

BLAST BIFX-550

S. Ravichandran, Ph.D.

Biology,

Hood College, Frederick, MD 21701

Goals

- Explore BLAST from NCBI
- How to carry out BLAST searches?
- Explain the BLAST search in detail.
- BLAST specific inputs/outputs
 - E-values & scores
- Discuss strategies for carrying out BLAST

Discussion ?s about Sequence Alignments

- Pairwise-Sequence Alignment?

```
GACCATA
||| |
GACTA--
```

If we see this first, we believe this is the best alignment

```
GACCATA
||| |
GAC--TA
```

What about this?

```
GACCATA
|| | |
GA-C-TA
```

What about this?

Discussion ?s about Sequence Alignments

- Why do we create alignments?
- What rules govern alignment?
 - Score (parameters)
 - matrices
 - Algorithm
 - Global, local, semi-global

Discussion ?s about Sequence Alignments

- Reality about Sequence Alignments
 - When sequences are similar
 - different algorithms produce similar alignments
 - When they are different
 - You need to explore different parameters, matrices
 - Caution
 - You cannot carry over the parameters/matrices/strategies from similar sequences to non-similar sequences

Discussion ?s about Sequence Alignments

- Sequence alignment displays
 - Gap character
 - “-”
 - Match
 - “|” and dot character for mis-match
 - What is the query? What is the convention for representing query sequence?
 - What is the sequence shown at the top of the pairwise alignments?

Discussion ?s about Sequence Alignments

- Universally best alignment?
 - Alignments depend on score

GATTACA

| | | . |

GATCA--

GATTACA

| | | | |

GAT--CA

GATTACA

| | | | |

GA-T-CA

Discussion ?s about Sequence Alignments

- Alignment scores?

Example, Match: +5; Mismatch: -5; Gap: -10;
Gapextension: -0.5

GATTACA

| | | . |

GATCA--

5 . 5

GATTACA

| | | | |

GAT--CA

14 . 5

GATTACA

| | | | |

GA-T-CA

5

Most biological methods produce errors towards the end of the sequences, To account for this most software apply error-correcting step that will not penalize the gaps at the end. So, the scores will be after correction:

16

14.5

5

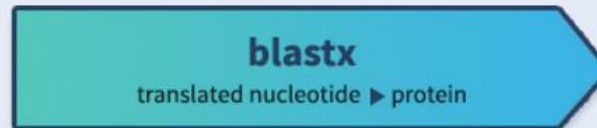
Discussion ?s about Sequence Alignments

- Difference between sequence similarity/identity

Basic Local Alignment Search Tool

- What is BLAST?
 - A query tool to retrieve similar (homologous) sequences from a DB
- Flavors
 - Blastp; blastn; Blastx; tblastn; Tblastx
 - BLAST2

Web BLAST



What is BLAST used for?

- Identifying homologs for proteins/DNA
 - Orthologs & paralog
- May have a sequence and would like to identify the identity
 - SMART BLAST
- Discover new genes
 - Genomic BLAST

Step1

- Sequence of interest
 - Query
 - DNA/Protein
- Input
 - Sequence or Accession number
 - FASTA
 - Etc.

Step 2

What Flavor of BLAST?

**Fig 3.12 from the
Pevsner Book
III Edition PLEASE
DO NOT DISTRIBUTE
Copyright figure**

**UniGene Uses all
nucleotides in its
DB to search
against all known
Protein sequences**

Program	Query	Number of database searches	Database
---------	-------	-----------------------------	----------

BLASTP

protein

1

protein

Use BLASTP to compare a protein query to a database of proteins.

BLASTN

DNA

1

DNA

Use BLASTN to compare both strands of a DNA query against a DNA database.

BLASTX

DNA

6

protein

BLASTX translates a DNA sequence into six protein sequences using all six possible reading frames, and then compares each of these proteins to a protein database.

TBLASTN

protein

6

DNA

TBLASTN is used to translate every DNA sequence in a database into six potential proteins, and then to compare your protein query against each of those translated proteins.

TBLASTX

DNA

36

DNA

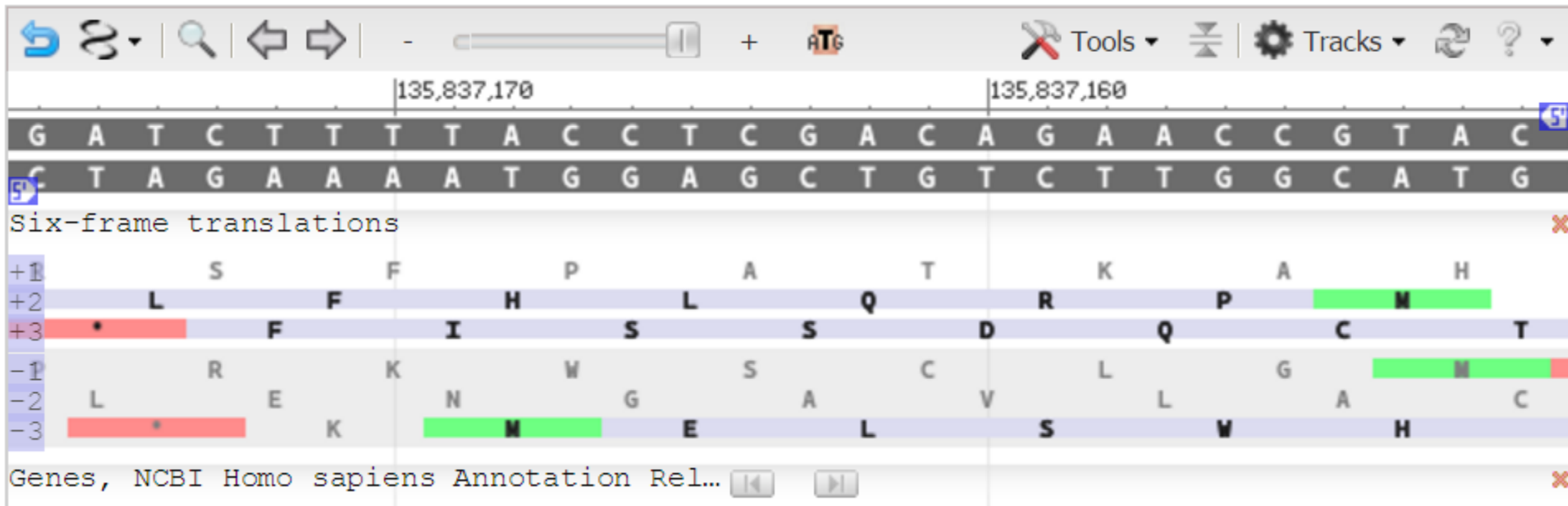
TBLASTX is the most computationally intensive BLAST algorithm. It translates DNA from both a query and a database into six potential proteins, then performs 36 protein-protein database searches.

BLASTX

Translation of a DNA to a protein
and searched against a Protein
DB. What is a translation?

LCT Gene → Protein

6 Frames of Translation



10 20 30 40 50
MELSWHVFI ALLSFSCWGS DWESDRNFIS TAGPLTNDLL HNLSGLLDQ

Step 3: Selecting a DB

- Protein
 - nr database
 - GenBank, PDB, SwissProt, PIR & PRF
 - RefSeq
- DNA
 - BLASTN, TBLASTN, TBLASTX
 - nr/nt (GenBank, EMBL, DDBJ, PDB & RefSeq)
 - Note that nr does not include, EST, STS, WGS, GSS, TSA, Patents or HTGS databases

Table 4.1 from Pevsner III Edition

Database	Title	# sequences	
nr	All nonredundant GenBank CDS translations + PDB + SwissProt + PIR + PRF excluding environmental samples from WGS projects	65 million	~146M
Reference proteins	NCBI protein reference sequences	50 million	~102M
UniProtKB/SwissProt	Nonredundant UniProtKB/SwissProt sequences	450,000	~468K
Patented protein sequences	Protein sequences derived from the Patent division of GenBank	1.3 million	~2.2M
Protein Data Bank	PDB protein database	77,000	~97K
Metagenomic proteins	Proteins from WGS metagenomic projects (env_nr)	6.5 million	~6.9M
Transcriptome	Transcriptome Shotgun Assembly (TSA) sequences	770,000	~2.42M

**Accessed Date
(2018/02/25)**

nr: formed by merging several main protein/DNA DBs

These often contain many identical sequences. Generally only one copy if kept during merging

Database	Title	# sequences
Human Genomic + Transcript	Homo sapiens NCBI Annotation Release 104 RNAs; Homo sapiens all assemblies	55,000
Mouse Genomic + Transcript	Mus musculus NCBI Annotation RNAs; Mus musculus all assemblies	N/A
nr/nt	All GenBank+EMBL+DDBJ+PDB+RefSeq sequences, but excludes EST, STS, GSS, WGS, TSA, patent sequences as well as phase 0, 1, and 2 HTGS sequences	25 million
refseq_rna	NCBI transcript reference sequences	3.5 million
refseq_genomic	NCBI genomic reference sequences	2.7 million
NCBI Genomes	NCBI chromosome sequences	28,000
Expressed sequence tags (EST)	Database of GenBank+EMBL+DDBJ sequences from EST Divisions	75 million
Genomic survey sequences (gss)	Genome survey sequence, includes single-pass genomic data, exon-trapped sequences, and Alu PCR sequences	36 million
High-throughput genomic sequences (HTGS)	Unfinished high-throughput genomic sequences; sequences: phases 0,1 and 2	153,000
Patent sequences	Nucleotide sequences derived from the Patent division of GenBank	21 million
Protein Data Bank	PDB nucleotide database	8000
alu	Human <i>Alu</i> repeat elements	325
Sequence tagged sites (STS)	Database of GenBank+EMBL+DDBJ sequences from STS Divisions	1.3 million
Whole-genome shotgun (wgs)	Whole-genome-shotgun contigs	116 million
Transcriptome Shotgun Assembly (TSA)	Transcriptome shotgun assembly (TSA) sequences	15 million
16S ribosomal RNA sequences (Bacteria and Archaea)	16S ribosomal RNA sequences (bacteria and archaea)	7500

**Table 4.2 from
Pevsner III Edition**

Why is this task difficult?

- Sep 2018
- nrDB

Title: All non-redundant GenBank CDS translations+PDB+SwissProt+PIR+PRF excluding environmental samples from WGS projects
Molecule Type: Protein
Update date: 2018/10/11
Number of sequences: 171418145

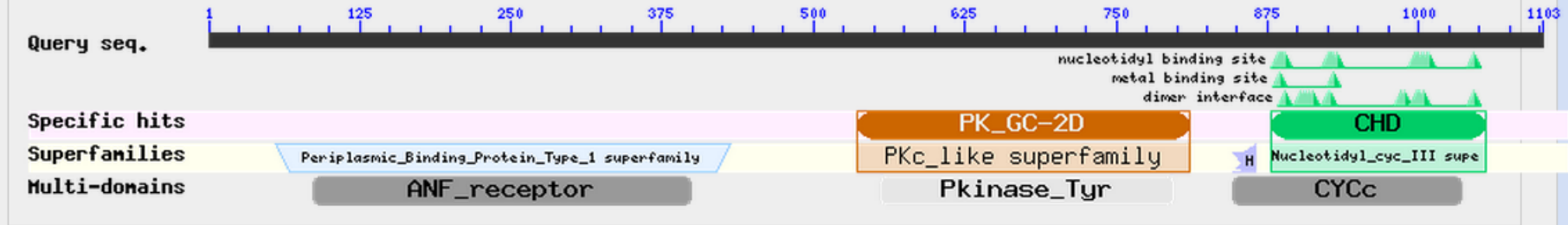
□ ~171M sequences Protein

□ RefSeq

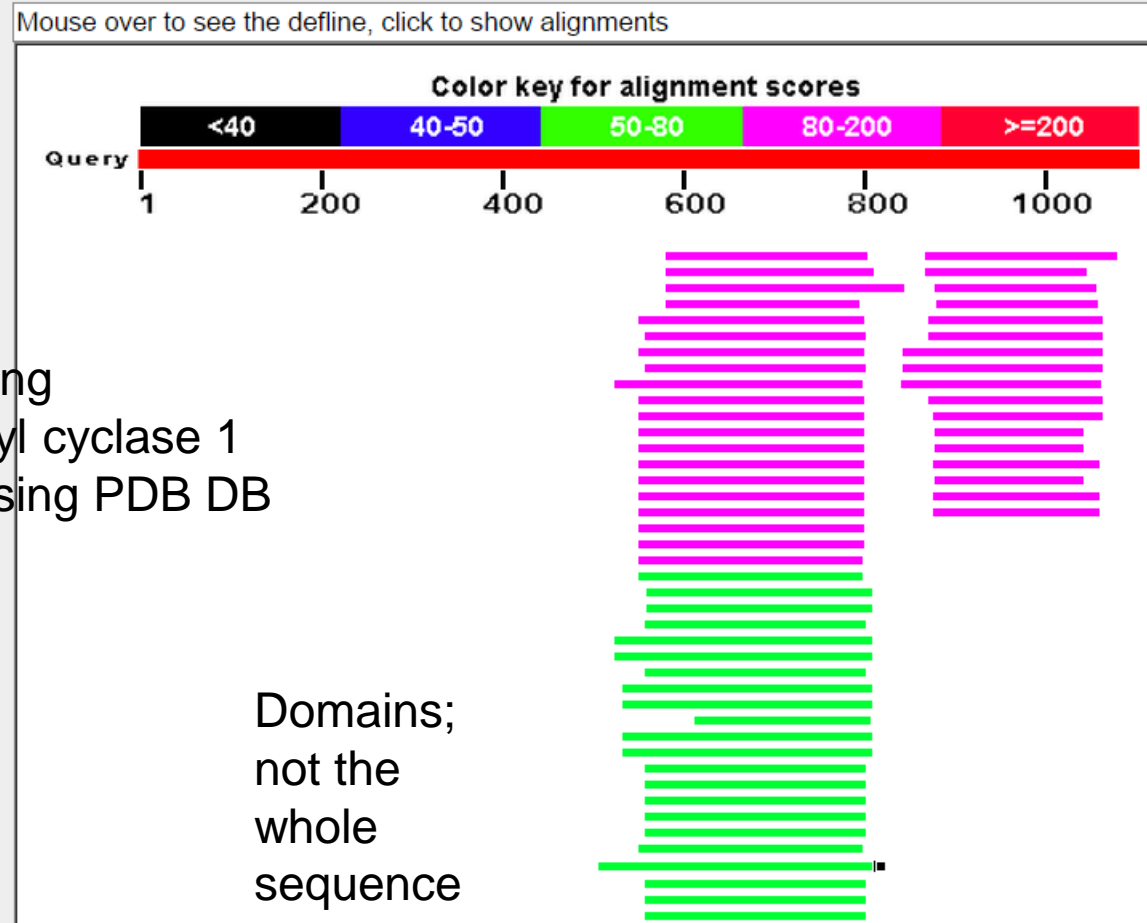
□ ~118M Sequences Protein

- Query

RID	96CT7C3E014 (Expires on 02-27 05:05 am)
Query ID	NP_000509.1
Description	hemoglobin subunit beta [Homo sapiens]
Molecule type	amino acid
Query Length	147



Distribution of 101 Blast Hits on the Query Sequence



Searching using
retinal guanylyl cyclase 1
isoform X1 using PDB DB

$$E = kmne^{-\lambda S}$$

Expected Threshold

- E-value

- “Number of different alignments with scores equal or greater than some score S that are expected to occur in a DB search by chance”

- Short queries

- Score is inversely prop to E-value

Query human insulin
Hit: insulin-like peptide3 from Drosophila

E value means that to get a score of 31.6 bits or better is expected by chance 1 in 20 times
(for a given DB/choice of parameters)

(a) Default: conditional compositional score matrix adjustment

Insulin-like peptide 3 [Drosophila melanogaster]

Sequence ID: [ref|NP_648360.2|](#) Length: 120 Number of Matches: 1

Range 1: 32 to 114 [GenPept](#) [Graphics](#)

Score	Expect	Method	Identities	Positives	Gaps
31.6 bits(70)	0.050	Compositional matrix adjust.	21/88(24%)	40/88(45%)	12/88(14%)
Query 29	HLCGSHLVEALYLVCGERGFFYTPKTRREAEDLQVGQVELGGGPGAGSLQPLALEGSLQ-				
	LCG L E L +C ++ T+R ++ Q++ G L+ L + S+Q				
Sbjct 32	KLCGRKLPETLSKLCV---YGFNAMTKRTLDPVNFNQID--GFEDRSLLERLLSDSSVQM				
Query 88	-----KRGIVEQCCTSICSLYQLENYC 109				
	+ G+ ++CC C++ ++ YC				
Sbjct 87	LKTRRLRDGVFDECCLKSCTMDEVLRVC 114				

E-values

$$E = kmne^{-\lambda S}$$

- Default value is 10
 - What this means is, At this E-value, 10 hits with score or better than the alignment score S are expected by chance.
 - Also assumes that you search using a random query with similar length of your actual query
- When you have a small query,

>Query
ELVIS

i Your search parameters were adjusted to search for a short input sequence.

Higher E values (200,000)
are set because shorter query
cannot get larger scores

Search Parameters	
Program	blastp
Word size	2
Expect value	200000
Hitlist size	100
Gapcosts	9,1
Matrix	PAM30
Filter string	F
Genetic Code	1
Window Size	40
Threshold	11

Short queries

☒ Automatically adjust parameters for short input sequences

Important Parameters

- E-value cut-off
 - Low (Example, $1\text{E}-06$)
- Word size (WS)
 - Low WS: Higher sensitivity
 - High WS: Higher Specificity
- Match/mismatch Score and Gap Costs
 - Extension scores
- Filters and masking
 - Substitute for RepeatMasker

NM_000518.4
HBB (homo sapiens) with nr
and other default options
BLASTN search

ref|NM_000518| (626 letters)

RID [TFNNRRFA01R](#) (Expires on 07-28 05:31 am)Query ID [gi|28302128|ref|NM_000518.4|](#)

Description Homo sapiens hemoglobin subunit beta (HBB), mRNA

Molecule type nucleic acid

Query Length 626

Database Name nr

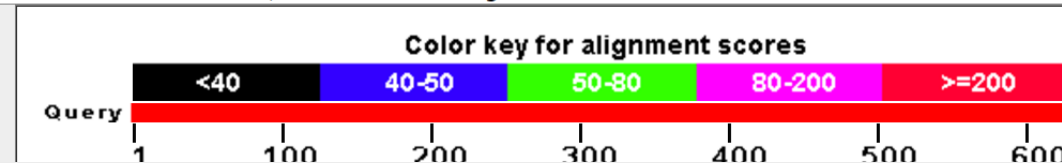
Description Nucleotide collection (nr)

Program BLASTN 2.4.0+ ▶ [Citation](#)Other reports: ▶ [Search Summary](#) [\[Taxonomy reports\]](#) [\[Distance tree of results\]](#)

Graphic Summary

Distribution of 103 Blast Hits on the Query Sequence

Mouse over to see the define, click to show alignments



Formatting options

[Reformat](#)

Show

Alignment as

HTML ▼

☐ Old View[Reset form to defaults](#)

Alignment View

Pairwise ▼

Display

☒ Graphical Overview☐ NCBI-gi☒ CDS feature

Masking

Character: Lower Case ▼

Color: Grey ▼

Limit results

Descriptions: 100 ▼

Graphical overview: 100 ▼

Line length: 60 ▼

Reformatting options

Range 1: 1 to 626 [GenBank](#) [Graphics](#) ▼ Next Match ▲ Previous Match

Score	Expect	Identities	Gaps	Strand
1157 bits(626)	0.0	626/626(100%)	0/626(0%)	Plus/Plus
CDS:hemoglobin subun	1			M V H
Query	1	ACATTTGCTTCTGACACAACTGTGTTCACTAGCAACCTCAAACAGACACCATGGTGCATC		60
Sbjct	1	ACATTTGCTTCTGACACAACTGTGTTCACTAGCAACCTCAAACAGACACCATGGTGCATC		60
CDS:hemoglobin subun	1			M V H
CDS:hemoglobin subun	4	L T P E E K S A V T A L W G K V N V D E		
Query	61	TGACTCCTGAGGAGAAGTCTGCCGTTACTGCCCTGTGGGGCAAGGTGAACGTGGATGAAG		120
Sbjct	61	TGACTCCTGAGGAGAAGTCTGCCGTTACTGCCCTGTGGGGCAAGGTGAACGTGGATGAAG		120
CDS:hemoglobin subun	4	L T P E E K S A V T A L W G K V N V D E		

Mismatches will be
in pink

CDS:hemoglobin subun	4	L T P E E K S A V T A L W G K V N V D E	
Query	61	TGACTCCTGAGGAGAAGTCTGCCGTTACTGCCCTGTGGGGCAAGGTGAACGTGGATGAAG	120
Sbjct	186	TGACTCCTGAGGAGAAGTCTGCCGTTACTGCCCTGTGGGGCAAGGTGAACGTGGATGAAG	245
CDS:PREDICTED: hemog	4	L T P E E K T A V T T L W G K V N V D E	

Standalone BLAST

- Hands on after lecture
- https://blast.ncbi.nlm.nih.gov/Blast.cgi?PAGE_TYPE=BlastDocs&DOC_TYPE=Download

Works like Google Search

BLAST ALGORITHM (using BLASTP as an example)

First Phase

>QUERY

MESADFYEAEP RP PMSSHLQSP PHAPSSAAFGFPRGAGPAQPPAPPAAPEPLGGICEHET
SIDI **SAYIDPAAFND** EFLADLFQHSRQQEKAKAAVGPTGGGGGGDFDYPGAPAGPGGAVM
PGGAHGPPPGYGCAAAGYLDGRLEPLYERVGAPALRPLVIKQEPREDEAKQLALAGLFP
YQPPPPPPSHPHPHPPPAHLAAPHLQFQIAHCGQTTMHLQPGHPTPPPTPVPSHPAPA
LGAAGLPGPGSALKGLGAAHPDLRASGGSGAGKAKKSVDKNSNEYRVRRENNIAVRKSR
DKAKQRNVETQQKVLELTSDNDRLRKRVEQLSRELDTLRGIFRQLPESSLVKAMGNCA

- Query

- Broken down into **word pairs**

- NT: 16-256; Proteins: 2-3:

- Threshold (T=11)

- Sliding window approach

- For each word (ex. SAY), the synonyms were formed and high scoring (BLOSUM matrix) words will be chosen.

SAYIDPAAFND

SAY **Score**

AYI **Score**

YID **Score**

IDP **Score**

DPA **Score**

PAA **Score**

AAF **Score**

AFN **Score**

FND **Score**

- Scores using Matrix for word pairs are collected

Word size 3, for 20 aa
there can be $20^3 = 8000$
possible words

How does BLAST work?

- SAY Score: $4+4+7 = 15$
- AYI Score: $4+7+4 = 15$
- YID Score: $7+4+6 = 17$

Scoring
Matrix to
calculate
Scores

- Now establish a cut-off (say 15)
 - Then only **YID** and other words that are above the cutoffs are retained

SAYIDPAAFND

SAY	Score
AYI	Score
YID	Score
IDP	Score
DPA	Score
PAA	Score
AAF	Score
AFN	Score
FND	Score

Phase 1: Setup: compile a list of words ($w=3$) above threshold T

- Query sequence: human beta globin NP_000509.1 (includes ...VTALWGKVNVD...). This sequence is read; low complexity or other filtering is applied; a “lookup” table is built.

- Words derived from query sequence (HBB): VTA TAL ALW **LWG** WGK GKV KVN VNV NVD

- Generate a list of words matching query (both above and below T). Consider **LWG** in the query and the scores (derived from a BLOSUM62 matrix) for various words.

- Generate similar lists of words spanning the query (e.g. words for **WGW**, **GWG**, **WGK**...).

examples of
words \geq
threshold 12

threshold

examples of
words below
threshold

LWG $4+11+6=21$

IWG $2+11+6=19$

MWG $2+11+6=19$

VWG $1+11+6=18$

FWG $0+11+6=17$

AWG $0+11+6=17$

LWS $4+11+0=15$

LWN $4+11+0=15$

LWA $4+11+0=15$

LYG $4+ 2+6=12$

LFG $4+ 1+6=11$

FWS $0+11+0=11$

AWS $-1+11+0=10$

CWS $-1+11+0=10$

IWC $2+11-3=10$

BLASTN

- First phase is slightly different than BLASTP
- Algorithm demands exact matches
 - Default word size is 11 (adjustable by user)
- Choosing a lower word length
 - Slower more accurate

Phase 2

- The selected words are now used to search for sequences that contain these words
- Create a hash table index with the locations of the hits for each word
- Perform two more searches
 - Un-Gapped extensions (first)
 - Followed by gapped extensions
 - Hits above certain score are saved

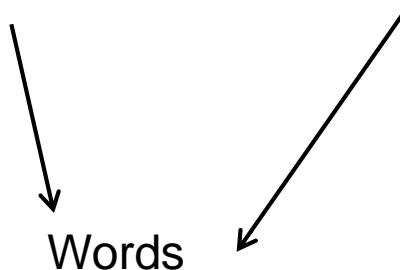
BLAST

- Scans the DB for matches to the words that are present above the Threshold Values
- Requires two hits within the target sequence (the new searched sequence)
- BLAST will set aside sequences with matches above Threshold for further analysis

Query 178 AFGWARVALVTAPQDLWVEAGRSLSTAL**RAR**GLPVASVTSMEPLDLSGA**REA**LRKVRDGP 237
Sbjct 148

RAR

REA



No need for exact match, but have to be in the list

Example:

Sequence = gaacgcctgcgcgatcagcataaaaaataa

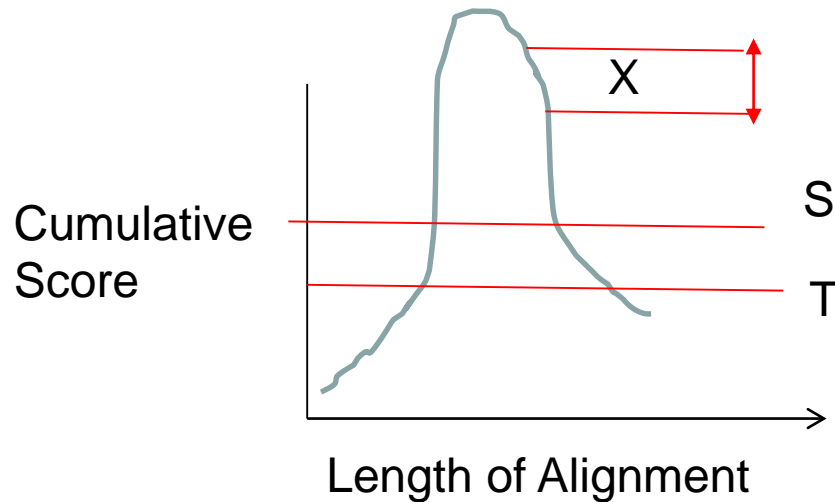
word length = $W = 7$ (there are 24 words possible)

For 'word' = ctgcgcg, a match in the database is found and then a local alignment done until a gap appears.

So...

query	=	ctgcgcg	7		ctgcgcgatcag	19
match	=	ctgcgcg	2598356		ctgcgcgatcag	2598367

Peak will tell us what
Length we are going to
focus



Based on Andy Baxevanis Book and
Seminars

High Scoring Segment

Extension

Three letter scores that are
greater than T will be carried
over to the next step

We keep adding in both
directions, more matches than
mismatches, so the score keeps
going up

Neighborhood threshold (S)
Everything above S will be
reported in BLAST results

When the score starts
decreasing then we go back and
pick a Window (X)

BLAST

Sequence ID: [pdb|3MJP|A](#) Length: 141 Number of Matches: 1

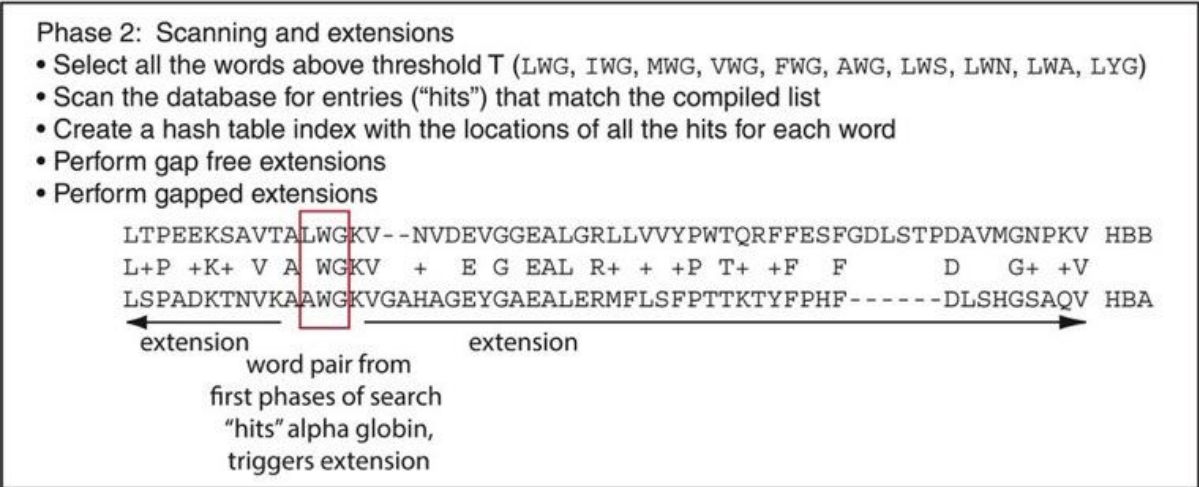
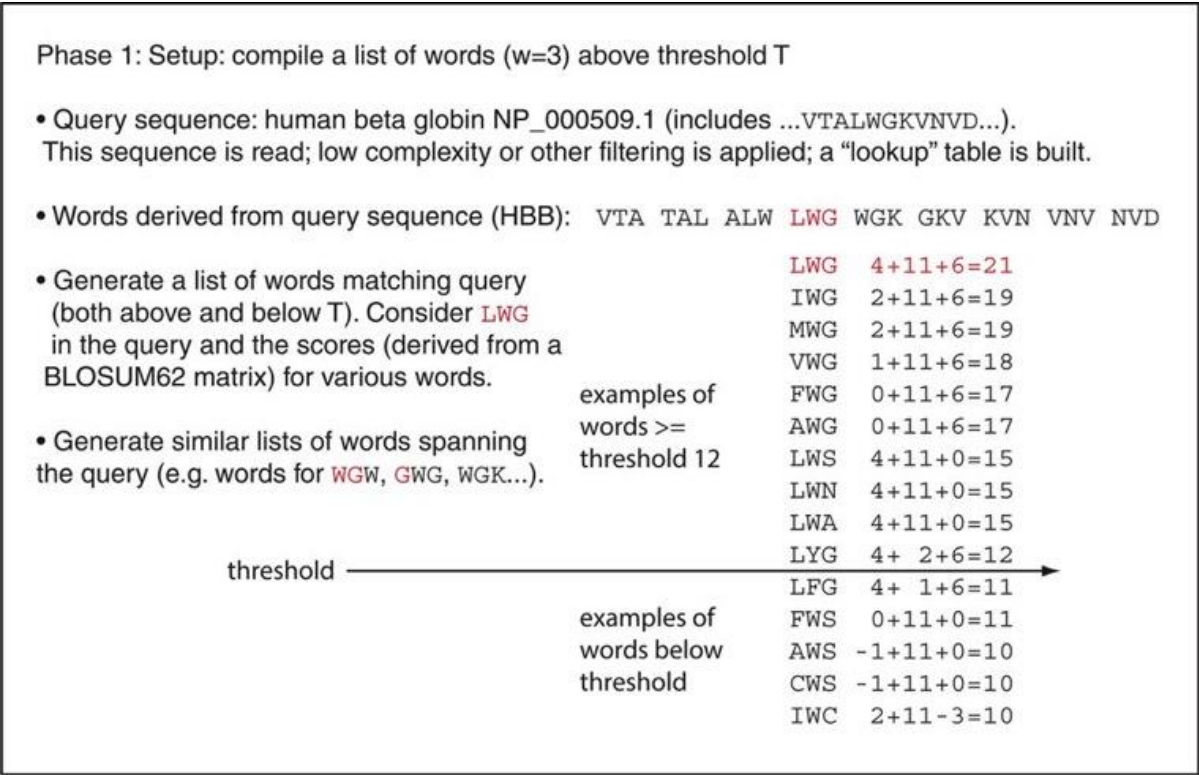
```
Query 178 AFGWARVALVTAPQDLWVEAGRSLSTALRARGLPVASVTSMEPLDLSGARREALRKVRDGP 237
Sbjct 148                                     RAR                               REA
```

```
Query 178 AFGWARVALVTAPQDLWVEAGRSLSTALRARGLPVASVTSMEPLDLSGARREALRKVRDGP 237
          G AR                      GR          RARGLPVA VTSMEP DLSGAREA          GP
Sbjct 148 --GAAR-----GR-----WRARGLPVALVTSMEPSDLSGARREAL--SASAGP 184
```



Extension until the
score drops

Fig 4.12 from the Pevsner Book III
Edition PLEASE DO NOT DISTRIBUTE
Copyright figure



Phase 3

- Traceback
 - Identify the locations of INDELS and matches from phase 2
 - If applicable, use composition-based statistics (BLASTP, TBLASN)
 - Generate final gapped alignment

Summary

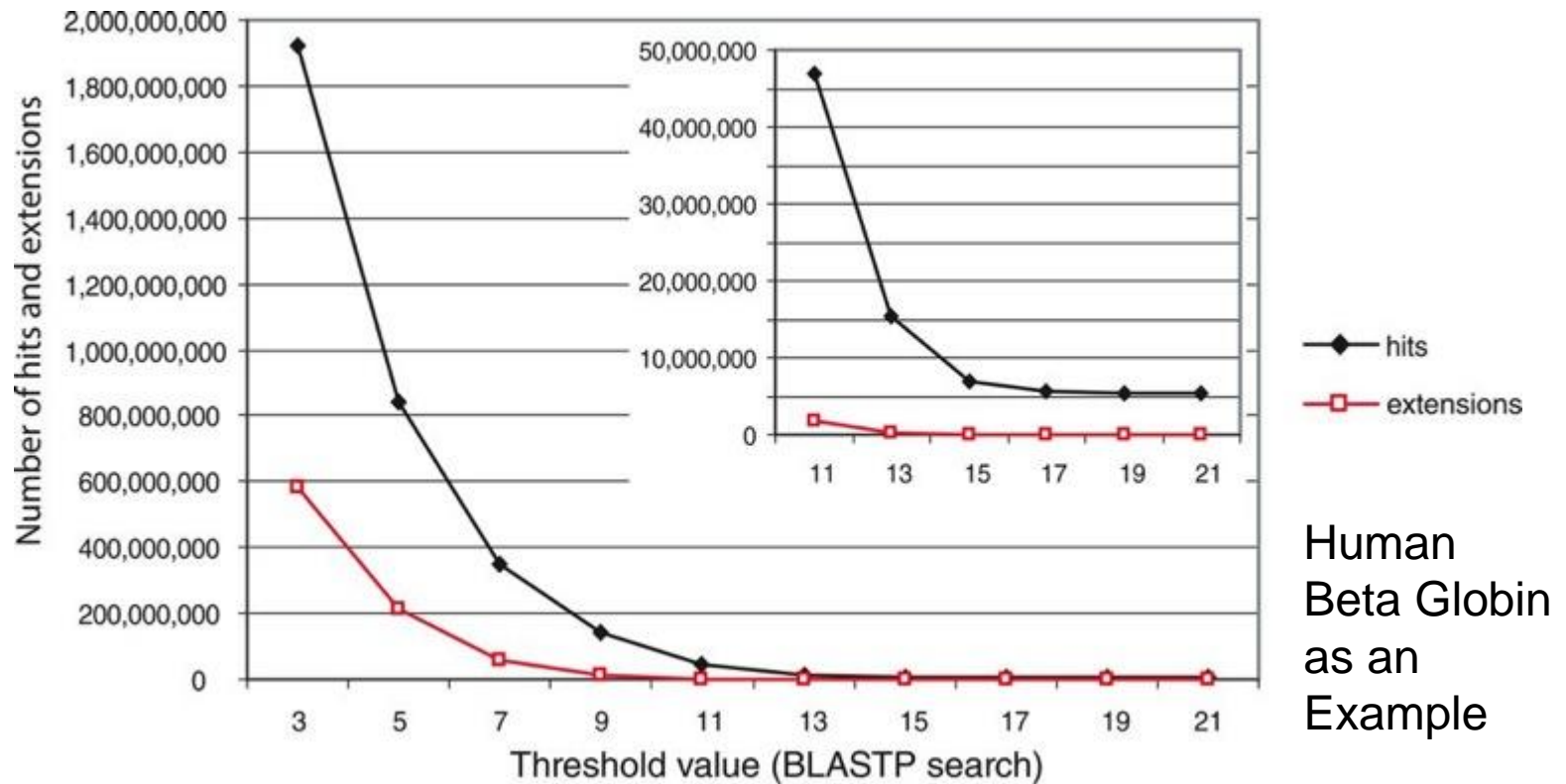
- BLAST
 - Heuristic algorithm (optimized for speed/sensitivity)
 - Threshold is increased
 - Speed is increased/fewer hits
 - Distantly related are missed
 - Threshold is lowered
 - Speed is lowered/large number of matches
 - Sensitivity is increased

Threshold size
T = 10 (BLASTP)
will compile words
≥ 10

Impact of Threshold Score(T)

General Parameters	
Max target sequences	100 <small>Select the maximum number of sequences to search</small>
Short queries	<input checked="" type="checkbox"/> Automatically adjust p
Expect threshold	10
Word size	6
Max matches in a query range	0

Fig 4.13 from the Pevsner Book III Edition PLEASE DO NOT DISTRIBUTE Copyright figure



What matrix to pick?

BLOSUM	Suitable For; Based on experience	% similarity
90	Short alignments; highly similar	70-90
80	Best for identifying family members	50-60
62	MOST EFFECTIVE for identifying all potential similarities (default in NCBI)	30-40
30	Longer/weaker local alignments	<30

Is One matrix is enough?

- David Altschul prescribes a “Triple Strategy”
- Pick the default and a higher/lower BLOSUM_n
- Analyze and pick the appropriate matrix

Statistics of Alignments

Let us begin with a simple
diagram that explains global
alignment

Global Alignment

- Distribution behavior of global alignments is not known (not Gaussian/normal)
 - Usually approximated using simulations

Local Alignment

- Statistics/distributions are known
 - Altschul many papers
- Start with Ungapped alignments
 - Random
 - Gapped alignments
 - Proteins

Local Alignment

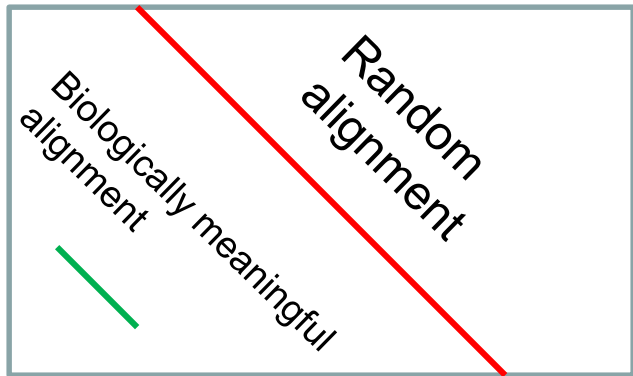
- To assess how high a score can occur by chance, we need random sequences and their scores
 - Conditions
 - **Expected score** for aligning random pair of amino acids has to be NEGATIVE

Expected Score matrix constraints

- **Condition I:** For alignment algorithms that seek to capture the local alignments of variable length should have a negative expected score (a necessary condition)

$$\sum_{i,j} p_i p_j s_{i,j} < 0$$

For local alignments of random sequences,
Negative Expected Score



Condition II: At **least** one of the score has to be positive ($s_{i,j}$)

Log-odds Scores

- “with the previous two assumptions, the scores of any substitution matrix (with a negative expected value and at least one positive score) can be written in the form” (Karlin & Altschul, PNAS, 872264(1990))

$$s_{i,j} = \frac{\left(\ln \frac{q_{i,j}}{p_i p_j} \right)}{\lambda} \stackrel{\text{Change of base}}{=} \log \left(\frac{q_{i,j}}{p_i p_j} \right)$$

λ : Scaling Parameter

x, a, b are all positive

$a \neq 1; b \neq 1$

$$\log_8 x = \frac{\ln x}{\ln 8} = \frac{1}{\ln 8} \ln x$$

What is a search space

- “Given a scoring system, how many **distinct local alignments** with score $\geq S$ (S is some number) can one **find by chance** by comparing **two random sequences** of length m and n ” S. Altschul

Random Subject or Database (length n residues); concatenating all the DB sequences

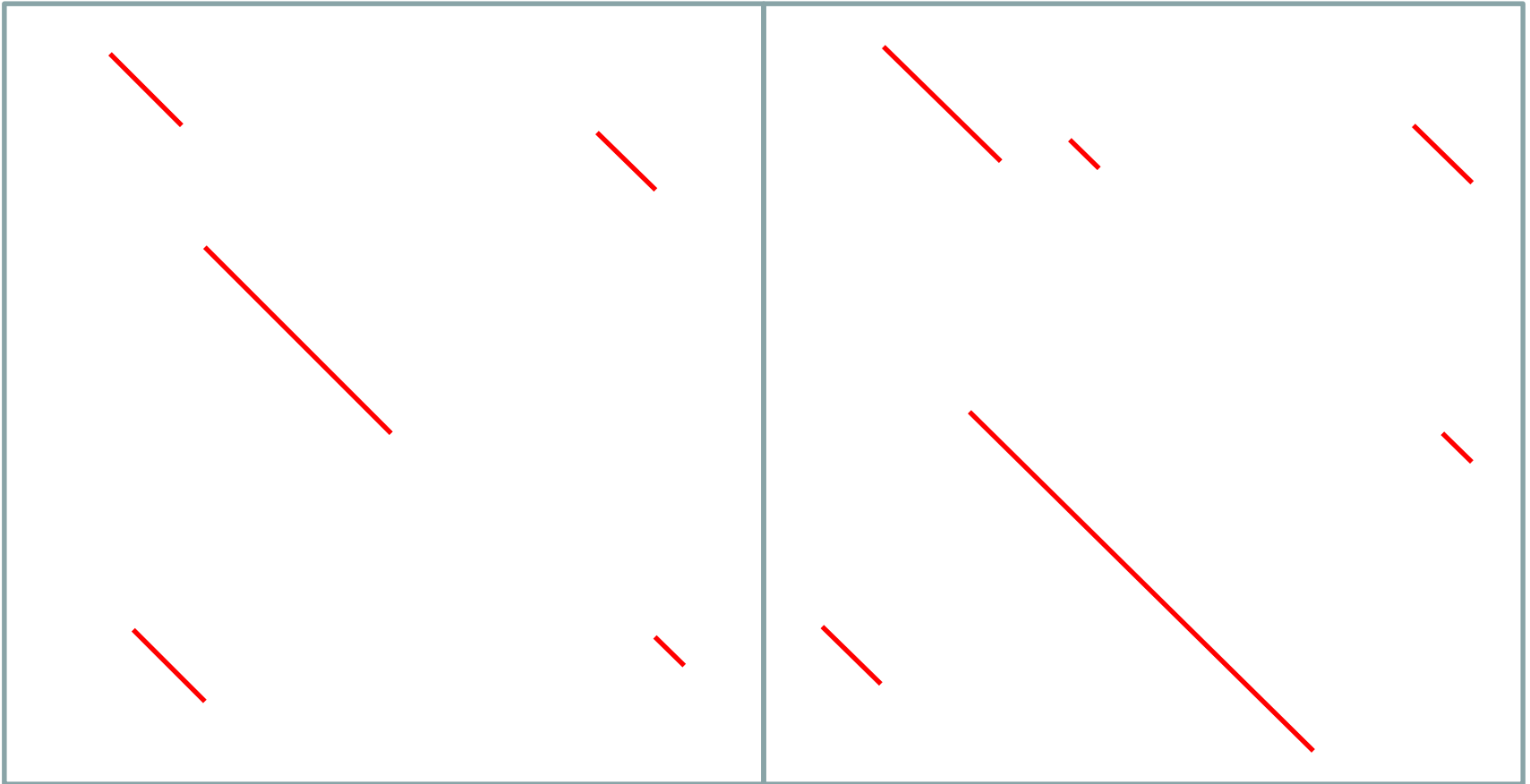
Random
Query
(length m)

Search Space, $N = m \cdot n$

- Answer: $E(S, m, n)$, E = expected score will depend on S , m and n

Number of Random High-Scoring Alignments ~ Search Space Size

$$E(S, m, n) \propto mn \quad \text{Asymptotic result}$$



Double the size then HAS will roughly double.

The # of random alignments with Score $\geq S$ should decrease Exponentially with S

Any scoring system, the probability that the optimal local alignment that starts at a particular position has a score $\geq S$ decreases exponentially with S

$$E(S, m, n) \propto e^{-\alpha S}$$

Alpha turns out to be the same as λ

$$\begin{array}{ccc} E(S, m, n) \propto e^{-\alpha S} & \longrightarrow & E = kmne^{-\lambda S} \\ E(S, m, n) \propto mn & & \end{array}$$

E: “Number of different alignments with scores equal or greater than some score S that are expected to occur in a DB search by chance”

Poisson

- Prob of finding 0 alignments (or none) with score $\geq S$ is

$$e^{-E}$$

$$P(k \text{ events in interval}) = \frac{e^{-L} L^k}{k!}$$

The average number of events in an interval is designated as L.

- Prob. of finding at least one alignment with score $\geq S$ is

$$p = 1 - e^{-E}$$

- This is called “p-value” associated with S.
- When $E \leq 0.1$, $p \sim E$

$$E = kmne^{-\lambda S}$$

Derivation of Normalized Scores

- To calculate E-value associated with a S, we need to know λ and K.

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

- Refer you to Altschul papers on derivation
- But, these values can be wrapped into a reduced form as shown above, then S' can be easily connected to E
- Refer back, $N = \text{Search Space}$

$$E = \frac{N}{2^{S'}}$$

Number of alignments with $\geq S$

$$E = (kmn)e^{-\lambda S}$$

$$E = Nke^{-\lambda S}$$

$$E = Ne^{-\ln_e k} e^{-\lambda S}$$

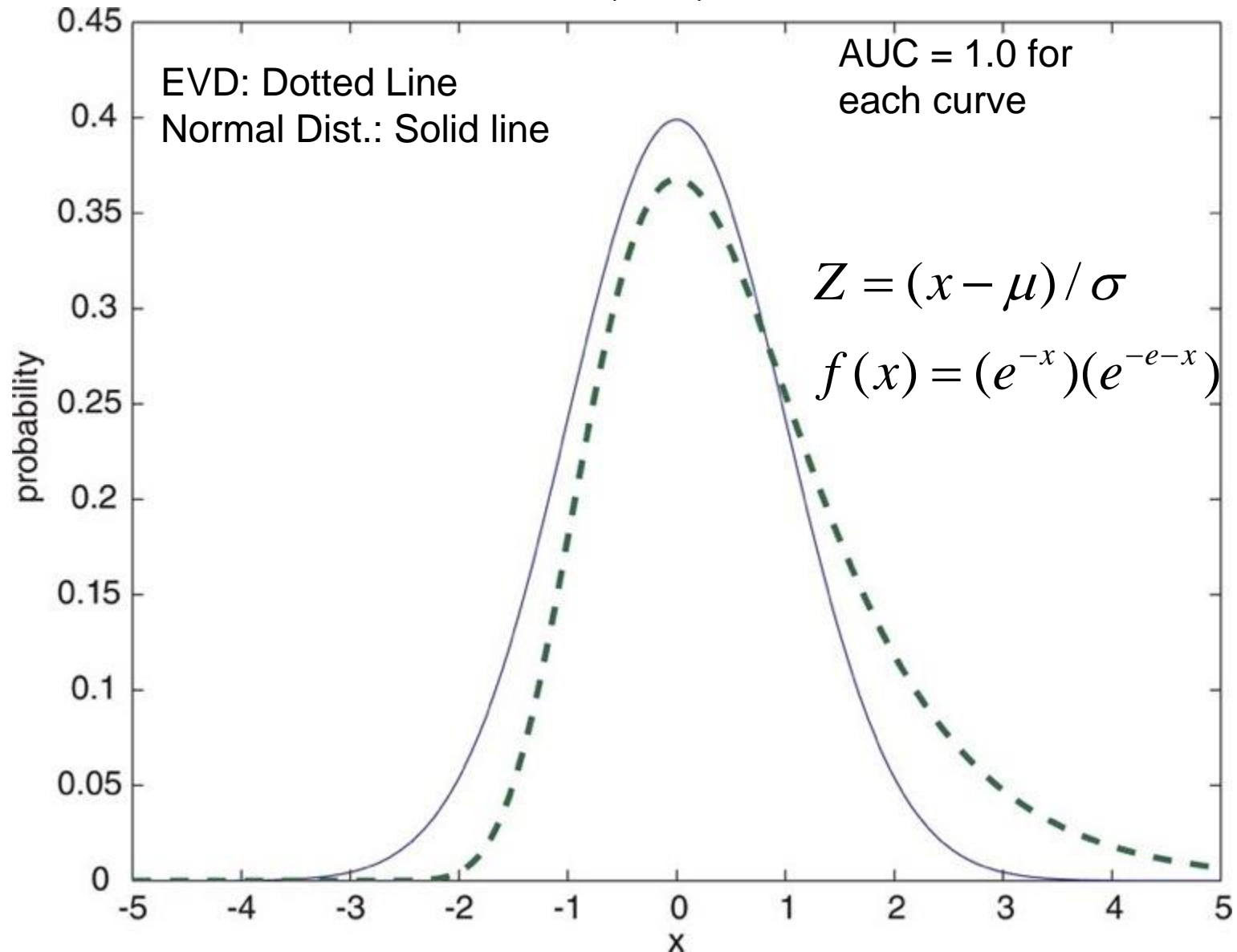
$$E = Ne^{-\left[\frac{(\lambda S - \ln_e k)}{\ln_e 2}\right] \ln_e 2}$$

$$E = N2^{-S'}$$

$$E = \frac{N}{2^{S'}}$$

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

Query compared to a set of random seq of same length as query. The alignment scores will take a extreme value distribution (EVD)



Ungapped → Gapped

- Everything discussed up to this apply to gapped alignment (as long gap scores are negative enough; not close to zero)
- **Not proved up until now!**
- According to Altschul, gapped alignments there is no way to statistical theory (?) to calculate the statistical parameters, K and λ but we can estimate them.

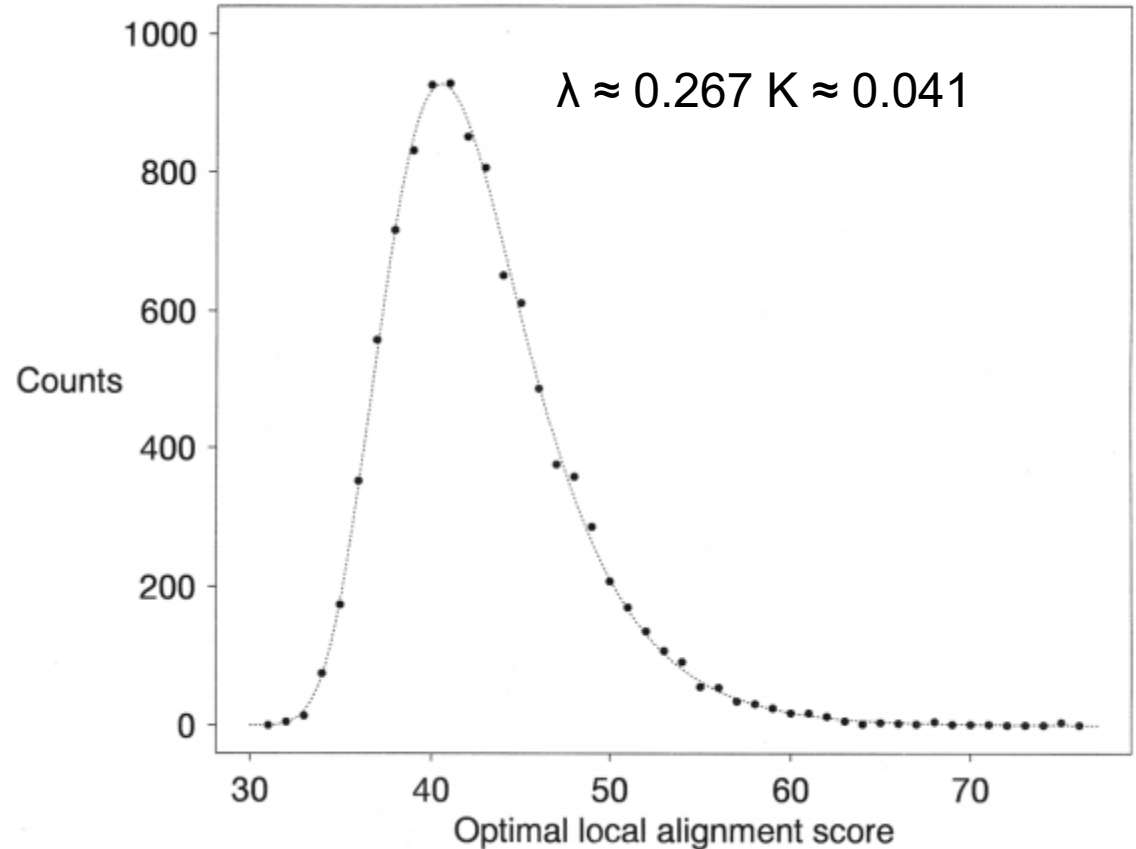
$$S' = \frac{(\lambda S - \ln K)}{\ln 2}$$

Local Alignment with Gaps

Simulation:

10,000 pairs of random protein sequence, each of length 1000 were compared using BLOSUM-62 substitution score, gap score of -11-k for a gap of length k

After several simulations, Altschul has plotted a histogram of how many times, he saw the scores. He fitted them to EVD and estimated Lambda and K



Random Sequence, Ungapped to Real proteins

- The theory still holds except for the following cases
 - Low-complexity filtering
 - Mask the sequence segments by giving a negative score
- Let us start with the following equation
 - $E = kmNe^{-\lambda S}$
 - E gives the # of HSPs found purely by chance

$$E = kmne^{-\lambda S'}$$

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

- Raw score (S)
 - Sub. Matrix (gap penalty etc)
- Bit Score (S', scaled value)
 - Bit scores can be compared even using different scoring matrices
- E values are derived from Bit Scores (S')
- Prob of chance alignment occurring with the score or better

$$p = 1 - e^{-E}$$

Why BLAST doesn't report P?

It is easier to think of the number of HSAs rather than the probability values; High-Scoring Alignment (HAS)

E	p
10	0.99995460
5	0.99326205
2	0.86466472
1	0.63212056
0.1	0.09516258
0.05	0.04877058
0.001	0.00099950
0.0001	0.0001000

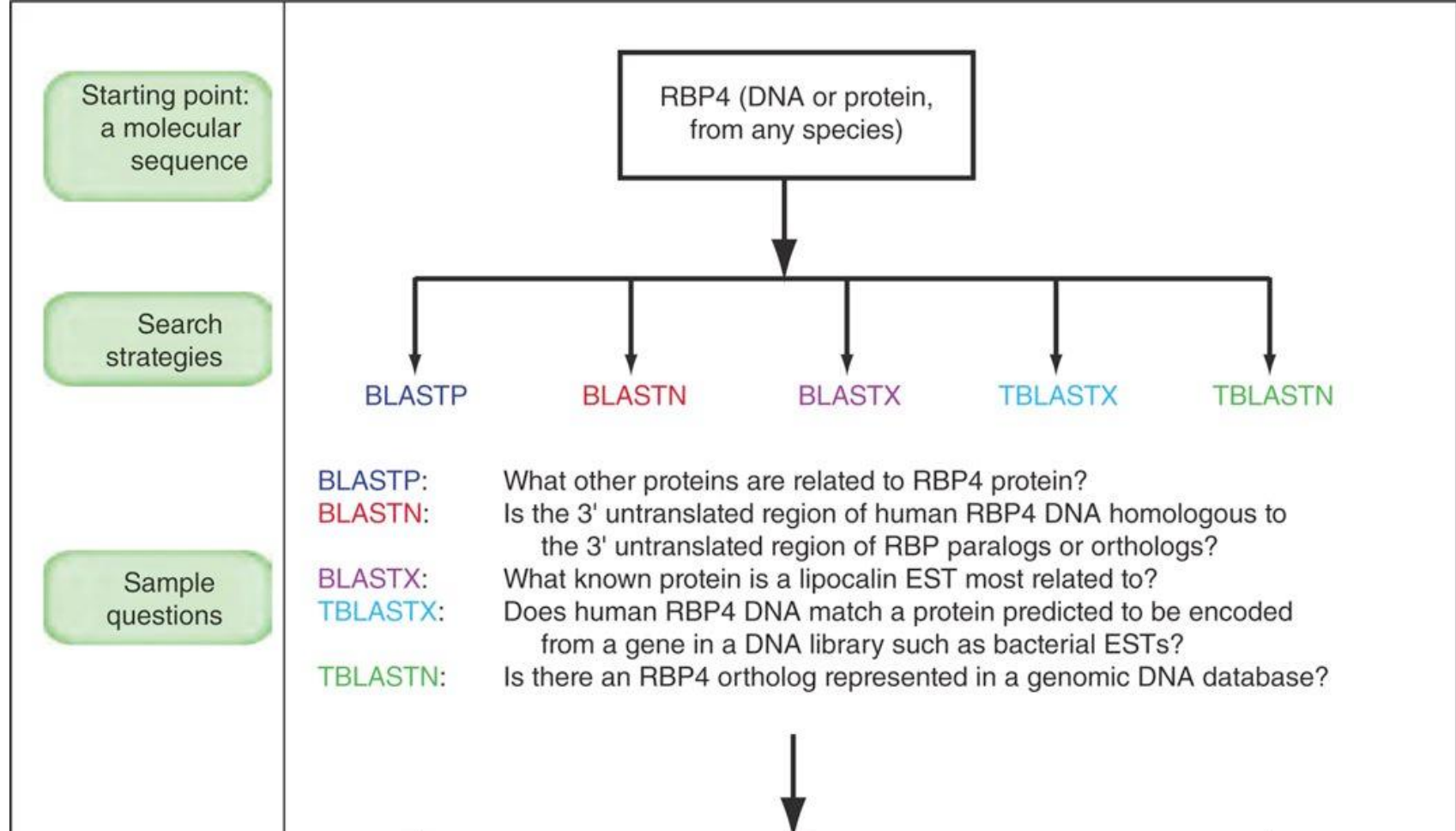
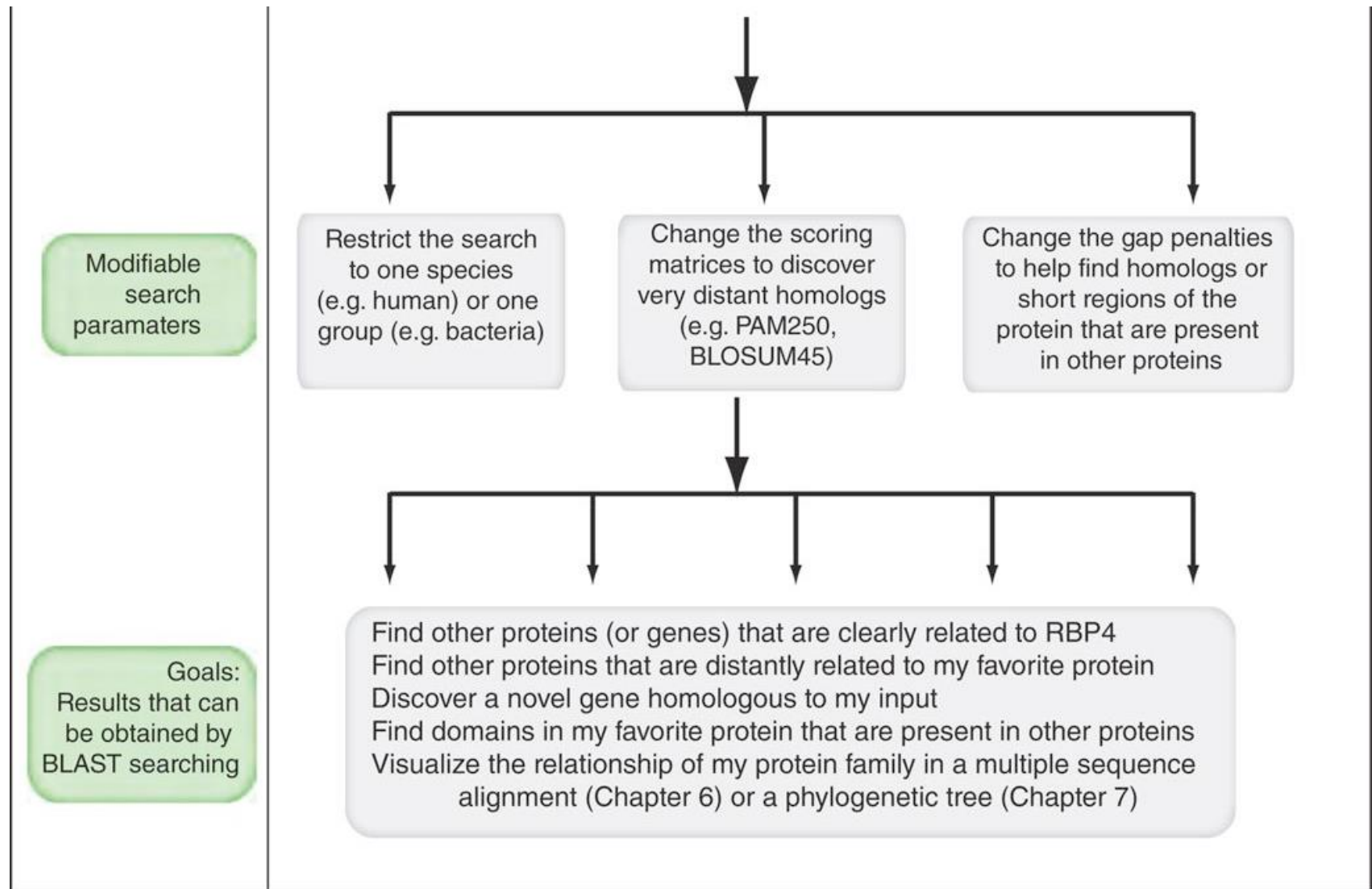


Fig 4.15 from Pevsner III edition

Fig 4.15 from Pevnsner III edition



Principles of DB searching using Calcin family as an example

- Lipocalcin family proteins in this example share very limited sequence similarity
 - RBP4, NP_006735.2
 - Odorant-binding protein (OBP)
- BLAST
 - DB: nr; organism: Homo sapiens; others: def
 - Restricting the output only to Human RefSeq proteins (how can we do this?)

Too many hits?

- Refseq
 - Nr database
 - Restrict to only specific organisms
 - Restrict to the domain of interest
 - How to find this?
 - UniProt
 - Adjust
 - Scoring matrix
 - Expect value

Too Small hits?

- How can this happen?
 - Exploring microbial/viral genomes
 - Only few are sequenced
 - Reset the BLAST page (to remove prev limits)
 - Matrices
 - High PAM or lower BLOSUM
 - Include all DBs (HTGS/GSS)
 - Search to include model sequences
 - Finally, use HMM based searches (PSI-BLAST etc)

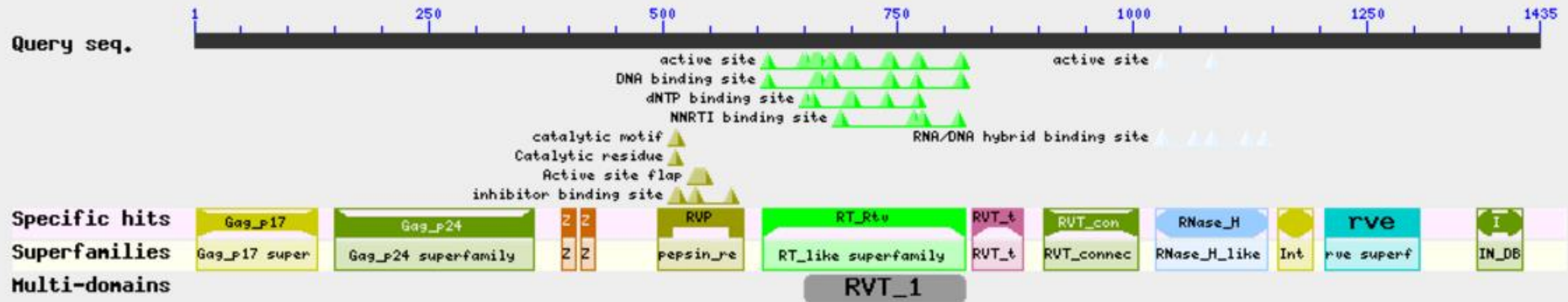
HIV-1 Pol: Second Example

- Multiple domain protein
- <http://www.ncbi.nlm.nih.gov/gene/155348>
- <http://www.uniprot.org/uniprot/P04585>
- 1435 aa
- What will happen when do a blastp search for this query?

NP_057849.4

Conserved Domains

Putative conserved domains have been detected, click on the image below for detailed results.



“New Gene” == discovery of some DNA sequence in a DB that has not been annotated yet

Find-a-Gene Tips

Start with the sequence
of a known protein

TBLASTN

Search a DNA database (e.g. HTGS,
dbEST, or genomic sequence
from a specific organism)

TBLASTN

protein

6



TBLASTN is used to translate every DNA sequence in a database into six potential proteins,
and then to compare your protein query against each of those translated proteins.

Inspect the output

Search your DNA or protein
against a protein database (nr)
to confirm you have
identified a novel gene

BLASTX or
BLASTP nr

Find matches...

- [1] to DNA encoding known proteins (not novel)
- [2] to DNA encoding related proteins (novel!)
- [3] to false positives

Things to Report

- Query sequence
- TBLASTN
 - What DB? What Matrix; what non-optional parameters?
 - Hits (follow the font and other details as Dr. Pevsner has suggested)

Things to Report

- Use additional BLASTX/BLASTP to confirm that the protein that you had identified is novel
 - *(follow the suggestions of Prof. Pevnser on what is novel; page 159 of the book)*
 - Again list DB, matrix; hits (top 10)
 - Name your protein, example “Anguillicola Globin”
 - Because of the organism and family it belongs to

Things to Report

- Carry out Multiple sequence alignment
 - Your novel protein + 5 or 10 (max 30) from the novel protein speculated family
- Create a phylogenetic tree
- Secondary/tertiary structure of your novel protein
- Provide whether the gene is under positive/negative selection (optional)
- Significance of the novel gene

Computer Lab

- 4-1, 4-2, 4-3, 4-5 and 4-9

Thanks

ravichandran@hood.edu