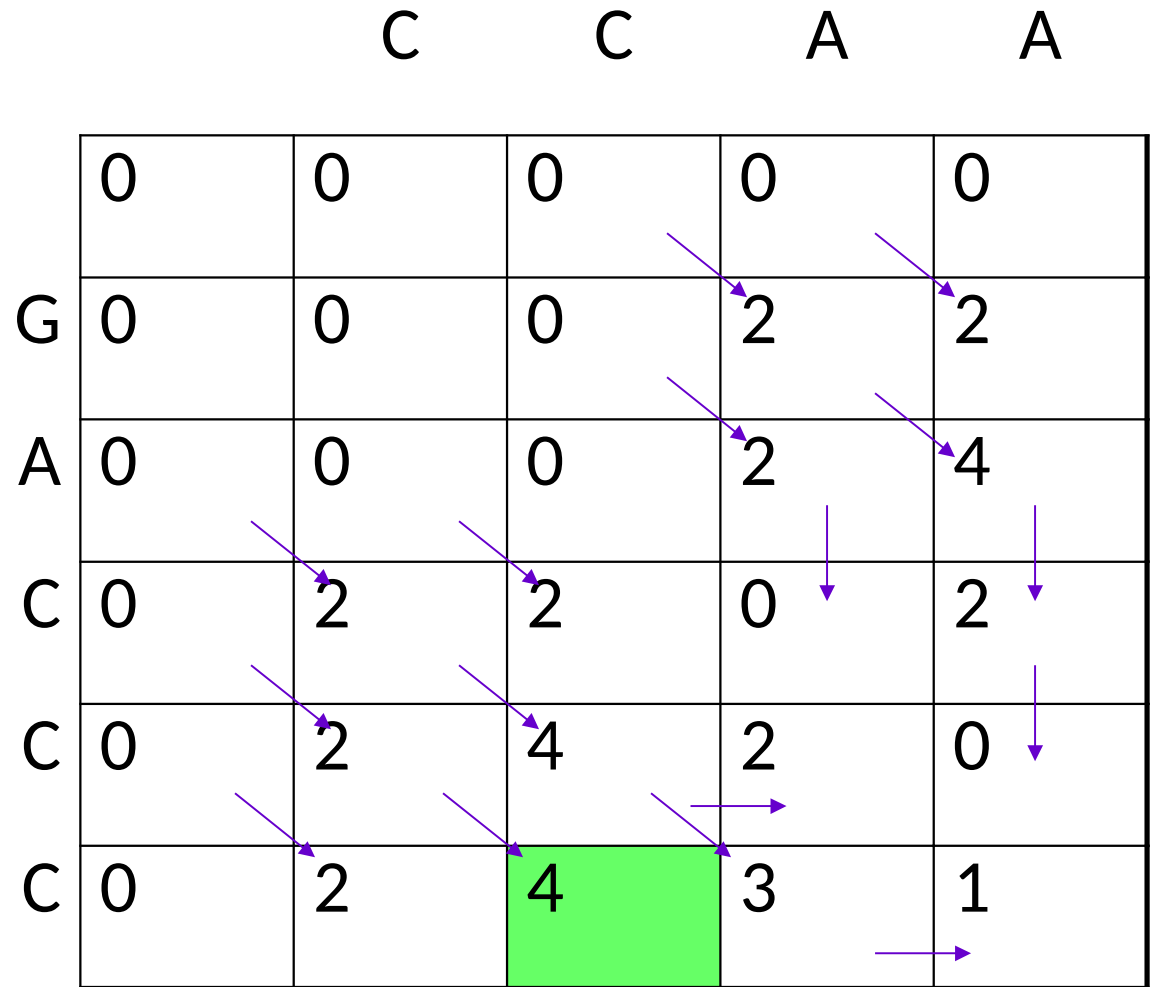
Largest value in the matrix  
is score of the best local  
alignment between  
subsequences

		C	C	A	A
	0	0	0	0	0
G	0	0	0	2	2
A	0	0	0	2	4
C	0	2	2	0	2
C	0	2	4	2	0
C	0	2	4	3	1

# Local pairwise sequence alignment

Traceback

		C	C	A	A
	0	0	0	0	0
G	0	0	0	2	2
A	0	0	0	2	4
C	0	2	2	0	2
C	0	2	4	2	0
C	0	2	4	3	1



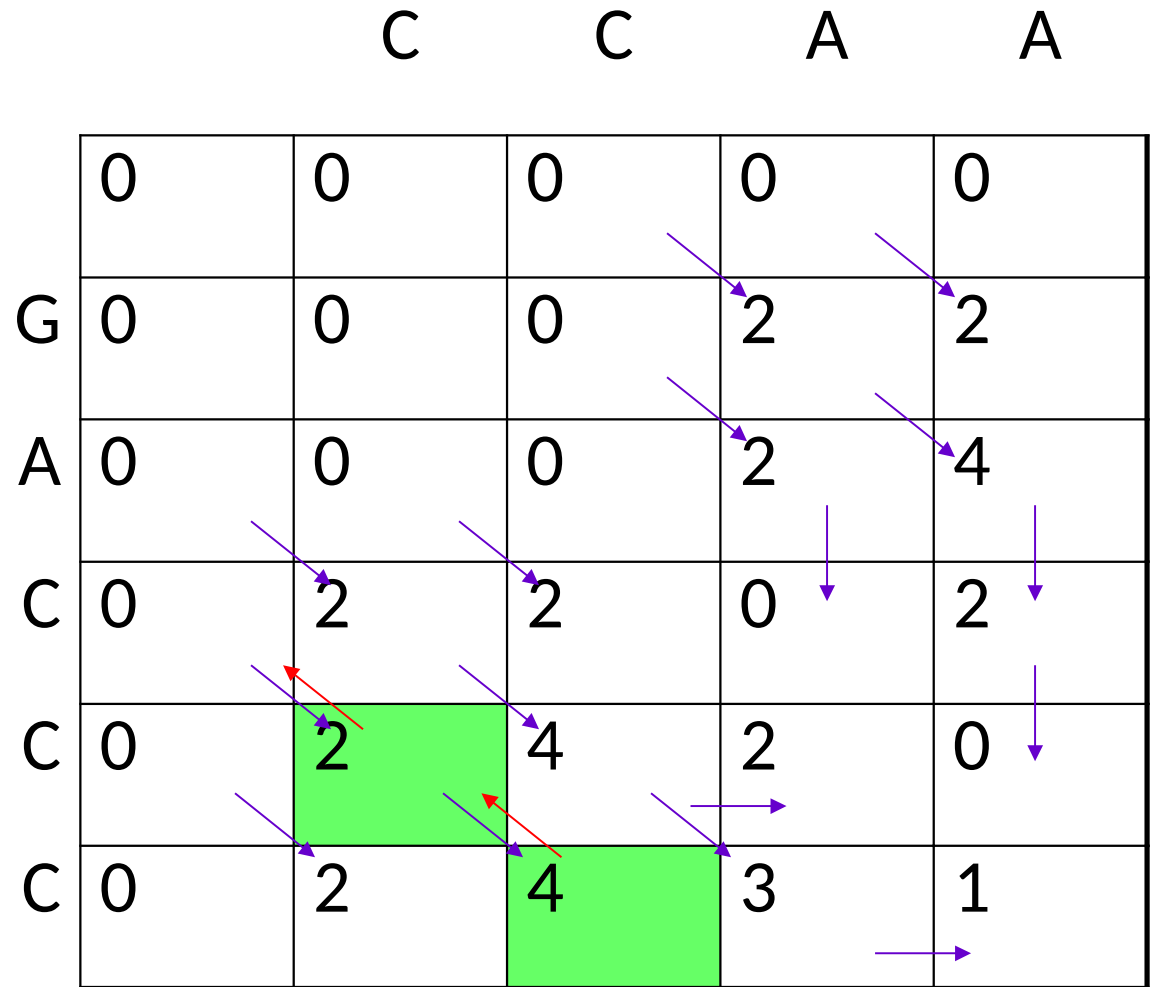
# Local pairwise sequence alignment

Traceback

		C	C	A	A
	0	0	0	0	0
G	0	0	0	2	2
A	0	0	0	2	4
C	0	2	2	0	2
C	0	2	4	2	0
C	0	2	4	3	1

# Local pairwise sequence alignment

Traceback



# Local pairwise sequence alignment

Traceback

What does this alignment look like?

		C	C	A	A
	0	0	0	0	0
G	0	0	0	2	2
A	0	0	0	2	4
C	0	2	2	0	2
C	0	2	4	2	0
C	0	2	4	3	1



# Local pairwise sequence alignment

Traceback

What does this alignment look like?

CC  
CC

		C	C	A	A
	0	0	0	0	0
G	0	0	0	2	2
A	0	0	0	2	4
C	0	2	2	0	2
C	0	2	4	2	0
C	0	2	4	3	1

Diagram illustrating a local pairwise sequence alignment between the sequences C C A A (columns) and G A C C C (rows). The alignment scores are shown in the table. The optimal alignment path is highlighted in green, showing a match between the second 'C' in the column sequence and the first 'C' in the row sequence, followed by a match between the third 'C' in the column sequence and the second 'C' in the row sequence. The alignment ends at the second 'C' in the row sequence.

# Local pairwise sequence alignment

Traceback

For this to work, what is requirement for expected value of a random match?

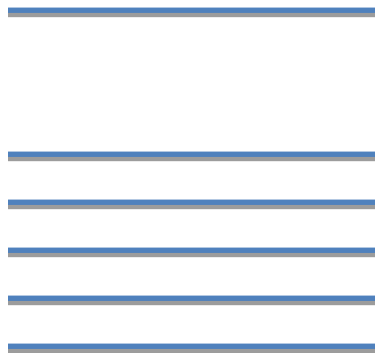
Why?

		C	C	A	A
	0	0	0	0	0
G	0	0	0	2	2
A	0	0	0	2	4
C	0	2	2	0	2
C	0	2	4	2	0
C	0	2	4	3	1

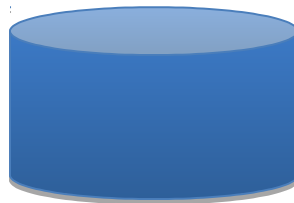
# Trends in homology-based pairwise sequence alignment



Align two sequences



Align a target sequence with everything in a sequence database



Exponential growth in database size

# Motivation for heuristic alignments

- For two sequences of length  $m$  and  $n$ 
  - Time to get optimal alignment score is  $\Theta(mn)$ 
    - $\Theta$  notation like  $O$  notation but provides upper and lower bounds on asymptotic behavior ( $O$  is only for upper)
  - Space required to recover the alignment (naively  $\Theta(mn)$  but can reduce to  $\Theta(m+n)$ )
- Not realistic for large scale problems

# BLAST

## Key insight:

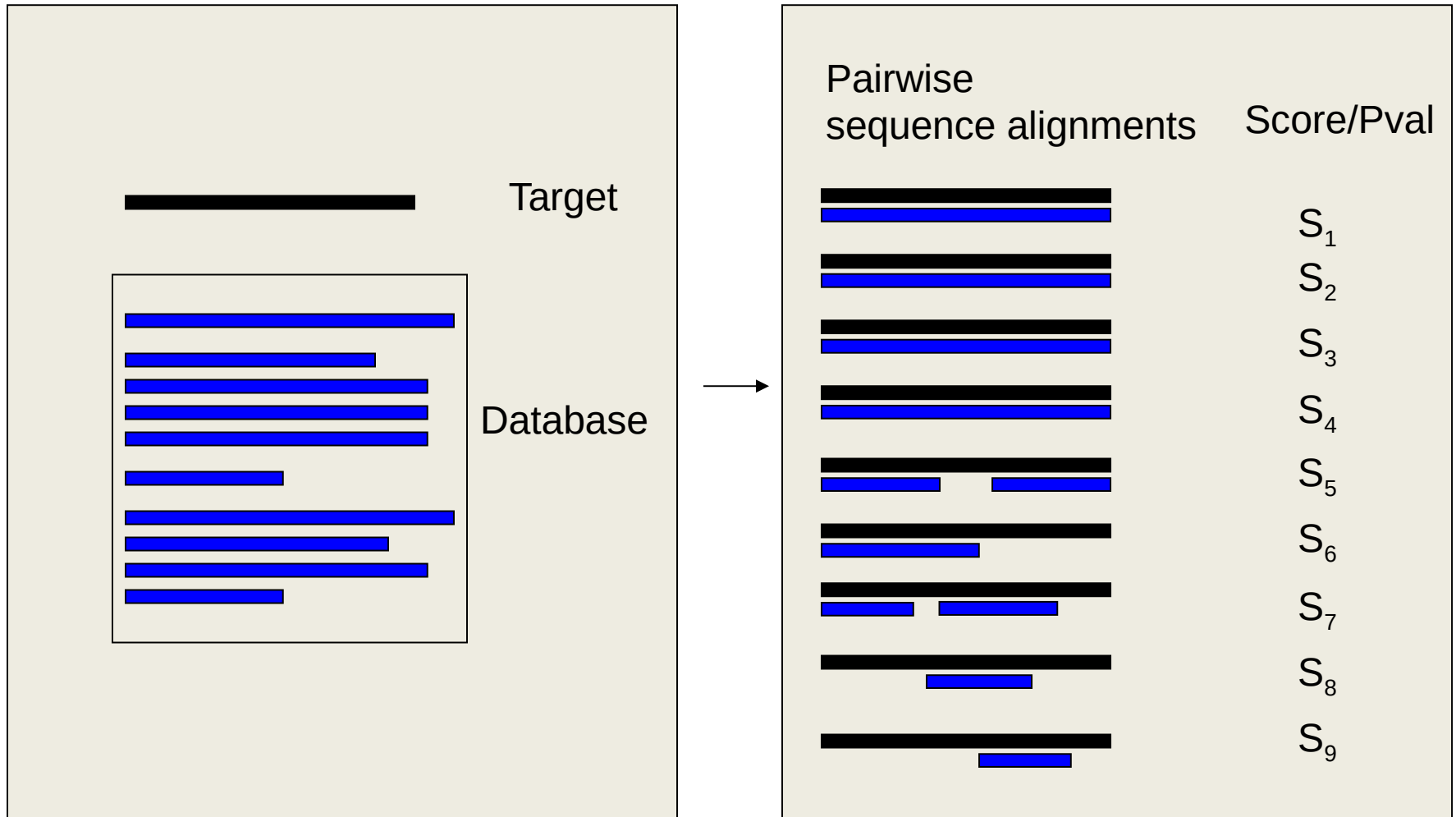
Biologically real alignments very likely to contain a short stretch of identities or very high scoring matches.



## Algorithm

1. Build "words" — find short statistically significant sub-sequences in query
2. Find "seeds" — scan sequences in database for matching words
3. Extend — use (nearby) seeds to form local alignments called HSPs
4. Score — combine groups of consistent HSPs into local alignment with best score

# BLAST

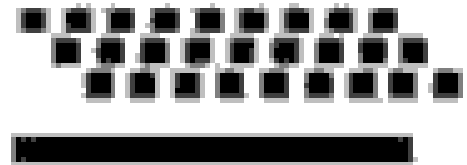


# BLAST

- Key parameters:
  - $W$  (word length)
  - $T$  (threshold for scores in the initial word list)

# BLAST word list

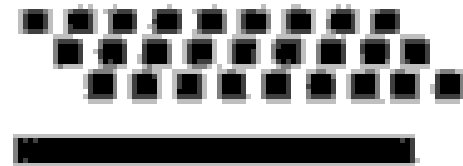
1. Build "words" — find short statistically significant sub-sequences in query





# BLAST word list

1. Build "words" — find short statistically significant sub-sequences in query



*Q* ACTGA

*D* GACTGC

Word List

$W=3$

ACT  
CTG  
TGA

# BLAST database seed finding

2. Find "seeds" — scan sequences in database for matching words



# BLAST database seed finding

2. Find “seeds” — scan sequences in database for matching words



*Q* ACTGA

*D* GACTGC

Word List

Seeds?

ACT  
CTG  
TGA

What are they?

# BLAST database seed finding

2. Find “seeds” — scan sequences in database for matching words



*Q* ACTGA

*D* GACTGC

Word List

ACT  
CTG  
TGA

Seeds?

ACT  
CTG

# BLAST database seed finding

2. Find “seeds” — scan sequences in database for matching words



*Q* ACTGA

*D* GACTGC

Word List

Seeds?

ACT

ACT

CTG

CTG

TGA

What kind of data structures might you use to store the words?  
The word/seed associations?

# BLAST seed extension

3. Extend — use (nearby) seeds to form local alignments called HSPs



- Extend seeds in both directions.
- No gaps allowed
- Stop when alignment score drops below threshold  $T$
- Extended seeds are called HSPs (high scoring pairs)

# BLAST seed extension

3. Extend — use (nearby) seeds to form local alignments called HSPs



Q ACTGA

Word List

ACT  
CTG  
TGA

D GACTGC

Seeds?

ACT  
CTG

HSP  
(high-scoring pair)

GACTGC  
ACT  
CTG

→

GACTGC  
ACTG

# BLAST scoring

4. Score—combine groups of consistent HSPs into local alignment with best score

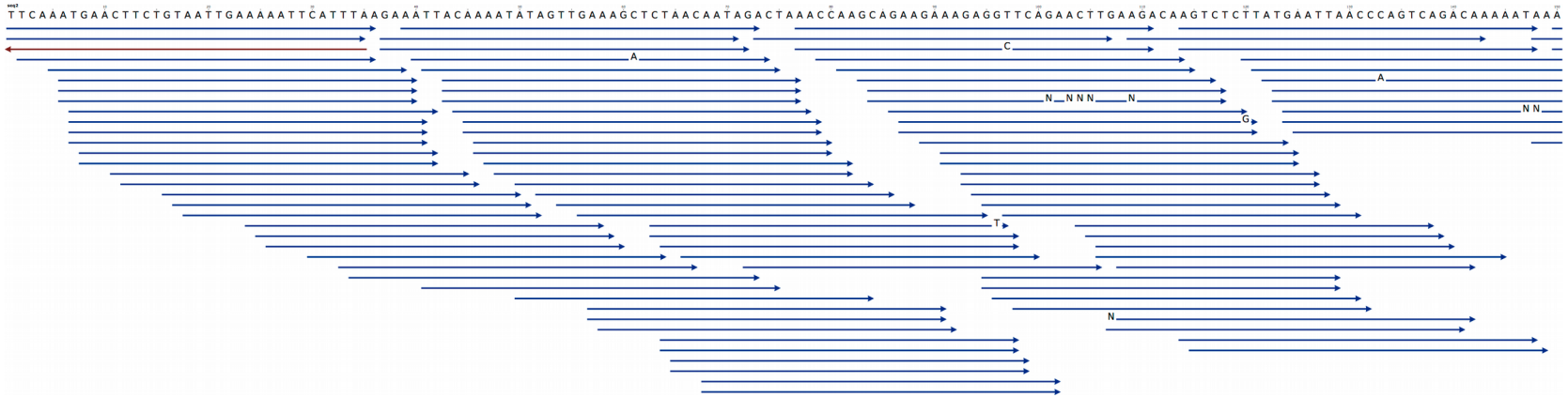


- Combine consistent HSPs
  - They can't overlap
  - They must be in same order in Q and D sequences
  - These “consistent local alignments” are assessed for statistical significance using an analytical null model



# Next-generation sequencing creates a new alignment problem

- Align short reads to a large genome



# Short-read aligners to a reference genome

- Need to be much faster than BLAST
- Reads or reference sequence or both can be indexed
- Types of indexes
  - Hash tables
  - Suffix arrays
  - Suffix trees
  - Burroughs-Wheeler Transform

# Next-generation sequencing alignment algorithms

- If you are interested, check out
  - Ben Langmead Coursera videos:  
<https://www.youtube.com/watch?v=IzXQVwWYFv4&index=9&list=PL2mpR0RYFQsBiCWVJSvVAO3OJ2t7DzoHA>
  - Steven Salzberg lecture posted on Class Schedule website