

BIOS 477/877 Bioinformatics and Molecular Evolution

Instructor: Etsuko Moriyama
(School of Biological Sciences)

Spring 2026 Lecture 11

BIOS477/877 L11 - 1

1

Today's topics

- Statistical Significance of Alignment Scores
- Similarity Search
 - FASTA and BLAST

BIOS477/877 L11 - 2

2

Significance of alignment scores

- $P(S \geq x | H_0)$: probability of getting the alignment score $S \geq x$ by chance

Karlin-Altschul equation (Karlin and Altschul 1990)

$$P(S \geq x) = 1 - \exp[-Kmn e^{-\lambda x}] \approx Kmn e^{-\lambda x}$$

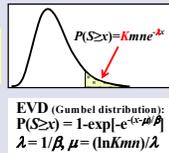
- K and λ : parameters based on a given scoring matrix and the amino acid composition of the sequences
- m and n : lengths of sequences aligned

- Solved for ungapped local alignments
- Can be applied for gapped local alignments

➤ **E-value** = $P(S \geq x | H_0) \times N$ **E-value \neq P-value**

N : the number of sequences in the dataset

Expected number of randomly selected sequences in the dataset that have alignment scores $\geq x$



BIOS477/877 L11 - 3

3

Estimation of K , λ , and P-value

- Estimation of K and λ from empirical distribution of random alignment scores (used in LALIGN and PRSS)

- 1) The second sequence is shuffled many times
→ Simulates random sequences
- 2) Smith-Waterman local alignment score is calculated from each alignment: $P(S \geq x | H_0)$
- 3) The distribution is fitted to an extreme value distribution to obtain estimates of K and λ
- 4) P-value is estimated based on the K and λ , and the original alignment score x : $P(S \geq x) \approx Kmn e^{-\lambda x}$

BIOS477/877 L11 - 4

4

Simulation of random score distribution (example)

[Input sequences]

- RECA_ECOLI (P0A7G6; 353 amino acids)
- RAD51_YEAST (P25454; 400 amino acids)

Smith-Waterman local alignment score = 293

(BLOSUM50, gap opening: -10, gap extension: -1)

RECA_ECOLI	3	IQNNE-QEALAAALQIENQFQGGSHMLGDSRDMQVETIQGLSILIA	51
RAD51_YEAST	124	IQNNE-QEALAAALQIENQFQGGSHMLGDSRDMQVETIQGLSILIA	170
RECA_ECOLI	52	LGAGLQPHGIVLQYVPSGQTTI-----TLQV-IAAAGQGGTCAPF	94
RAD51_YEAST	171	LG-GGVYGTITELGFRFRKQGLQNTLAVTQIPLDGGGK-CLYI	218
RECA_ECOLI	95	DAHALDPI---VARELQVDFD---NLLSQPQFQDALEICDALAR-	135
RAD51_YEAST	219	DTEQFRFVRLVSIAGRFQDLDALNVAATATADNADLALDAAGM	268
RECA_ECOLI	136	--SGAVVIVDVSVAALTFKAEI--EGEIGDSBGLAAMMSQMRKLAG	181
RAD51_YEAST	269	MSEKRFSLIVDSVVALY-RDFQSGEISANQMLAKFM--RILQRLA-	314
RECA_ECOLI	182	NLQSNYLLPFIQIRKMI--GVNFG-NPFTTGGNALEFYASVLDLRR	228
RAD51_YEAST	315	--GQVAVVYVYVAVYVGGGNAFQDFEETIGQNIHNSSETTL---	358
RECA_ECOLI	229	IGAVKEGVVYVYVYVYVYVYVYVYVYVYVYVYVYVYVYVYVYVYV	269
RAD51_YEAST	359	--GPEKGE--GQRLCKVY-DEPKLFEAKYVAL-YEDVY	392

Using [EMBOSS WATER @ EBI](#)

BIOS477/877 L11 - 5

5

Simulation of random score distribution (example)

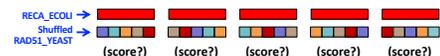
[Input sequences]

- RECA_ECOLI (P0A7G6; 353 amino acids)
- RAD51_YEAST (P25454; 400 amino acids)

Smith-Waterman local alignment score = 293

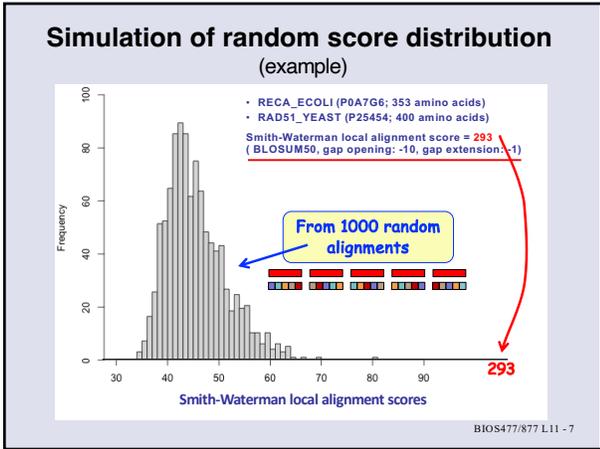
(BLOSUM50, gap opening: -10, gap extension: -1)

1. Shuffle RAD51_YEAST 1000 times
(generate 1000 random sequences)
2. Align with RECA_ECOLI
(generate 1000 random local alignments)

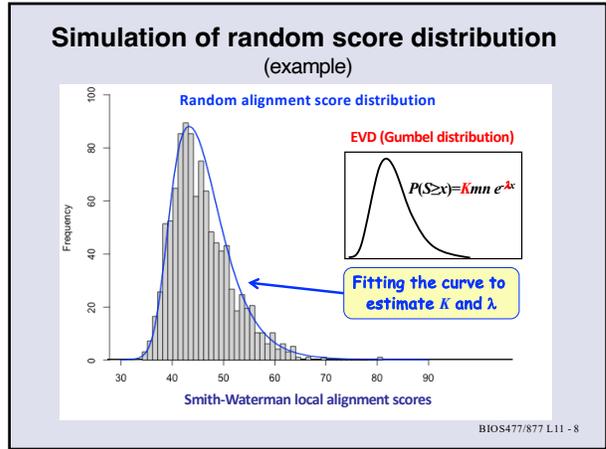


BIOS477/877 L11 - 6

6



7



8

FASTA Web server by William Pearson

FASTA Sequence Comparison at the U. of Virginia

UvA FASTA Server

William R. Pearson

- Original FASTA package was released on 1988 (earlier than BLAST)
- The origin of FASTA format

BIOS477/877 L11 - 9

9

Statistical analysis of alignments by PRSS

Statistical Significance from Shuffles

Start shuffling

The second sequence is shuffled 1000 times to generate random sequences. → Amino acid composition can be maintained. → Window size can be set so that shuffling is done locally

BIOS477/877 L11 - 10

10

Statistical analysis of alignments by PRSS (result example)

```

# search for P: 4 -w 80 -m 8 -2 10000 -E 10000 -f 10 -g 2 @ TMPq2
SEARCH performed a Smith-Waterman search
P: sp|P0A7G6|RECA_ECOLI_Protein RecA; Recombinase A - 353 aa
P: sp|P25454|RAD51_YEAST_DNA repair protein RAD51 (400 aa)
...

```

$\lambda = 0.1376$, $K = 0.01032$
 $x = 236$, $m = 353$, $n = 400$
 $E(10000) = 10000 \times P = 1.1 \times 10^{-11} \times 10000 = 1.1 \times 10^{-7}$

BIOS477/877 L11 - 11

11

Statistical analysis of alignments by PRSS (result example)

```

Query: sp|P0A7G6|RECA_ECOLI_Protein RecA; Recombinase A - 353 aa
Library: TMP_42
400 residues in 1 sequences
Statistics: (shuffled 1000) MLE statistics: Lambda= 0.1376; K=0.01032
Algorithm: Smith-Waterman (SSE2, Michael Farrar 2006) (7/7 Nov 2010)
Parameters: BLOSUM matrix (35-35), open/ext: -10/-2
Scan time: 0.080

```

$\lambda = 0.1376$, $K = 0.01032$
 $x = 236$, $m = 353$, $n = 400$
 $E(10000) = 10000 \times P = 1.1 \times 10^{-11} \times 10000 = 1.1 \times 10^{-7}$

$P(S \geq x) \approx Kmne^{-\lambda x}$
 λ and K are estimated based on simulated random alignments → Each shuffling can generate slightly different values!

BIOS477/877 L11 - 12

12

Similarity search as multiple pairwise alignments

Similarity search

Multiple pairwise alignments

Pairwise alignment search space

X multiple comparisons

BIOS477/877 L11 - 19

19

Similarity search using Smith-Waterman algorithm (SSEARCH)

Searches the entire alignment space (X multiple comparisons)

Search time needs to be reduced!

For a large database, it requires lots of time

BIOS477/877 L11 - 20

20

Heuristic algorithm (FASTA/BLAST)

Minimizes the search space

Faster search speed
But search is not thorough

BIOS477/877 L11 - 21

21

Search speed and sensitivity

Slower → Faster

SSEARCH (Smith-Waterman dynamic programming)

Guarantees to find the optimal alignments

BLAST (Basic Local Alignment Search Tool; Altschul *et al.* 1990, 1997)

FASTA (Pearson and Lipman 1988)

More sensitive ← Less sensitive

More sensitive = detects more distantly related sequences

Heuristic algorithm (fast approximation)

BIOS477/877 L11 - 22

22

Word/k-tuple matching

To reduce the search space:

Break each sequence into words (or k-tuples/k-mers)
e.g., word length (k-tuple) = 2

Finding short identical words can be done very rapidly

BIOS477/877 L11 - 23

23

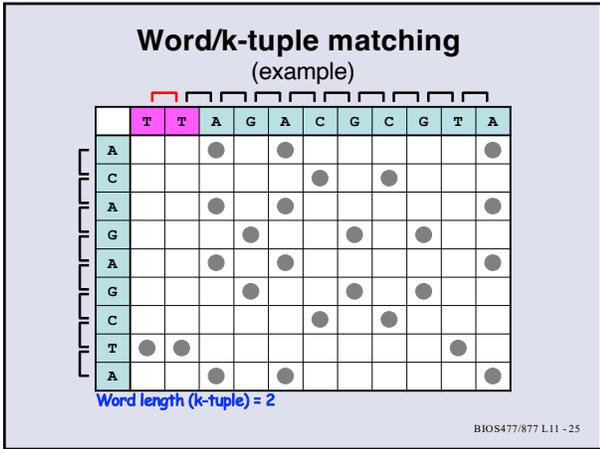
Word/k-tuple matching

Break each sequence into words
e.g., word length (k-tuple) = 2

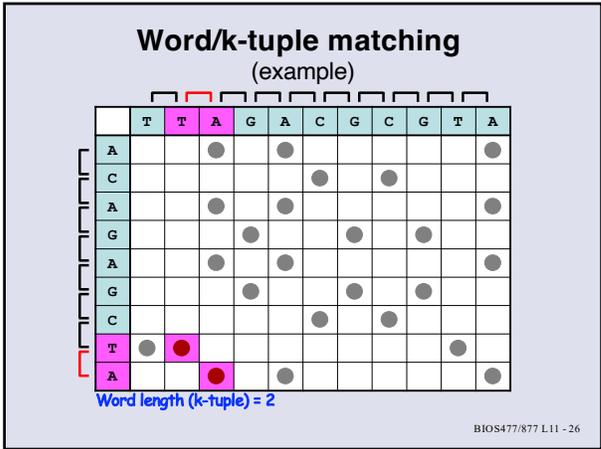
Larger the word size, more rapid, but less thorough (less sensitive)

BIOS477/877 L11 - 24

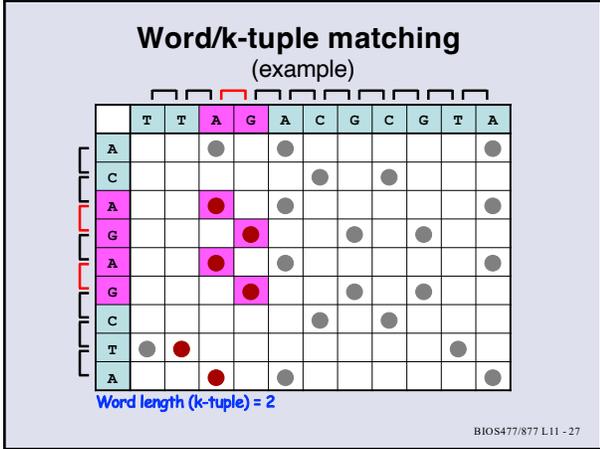
24



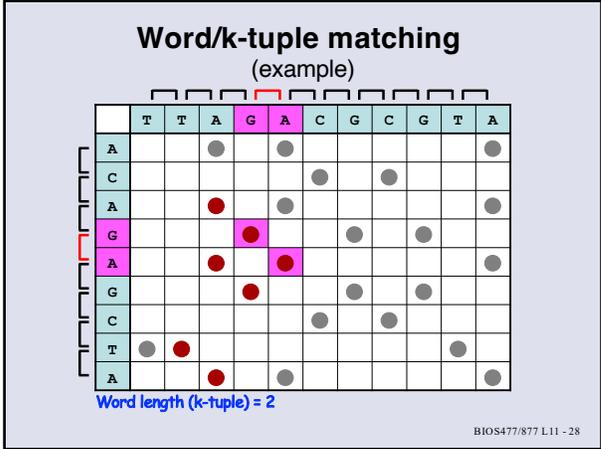
25



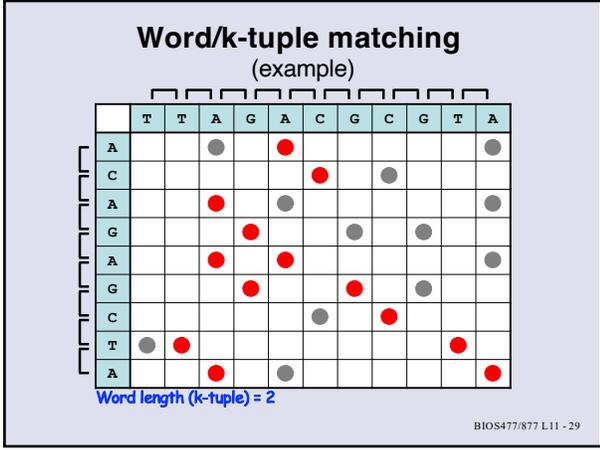
26



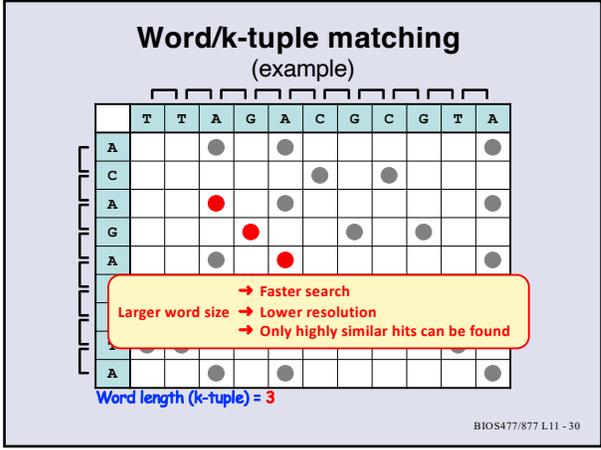
27



28



29

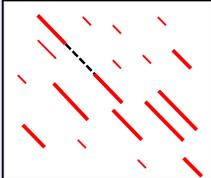


30

FASTA algorithm

1. Find identities using k-tuples
(k-tuples = words)

2. Join diagonals without gaps
3. Choose top 10 diagonals using a scoring matrix (e.g., BLOSUM62)
(init1: the top diagonal score)

BIOS477/877 L11 - 31

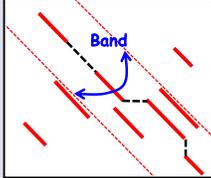
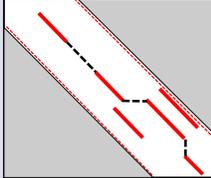
31

FASTA algorithm

4. Join again with gaps
(initn: score of the longer diagonal)

5. A diagonal band is defined
(width: 32 if k=1, 16 if k=2 for protein)

6. Find optimal local alignment using dynamic programming algorithm within the band
(opt: the final score)

BIOS477/877 L11 - 32

32

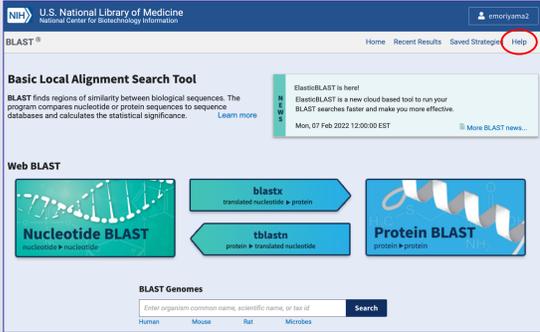
FASTA similarity search

- **Ranking**
 - Database sequences are ranked based on z-values or OPT scores
 - z-value = the number of standard deviations from the mean (standardized score)
- **Alignment**
 - The high scored sequence pairs are aligned using the full Smith-Waterman dynamic programming algorithm
 - Not just from the diagonal band → better alignment
- **FASTA/SSEARCH website (more on CANVAS)**
[Original FASTA server @ U Virginia \(By W. Pearson\)](#)
[FASTA website @ EBI](#)

BIOS477/877 L11 - 33

33

BLAST similarity search



BIOS477/877 L11 - 34

34

BLAST resources

- **BLAST Help**
<https://blast.ncbi.nlm.nih.gov/doc/blast-help/>
 - Guide to BLAST home and search pages
 - BLAST Command Line User Manual
 - BLAST Tutorials on YouTube
- **BLAST Statistics**
<https://www.ncbi.nlm.nih.gov/BLAST/tutorial/Altschul-1.html>
- **"BLAST" (online version)**
(see more on CANVAS)



BIOS477/877 L11 - 35

35

BLAST algorithm: seeding

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

Word hits

Word length = 2

BIOS477/877 L11 - 36

36

BLAST algorithm: seeding (including neighborhood)

- Using **words** reduces the search space
- Neighborhood** increases the sensitivity

Match = 2
Mismatch (Ts) = -1
Mismatch (Tv) = -5

TC
TC = 2 + 2 = 4

TC
TC = 2 - 1 = 1

TC
TC = -1 + 2 = 1

...

TC
TC = 2 - 5 = -3

TC
TC = 2 - 5 = -3

TC
TC = -5 + 2 = -3

TC
TC = -5 - 1 = -6

Neighborhood
Neighborhood Threshold (T) = 1
[minimum score allowed to be the neighborhood]

BIOS477/877 L11 - 37

37

BLAST algorithm: seeding (including neighborhood)

BIOS477/877 L11 - 38

38

BLAST algorithm: seeding (including neighborhood)

BIOS477/877 L11 - 39

39

BLAST algorithm

BIOS477/877 L11 - 40

40

BLAST algorithm: extension

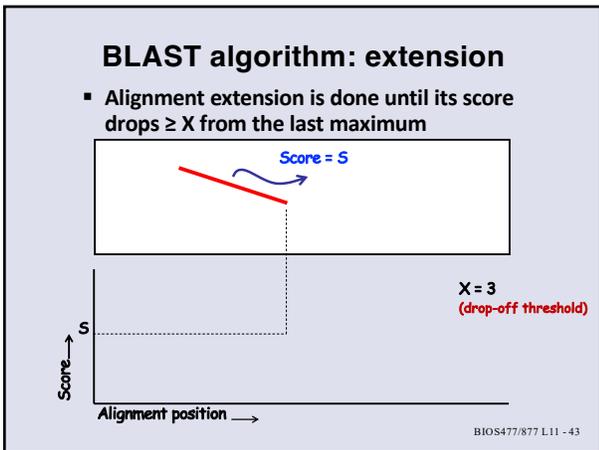
BIOS477/877 L11 - 41

41

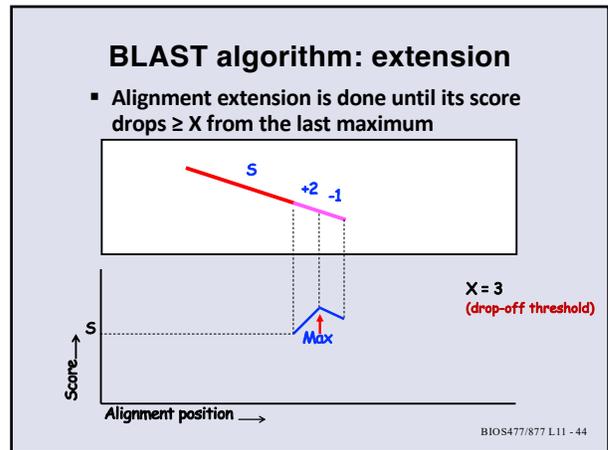
BLAST algorithm: extension

BIOS477/877 L11 - 42

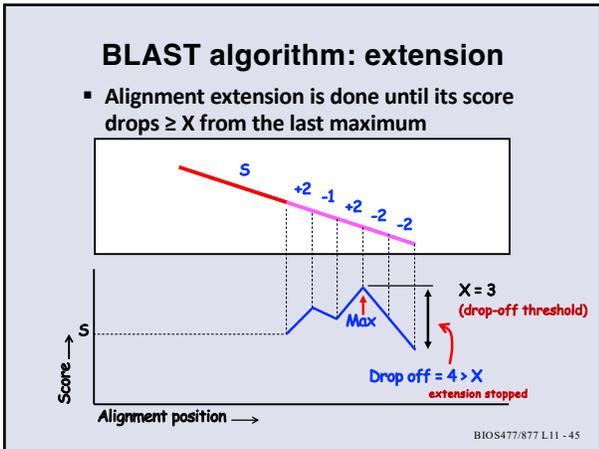
42



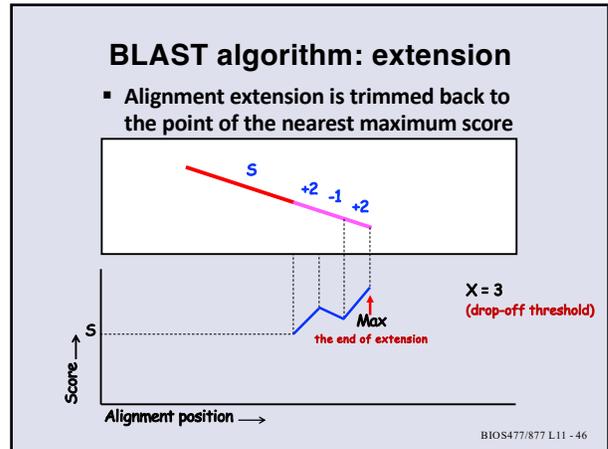
43



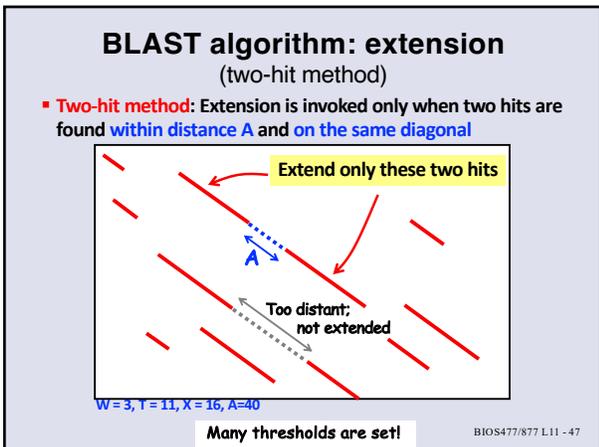
44



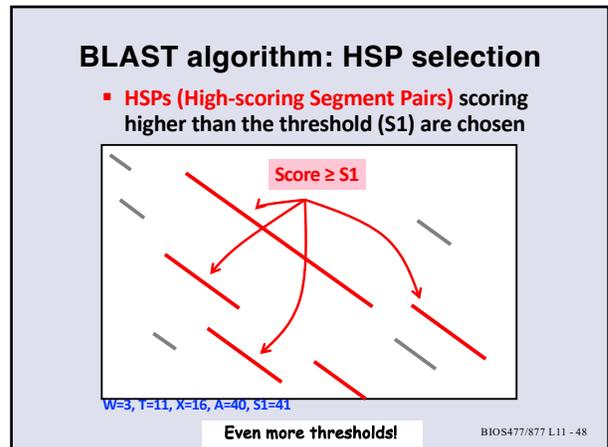
45



46



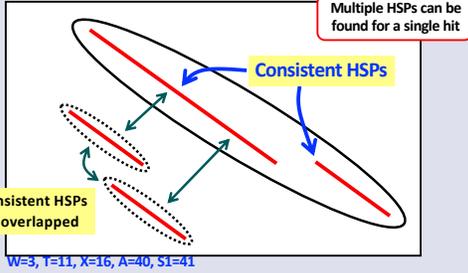
47



48

BLAST algorithm: HSP selection

- Consistent (non-overlapping) HSPs are grouped for E-value calculation

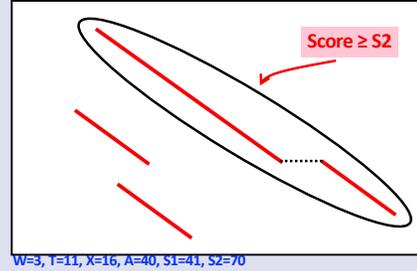


BIOS477/877 L11 - 49

49

BLAST algorithm: gapped extension

- Gapped extension is done only for selected HSPs
- Gapped alignments scoring higher than the threshold ($S2$) is chosen as the final set

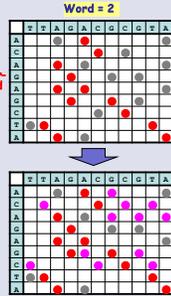
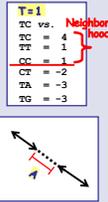


BIOS477/877 L11 - 50

50

BLAST algorithm: summary

- Word-matching size (W)
 - Longer words: faster but less sensitive
- Neighborhood threshold (T)
 - Lower T :
 - detects weaker similarities
 - slower but more sensitive
- Extension
 - Drop-off score threshold (X)
 - Two-hit method (A : distance b/w 2 hits)
- HSP selection (ungapped alignment threshold: $S1$)
- Gapped HSP extension (gapped alignment threshold: $S2$)



BIOS477/877 L11 - 51

51