

Algorithms for Bioinformatics: Basic Sequence Analysis

Jan Michael C. Yap

Core Facility for Bioinformatics

jcyap@up.edu.ph

Outline

- Sequence Alignment
- Sequence Database Searching: BLAST



Outline

- Sequence Alignment
- Sequence Database Searching: BLAST



Basis for Sequence Similarity

begin

A C G T C A T C A

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

A ~~X~~ G T G ~~X~~ T C A

A G T G T C A

T A G T G T C A

end

T A G T G T C A

mutation

deletion

insertion

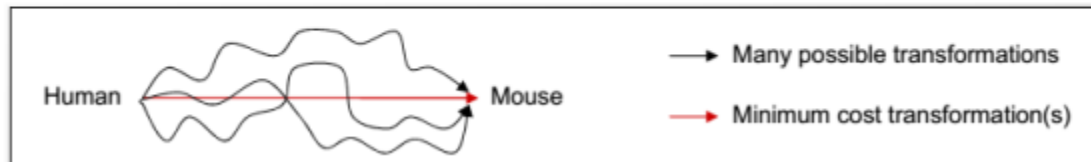
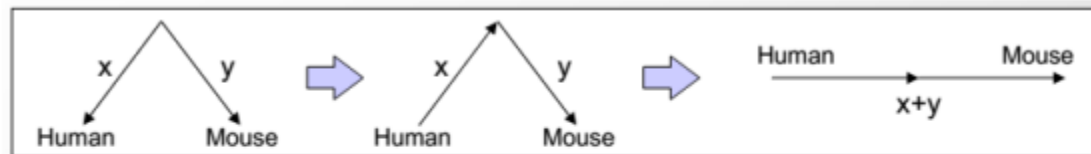
begin

A C G T C A T C A

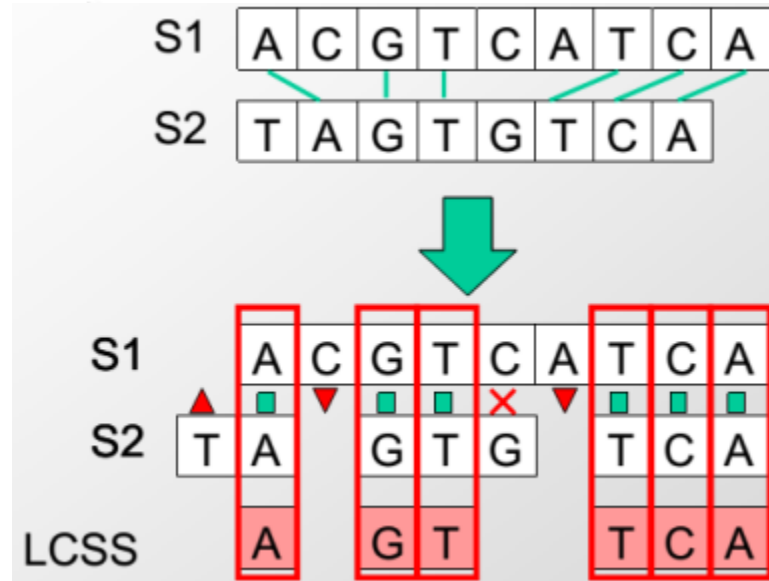
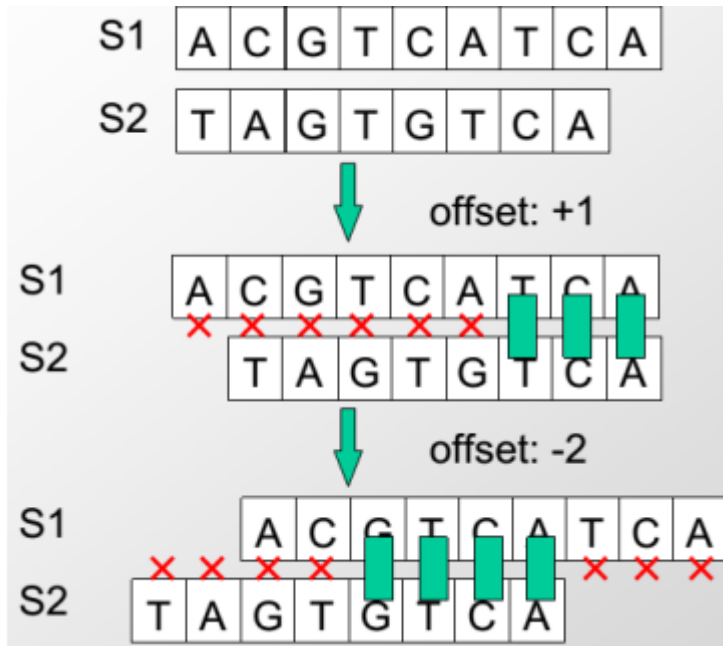
?

end

T A G T G T C A



Viewing the Sequence Similarity Problem Computationally



Scoring function:

Match(x,x) = +1

Mismatch(A,G) = -1/2

Mismatch(C,T) = -1/2

Mismatch(x,y) = -1

	A	G	T	C
A	+1	-1/2	-1	-1
G	-1/2	+1	-1	-1
T	-1	-1	+1	-1/2
C	-1	-1	-1/2	+1

purine pyrimid.

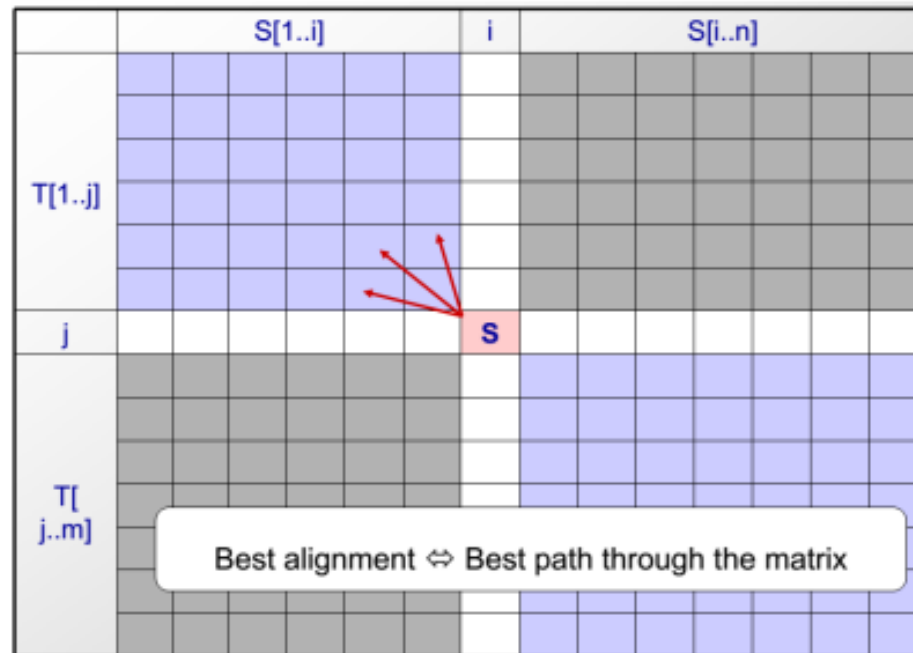
Transitions:

A ↔ G, C ↔ T common
(lower penalty)

Transversions:

All other operations

Setting Things Up: Data Structure and Scoring



- **Local update:** Compute **next alignment** based on **previous alignment** (table look-up)
- Compute **scores for prefixes of increasing length**
- Compute the **score of a cell from its neighbors**

$$F(i, j) = \max(F(i-1, j) - gap, F(i, j) + score, F(i, j-1) - gap)$$

Global Alignment: Needleman-Wunsch Algorithm

		A	T	G
	0			
G				
G				
A				
A				
T				
G				
T				
I/D	-2			
Mismatch	-1			
Match	1			

		A	T	G
	0	-2	-4	-6
G	-2			
G	-4			
A	-6			
A	-8			
T	-10			
G	-12			
T	-14			
I/D	-2			
Mismatch	-1			
Match	1			

Global Alignment: Needleman-Wunsch Algorithm

		A	T	G
	0	-2	-4	-6
G	-2			
G	-4			
A	-6			
A	-8			
T	-10			
G	-12			
T	-14			
I/D	-2			
Mismatch	-1			
Match	1			

		A	T	G
	0	-2	-4	-6
G	-2			
G	-4			
A	-6			
A	-8			
T	-10			
G	-12			
T	-14			
I/D	-2	Direction	Align	Score
Mismatch	-1	D:	A	0 + -1 = -1
Match	1		G	
		U:	A-	-2 + -2 = -4
			-G	
		L:	-A	-2 + -2 = -4
			G-	

Global Alignment: Needleman-Wunsch Algorithm

		A	T	G
	0	-2	-4	-6
G	-2			
G	-4			
A	-6			
A	-8			
T	-10			
G	-12			
T	-14			

I/D	-2	Direction	Align	Score
Mismatch	-1	D:	A	$0 + -1 = -1$
Match	1		G	
		U:	A-	$-2 + -2 = -4$
			-G	
		L:	-A	$-2 + -2 = -4$
			G-	

		A	T	G
	0	-2	-4	-6
G	-2	-1(D)		
G	-4			
A	-6			
A	-8			
T	-10			
G	-12			
T	-14			

I/D	-2	Direction	Align	Score
Mismatch	-1	D:	A	$0 + -1 = -1$
Match	1		G	
		U:	A-	$-2 + -2 = -4$
			-G	
		L:	-A	$-2 + -2 = -4$
			G-	

Global Alignment: Needleman-Wunsch Algorithm

		A	T	G
	0	-2	-4	-6
G	-2			
G	-4			
A	-6			
A	-8			
T	-10			
G	-12			
T	-14			

I/D	-2	Direction	Align	Score
Mismatch	-1	D:	A	$0 + -1 = -1$
Match	1		G	
		U:	A-	$-2 + -2 = -4$
			-G	
		L:	G-	$-2 + -2 = -4$
			-A	

		A	T	G
	0	-2	-4	-6
G	-2	-1 (D)		
G	-4			
A	-6			
A	-8			
T	-10			
G	-12			
T	-14			

I/D	-2	Direction	Align	Score
Mismatch	-1	D:	A	$0 + -1 = -1$
Match	1		G	
		U:	A-	$-2 + -2 = -4$
			-G	
		L:	G-	$-2 + -2 = -4$
			-A	

Global Alignment: Needleman-Wunsch Algorithm

		A	T	G
	0	-2	-4	-6
G	-2	-1 (D)		
G	-4			
A	-6			
A	-8			
T	-10			
G	-12			
T	-14			

I/D	-2	Direction	Align	Score
Mismatch	-1	D:	A	0 + -1 = -1
Match	1		G	
		U:	A-	-2 + -2 = -4
			-G	
		L:	G-	-2 + -2 = -4
			-A	

		A	T	G
	0	-2	-4	-6
G	-2	-1(D)	-3(D/L)	
G	-4			
A	-6			
A	-8			
T	-10			
G	-12			
T	-14			

I/D	-2	Direction	Align	Score
Mismatch	-1	D:	AT	-2 + -1 = -3
Match	1		-G	
		U:	AT-	-4 + -2 = -6
			--G	
		L:	AT	-1 + -2 = -3
			G-	

Global Alignment: Needleman-Wunsch Algorithm

		A	T	G
	0	-2	-4	-6
G	-2	-1(D)	-3(D/L)	
G	-4			
A	-6			
A	-8			
T	-10			
G	-12			
T	-14			

I/D	-2	Direction	Align	Score
Mismatch	-1	D:	AT	$-2 + -1 = -3$
Match	1		-G	
		U:	AT-	$-4 + -2 = -6$
			--G	
		L:	AT	$-1 + -2 = -3$
			G-	

		A	T	G
	0	-2	-4	-6
G	-2	-1(D)	-3(D/L)	-3(D)
G	-4			
A	-6			
A	-8			
T	-10			
G	-12			
T	-14			

I/D	-2	Direction	Align	Score
Mismatch	-1	D:	ATG	$-4 + 1 = -3$
Match	1		--G	
		U:	ATG-	$-6 + -2 = -8$
			---G	
		L:	ATG	$-3 + -2 = -5$
			-G-	

Global Alignment: Needleman-Wunsch Algorithm

		A	T	G
	0	-2	-4	-6
G	-2	-1(D)	-3(D/L)	-3(D)
G	-4	-3(D/U)		
A	-6			
A	-8			
T	-10			
G	-12			
T	-14			

		A	T	G
	0	-2	-4	-6
G	-2	-1(D)	-3(D/L)	-3(D)
G	-4	-3(D/U)	-2(D)	-2(D)
A	-6	-3(D)	-4(D/U)	-3(D)
A	-8	-5(D/U)	-4(D)	-5(D/U)
T	-10	-7(U)	-4(D)	-5(D)
G	-12	-9(U)	-6(U)	-3(D)
T	-14	-11(U)	-8(U)	-5(U)

I/D	-2	Direction	Align	Score
Mismatch	-1	D:	-A	$-2 + -1 = -3$
Match	1		GG	
		U:	A-	$-1 + -2 = -3$
			GG	
		L:	--A	$-4 + -2 = -6$
			GG-	

Global Alignment: Needleman-Wunsch Algorithm

		A	T	G
	0	-2	-4	-6
G	-2	-1(D)	-3(D/L)	-3(D)
G	-4	-3(D/U)	-2(D)	-2(D)
A	-6	-3(D)	-4(D/U)	-3(D)
A	-8	-5(D/U)	-4(D)	-5(D/U)
T	-10	-7(U)	-4(D)	-5(D)
G	-12	-9(U)	-6(U)	-3(D)
T	-14	-11(U)	-8(U)	-5(U)

		A	T	G
	0	-2	-4	-6
G	-2	-1(D)	-3(D/L)	-3(D)
G	-4	-3(D/U)	-2(D)	-2(D)
A	-6	-3(D)	-4(D/U)	-3(D)
A	-8	-5(D/U)	-4(D)	-5(D/U)
T	-10	-7(U)	-4(D)	-5(D)
G	-12	-9(U)	-6(U)	-3(D)
T	-14	-11(U)	-8(U)	-5(U)

BLUE:	G	G	A	A	T	G	T
	-	-	-	A	T	G	-
RED:	G	G	A	A	T	G	T
	-	-	A	-	T	G	-

Semi-Global Alignment: Modified Needleman-Wunsch Algorithm

		A	T	G
	0			
G				
G				
A				
A				
T				
G				
T				
I/D	-2			
Mismatch	-1			
Match	1			

		A	T	G
	0	0	0	0
G	0			
G	0			
A	0			
A	0			
T	0			
G	0			
T	0			
I/D	-2			
Mismatch	-1			
Match	1			

Semi-Global Alignment: Modified Needleman-Wunsch Algorithm

		A	T	G
	0	0	0	0
G	0			
G	0			
A	0			
A	0			
T	0			
G	0			
T	0			
I/D	-2			
Mismatch	-1			
Match	1			

		A	T	G
	0	0	0	0
G	0	-1(D)	-1(D)	-1(D)
G	0	-1(D)	-2(D)	-2(D)
A	0	1(D)	-2(D/L)	-3(D)
A	0	1(D)	0(D)	-2(L)
T	0	-1(D)	2(D)	0(L)
G	0	-1(D)	0(U)	3(D)
T	0	-1(D)	0(D)	1(U)
I/D	-2			
Mismatch	-1			
Match	1			

BLUE:

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

Local Alignment: Smith-Waterman Algorithm

		A	T	G
	0	0	0	0
G	0			
G	0			
A	0			
A	0			
T	0			
G	0			
T	0			
I/D	-2			
Mismatch	-1			
Match	1			

		A	T	G
	0	0	0	0
G	0			
G	0			
A	0			
A	0			
T	0			
G	0			
T	0			
I/D	-2	Direction	Align	Score
Mismatch	-1	D:	A	0 + -1 = -1
Match	1		G	
		U:	A-	-2 + -2 = -4
			-G	
		L:	-A	-2 + -2 = -4
			G-	
		IF D, U, and L < 0, SET TO 0		

Local Alignment: Smith-Waterman Algorithm

		A	T	G
	0	0	0	0
G	0	0		
G	0			
A	0			
A	0			
T	0			
G	0			
T	0			

I/D	-2	Direction	Align	Score
Mismatch	-1	D:	A	0 + -1 = -1
Match	1		G	
		U:	A-	-2 + -2 = -4
			-G	
		L:	-A	-2 + -2 = -4
			G-	
IF D, U, and L < 0, SET TO 0				

		A	T	G
	0	0	0	0
G	0	0	0	0
G	0	0	0	0
A	0	1		
A	0			
T	0			
G	0			
T	0			

I/D	-2	Direction	Align	Score
Mismatch	-1	D:	--A	0 + 1 = 1
Match	1		GGA	
		U:	-A-	0 + -2 = -2
			GGA	
		L:	---A	0 + -2 = -2
			GGA-	
		IF D, U, and L < 0, SET TO 0		

Local Alignment: Smith-Waterman Algorithm

		A	T	G
	0	0	0	0
G	0	0	0	0
G	0	0	0	0
A	0	1	0	0
A	0	1		
T	0			
G	0			
T	0			

I/D	-2	Direction	Align	Score
Mismatch	-1	D:	---A	0 + 1 = 1
Match	1		GGAA	
		U:	-A--	1 + -2 = -1
			GGAA	
		L:	----A	0 + -2 = -2
			GGAA-	
		IF D, U, and L < 0, SET TO 0		

		A	T	G
	0	0	0	0
G	0	0	0	0
G	0	0	0	0
A	0	1	0	0
A	0	1	0	
T	0			
G	0			
T	0			

I/D	-2	Direction	Align	Score
Mismatch	-1	D:	---AT	1 + -1 = 0
Match	1		GGAAA	
		U:	-AT-	0 + -2 = -2
			GGAA	
		L:	---AT	0 + -2 = -2
			GGAA-	
		IF D, U, and L < 0, SET TO 0		

Local Alignment: Smith-Waterman Algorithm

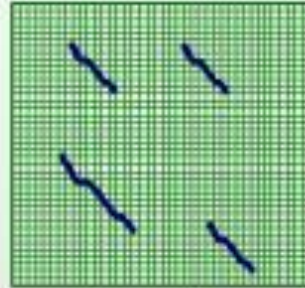
		A	T	G
	0	0	0	0
G	0	0	0	0
G	0	0	0	0
A	0	1(D)	0	0
A	0	1(D)	0(D)	0
T	0	0	2(D)	0
G	0	0	0	3(D)
T	0	0	1(D)	1(U)

A, ATG, T

Pairwise Alignment Algorithms: Summary



Global



Local



Semi-global

Initialization

Top left

Top row/left col.

Top row

Iteration: max

$$F(i-1, j) - d$$

$$F(i, j-1) - d$$

$$F(i-1, j-1) + s(x_i, y_j)$$

$$0$$

$$F(i-1, j) - d$$

$$F(i, j-1) - d$$

$$F(i-1, j-1) + s(x_i, y_j)$$

$$F(i-1, j) - d$$

$$F(i, j-1) - d$$

$$F(i-1, j-1) + s(x_i, y_j)$$

Termination

Bottom right

Anywhere

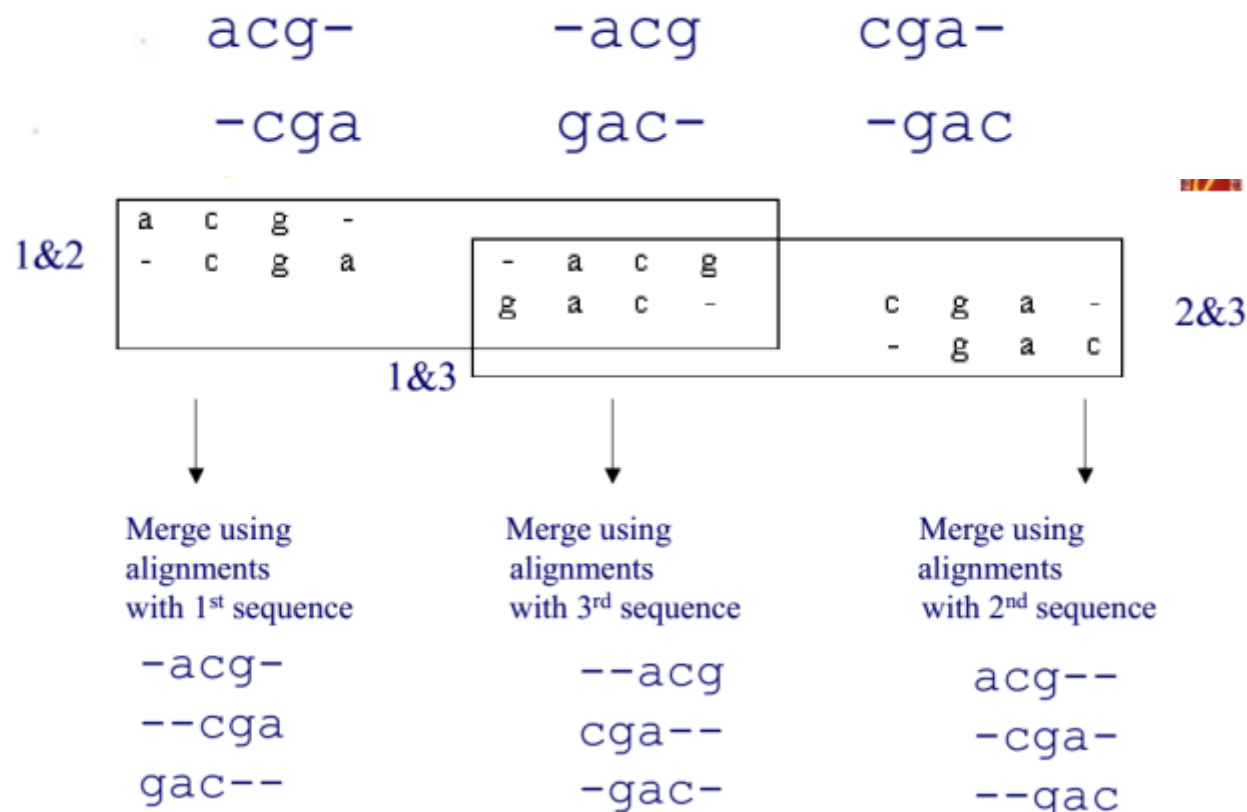
Right column

Multiple Sequence Alignment: Progressive Alignment

```
--T--CC-C-AGT--TATGT-CAGGGGACACG--A-GCATGCAGA-GAC
  |  || |  ||  | | | |||      || |  | |  | |||  |
AATTGCCGCC-GTCGT-T-TTCAG----CA-GTTATG--T-CAGAT--C
  ||||| |   X|||| |           || XXX|||  | ||| |  |
-ATTGC-G--ATTCGTAT-----GGGACA-TGGATGCATGCAG-TGAC
```

Multiple Sequence Alignment: Progressive Alignment

- Align the following: acg, cga, gac



Outline

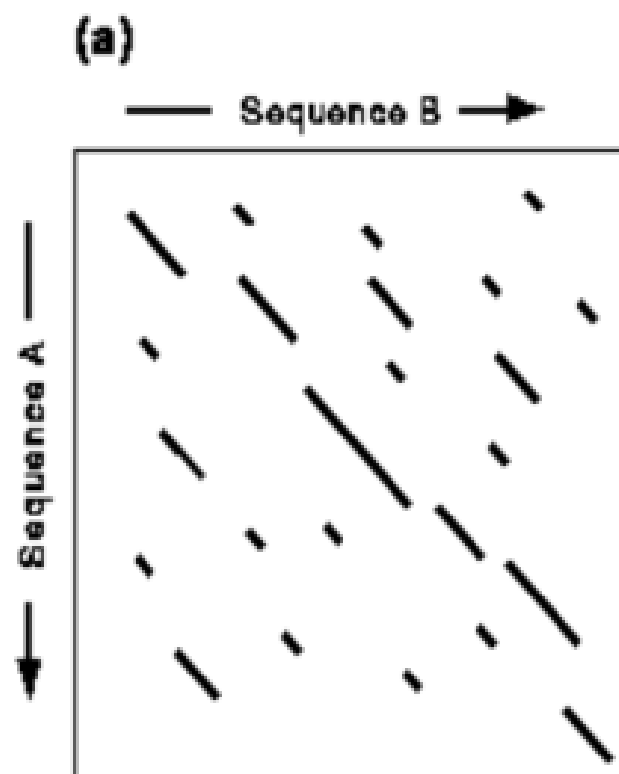
- Sequence Alignment
- Sequence Database Searching: BLAST



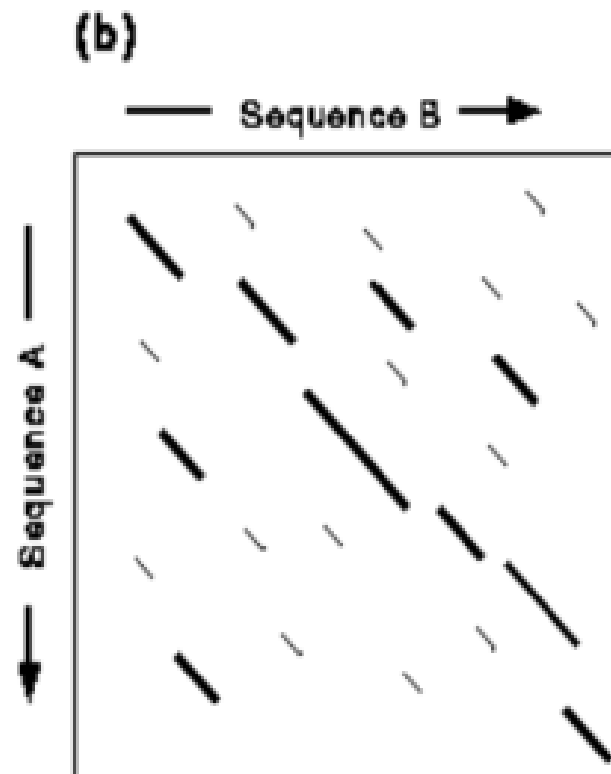
BLAST

- BLAST is a refinement of **FASTA**
- Employs **heuristics** for local alignment
 - Discard irrelevant sequences and perform exact local alignment on the remaining sequences
- Designed **specifically for database searches**
- BLAST produces several short segments called **High Scoring Segment Pairs (HSPs)**

FASTA

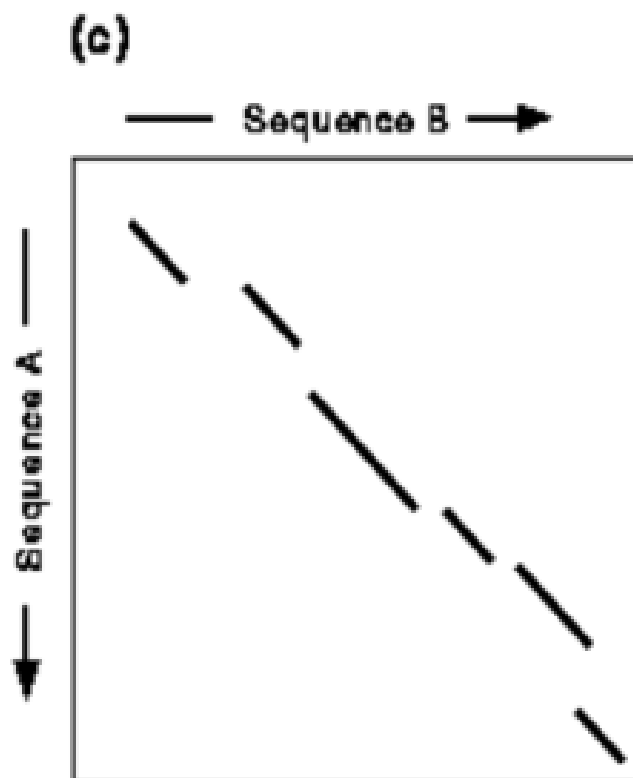


Find runs of identical words

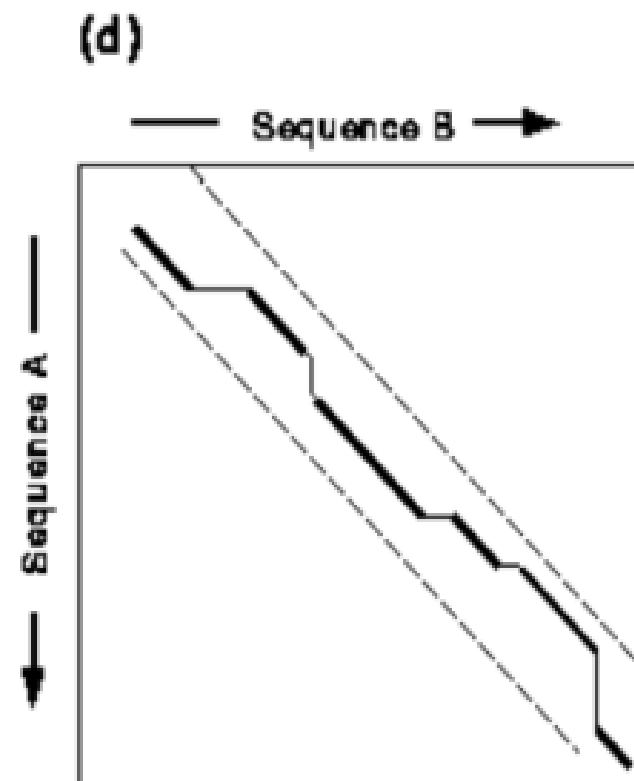


Re-score using PAM matrix
Keep top scoring segments

FASTA



Join segments using gaps,
eliminate other segments



Use dynamic programming to
create an optimal alignment

BLAST Preprocessing

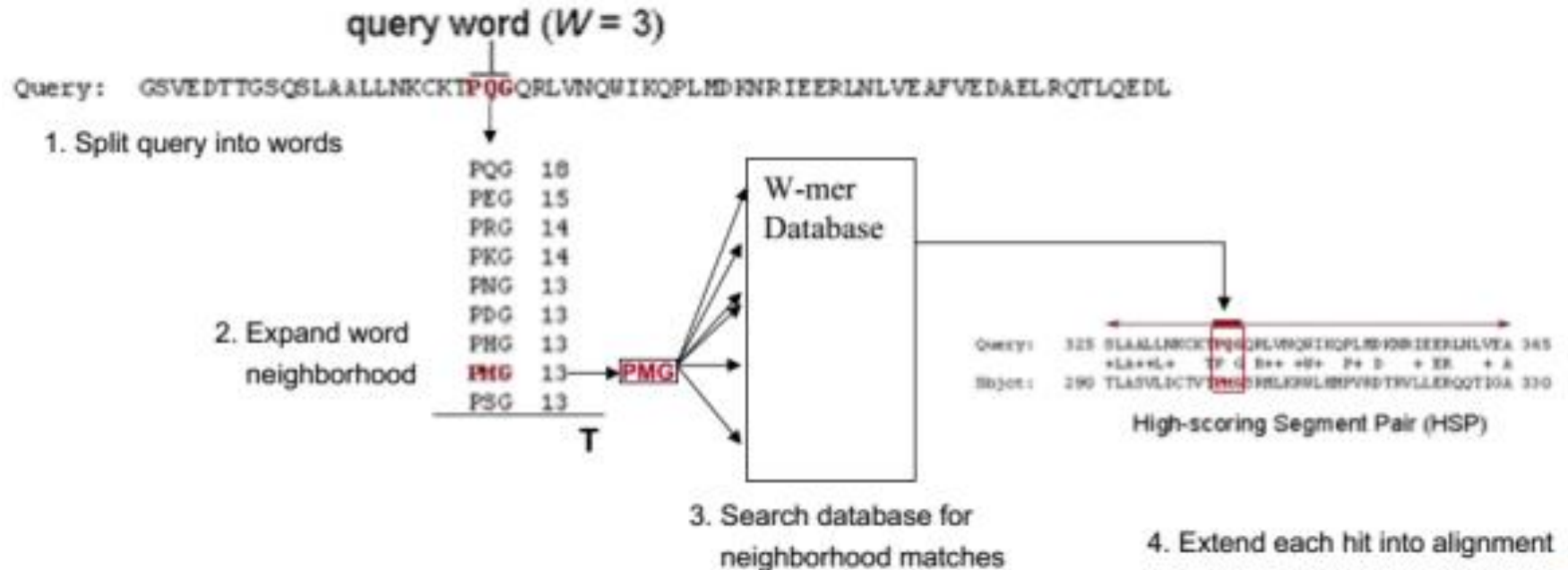
MEAAVKEEISVEDEAVDKNI

MEA
EAA
AAV
AVK
VKE
KEE
EEI
EIS
ISV
...

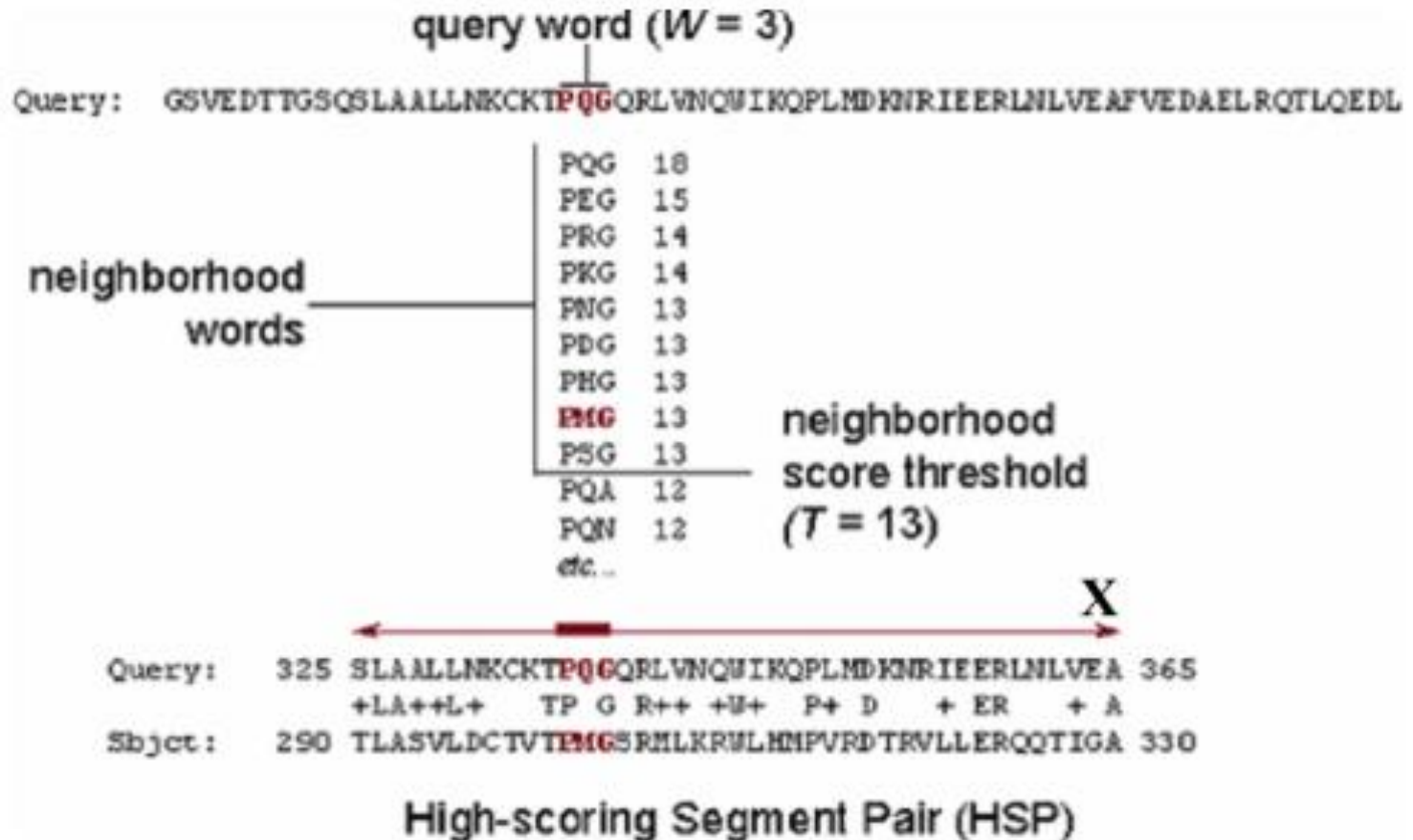
Break query
into words:

Break database
sequences
into words:

Sequence Database Searching: BLAST Algorithm



Sequence Database Searching: BLAST Heuristic



BLAST Alignment Extension

ASKIOPLLWLAASFLHNEQAPALSDAN

JWQEOPLWPLAASOIHLFACNSIFYAS

Score=15



Score=17



Score=14



- Stop extending if the score of the current alignment drops at least X points below the maximum score (obtained so far)
 - X is called the **alignment extension decrement threshold**
- Retain HSPs by discarding segments with score less than the **segment score threshold**

BLAST Scoring

- E-value¹

$$E - value = KMNe^{-\lambda S}$$

- K and λ are “normalization” parameters, as determined by the scoring matrix and gap penalties used
- M is the length of the query sequence
- N is the total length of the sequences stored in the database

- Bit score²

- Bit score S'

$$S' = \frac{\lambda S - \ln K}{\ln 2}$$

- Formula for E-value using S'

$$E - value = MN2^{-S'}$$

1 http://www.fing.edu.uy/inco/grupos/bioinf/bioinfo1/material/Blast_bioinfo1_set09.pdf

2 <https://www.ncbi.nlm.nih.gov/BLAST/tutorial/Altschul-1.html#head3>

THANK YOU VERY MUCH! 😊