

Programming for Biology
Similarity Searching II –

Practical search strategies

Bill Pearson
wrp@virginia.edu

CSHL Programming for Biology

1

1

Why is this material important?

- You might be asked to find a homolog
- You might be asked to what your gene/protein does
 - Annotated homologs are missed because databases are large and redundant
 - Short domains and short exons are missed because the “standard” matrix needs long alignments
 - Sometimes, alignments include non-homologous regions

CSHL Programming for Biology

2

2

Effective Similarity Searching

1. Always search protein databases (possibly with DNA – blastx, fastx)
 2. Use E()-values, not percent identity, to infer homology
 - $E() < 0.001$ is significant in a single search
-
1. Search smaller (comprehensive) databases
 - Less redundancy; higher sensitivity
 2. Change the scoring matrix for:
 - short sequences (exons, reads)
 - short evolutionary distances (mammals, vertebrates, a-proteobacteria)
 - high identity (>50% alignments) to reduce over-extension

CSHL Programming for Biology

3

3

Review – Sequence Similarity - Conclusions

- Homologous sequences share a common ancestor, but most sequences are non-homologous
- Always compare Protein Sequences
- Sequence Homology can be reliably inferred from statistically significant similarity (non-homology cannot from non-similarity)
- Homologous proteins share common structures, but not necessarily common functions
- Sequence statistical significance estimates are accurate (verify this yourself) $10^{-6} < E() < 10^{-3}$ is statistically significant

CSHL Programming for Biology

4

4

Similarity Searching II

1. What question to ask?
2. What program to use?
3. What database to search?
4. When to do something different (changing scoring matrices)
5. Is every aligned domain homologous?
6. (Tomorrow) – more sensitive methods (PSI-BLAST, HMMER)

CSHL Programming for Biology

5

5

1. What question to ask?

- Is there an homologous protein (a protein with a similar structure)?
- Does that homologous protein have a similar function?
- Does XXX genome have YYY (kinase, GPCR, ...)?

Questions not to ask:

- Does this DNA sequence have a similar regulatory element (too short – never significant)?
- Does (non-significant) protein have a similar function/modification/antigenic site?

CSHL Programming for Biology

6

6

2. What program to run?

- What is your query sequence?
 - protein – BLASTP (NCBI), SSEARCH (EBI)
 - protein coding DNA (EST) – BLASTX (NCBI), FASTX (EBI)
 - DNA (structural RNA, repeat family) – BLASTN (NCBI), FASTA (EBI)
- Does XXX genome have YYY (protein)?
 - TBLASTN YYY vs XXX genome
 - TFASTX YYY vs XXX genome
- Does my protein contain repeated domains?
 - LALIGN (UVA <http://fasta.bioch.virginia.edu>, EBI)

CSHL Programming for Biology

7

7

NCBI BLAST Server blast.ncbi.nlm.nih.gov

BLAST®

Home Recent Results Saved Strategies Help

Basic Local Alignment Search Tool

BLAST finds regions of similarity between biological sequences. The program compares nucleotide or protein sequences to sequence databases and calculates the statistical significance.

[Learn more](#)

Web BLAST

Nucleotide BLAST
nucleotide → nucleotide

blastx
translated nucleotide → protein

tblastn
protein → translated nucleotide

Protein BLAST
protein → protein

Always compare protein sequences

Enter organism common name, scientific name, or tax id

Human Mouse Rat Microbes

Search

CSHL Programming for Biology

8

8

NCBI BLAST Server

BLAST® » blastp suite Home Recent Results Saved Strategies

blastn **blastp** blastx tblastn tblastx

Standard Protein BLAST

BLASTP programs search protein databases using a protein query. more...

[Reset page](#) [Bookmark](#)

Enter Query Sequence

Enter accession number(s), gi(s), or FASTA sequence(s) [?](#) [Clear](#)

Query subrange [?](#)

From

To

Or, upload file no file selected [?](#)

Job Title

Enter a descriptive title for your BLAST search [?](#)

☐ Align two or more sequences [?](#)

Choose Search Set

Databases ☒ Standard databases (nr etc.): [New](#) ☐ Experimental databases [Try experimental clustered nr database](#) [?](#)
For more info see [What is clustered nr?](#)

Compare ☐ Select to compare standard and experimental database [?](#)

Standard

Database [?](#)

Organism [Optional](#) ☐ exclude [Add organism](#) [?](#)
Enter organism common name, binomial, or tax id. Only 20 top taxa will be shown. [?](#)

Exclude [Optional](#) ☐ Models (XM/XP) ☐ Non-redundant RefSeq proteins (WP) ☐ Uncultured/environmental sample sequences

CSHL Programming for Biology

9

9

NCBI BLAST Server

BLAST Basic Local Alignment Search Tool

[Home](#) [Recent Results](#) [Saved Strategies](#) [Help](#)

► [NCBI BLAST Home](#)

BLAST finds regions of similarity between biological sequences. [more...](#)

[New](#) Aligning Multiple Protein Sequences? Try the [COBALT Multiple Alignment Tool](#). [Go](#)

BLAST Assembled Genomes

Choose a species genome to search, or [list all genomic BLAST databases](#).

| | | |
|---|--|--|
| <input type="checkbox"/> Human | <input type="checkbox"/> Oryza sativa | <input type="checkbox"/> Gallus gallus |
| <input type="checkbox"/> Mouse | <input type="checkbox"/> Bos taurus | <input type="checkbox"/> Pan troglodytes |
| <input type="checkbox"/> Rat | <input type="checkbox"/> Danio rerio | <input type="checkbox"/> Microbes |
| <input type="checkbox"/> Arabidopsis thaliana | <input type="checkbox"/> Drosophila melanogaster | <input type="checkbox"/> Apis mellifera |

Basic BLAST

Choose a BLAST program to run.

| | |
|----------------------------------|--|
| nucleotide blast | Search a nucleotide database using a nucleotide query <i>Algorithms: blastn, megablast, discontinuous megablast</i> |
| protein blast | Search protein database using a protein query <i>Algorithms: blastp, psi-blast, phi-blast</i> |
| blastx | Search protein database using a translated nucleotide query |
| tblastn | Search translated nucleotide database using a protein query |
| tblastx | Search translated nucleotide database using a translated nucleotide query |

Specialized BLAST

Choose a type of specialized search (or database name in parentheses.)

| |
|--|
| <input type="checkbox"/> Make specific primers with Primer-BLAST |
| <input type="checkbox"/> Search trace archives |
| <input type="checkbox"/> Find conserved domains in your sequence (cds) |
| <input type="checkbox"/> Find sequences with similar conserved domain architecture (cdart) |

10

3. What database to search?

- Search the smallest comprehensive database likely to contain your protein
 - vertebrates – human proteins (40,000)
 - NCBI Landmark sequences (human, mouse, no rat)
 - Quest for Orthologs reference proteomes (1,000,000)
- Search a richly annotated protein set (SwissProt: 500,000, NCBI Landmark:)
- Always search NR (> 500 million) *LAST*
- Never Search “GenBank” (DNA)

CSHL Programming for Biology

11

11

Effective Similarity Searching

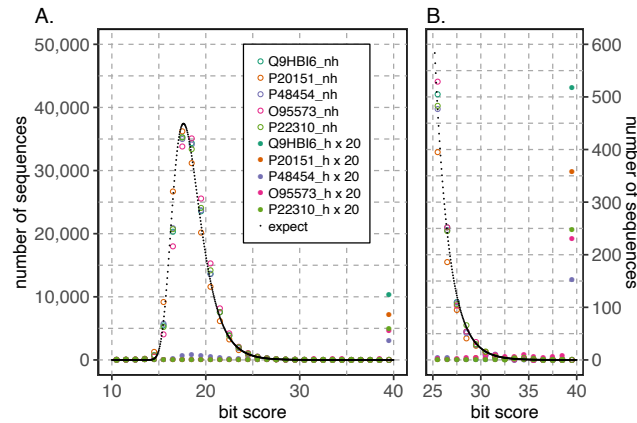
1. Always search protein databases (possibly with translated DNA)
 2. Use E()-values, not percent identity, to infer homology
 - $E() < 0.001$ is significant in a single search
-
1. Search smaller (comprehensive) databases
 2. Change the scoring matrix for:
 - short sequences (exons, reads)
 - short evolutionary distances (mammals, vertebrates, a-proteobacteria)
 - high identity (>50% alignments) to reduce over-extension
 3. Is every aligned residue homologous?
 - alignment overextension
 4. (Tomorrow) All methods (pairwise, HMM, PSSM) miss homologs, and find homologs the other methods miss

CSHL Programming for Biology

12

12

Homology inferences are reliable because
similarity statistics are accurate (I)
(we know how unrelated sequences behave)



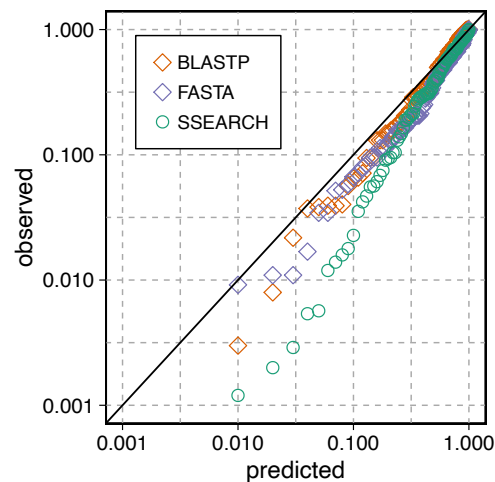
Distributions of similarity scores in searches with 5 human enzymes. Open circles (_nh) show scores for non-homologs. Closed circles show homolog (_h) scores.

CSHL Programming for Biology

13

13

Homology inferences are reliable because
similarity statistics are accurate (II)
(we know how unrelated sequences behave)



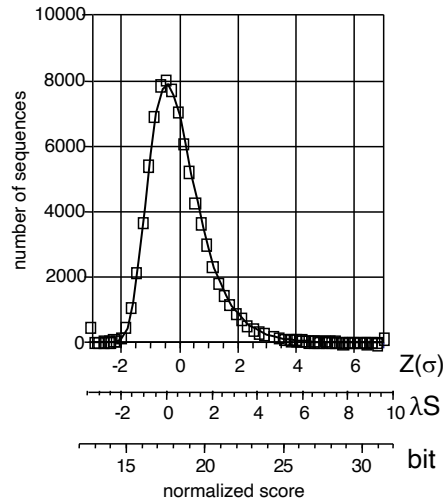
Reported (observed) and expected probabilities of the highest scoring unrelated sequence in searches with 100 human enzymes vs 78 complete proteomes (~1 million sequences).

CSHL Programming for Biology

14

14

Why smaller databases are better – statistics



$$S' = \lambda S_{\text{raw}} - \ln K m n$$

$$S_{\text{bits}} = (\lambda S_{\text{raw}} - \ln K) / \ln(2)$$

$$P(S' > x) = 1 - \exp(-e^{-x})$$

$$P(S_{\text{bits}} > x) = 1 - \exp(-mn2^{-x})$$

$$E(S' > x \text{ ID}) = P D$$

Bonferroni correction

$$P(B \text{ bits}) = m n 2^{-B}$$

$$P(40 \text{ bits}) = 1.5 \times 10^{-7}$$

$$E(40 \mid D=4000) = 6 \times 10^{-4}$$

$$E(40 \mid D=500E6) = 75$$

CSHL Programming for Biology

15

15

Local similarity statistics

$S' = \lambda S_{\text{raw}} - \ln K m n$ m : query length, n : subj length

$$S_{\text{bit}} = (\lambda S_{\text{raw}} - \ln K) / \ln(2)$$

$$P(S' > x) = 1 - \exp(-e^{-x})$$

$$P(S' > x) = e^{-x} \quad (\text{for } P < 0.1)$$

$$P(S_{\text{bits}} > \text{bits}) = 1 - \exp(-mn2^{-x})$$

$$P(S_{\text{bits}} > \text{bits}) = mn2^{-\text{bits}} \quad (\text{for } P < 0.1)$$

$$E(S', S_{\text{bits}} \text{ ID}) = P D$$

$$E(S_{\text{bits}} \text{ ID}) = D mn2^{-\text{bits}} \quad \text{Bonferroni correction}$$

$$\text{dblength} = D n$$

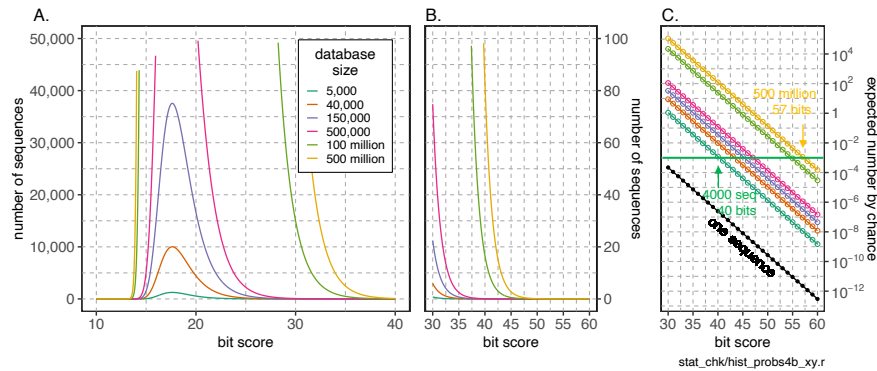
$$E(S_{\text{bit}}) = m \text{ dblength } 2^{-\text{bits}} \quad (\text{BLAST})$$

CSHL Programming for Biology

16

16

Smaller databases increase sensitivity



(More sophisticated algorithms – PSIBLAST, JACKHMMER – also improve sensitivity, and they work better with large databases)

CSHL Programming for Biology

17

17

NCBI – selecting sequences with Entrez

NCBI/ BLAST/ blastp suite

blastn blastp blastx tblastn tblastx

BLASTP programs search protein databases using a protein query. [more...](#)

Enter Query Sequence

Enter accession number, gi, or FASTA sequence [Clear](#)

Query subrange

From

To

Or, upload file no file selected [?](#)

Job Title

Enter a descriptive title for your BLAST search [?](#)

☐ Align two or more sequences [?](#)

Choose Search Set

Database [?](#)

Organism ☐ Exclude [+](#)

Optional Enter organism common name, binomial, or tax id. Only 20 top taxa will be shown. [?](#)

Entrez Query

Optional Enter an Entrez query to limit search [?](#)

CSHL Programming for Biology

18

18

3. What database to search?

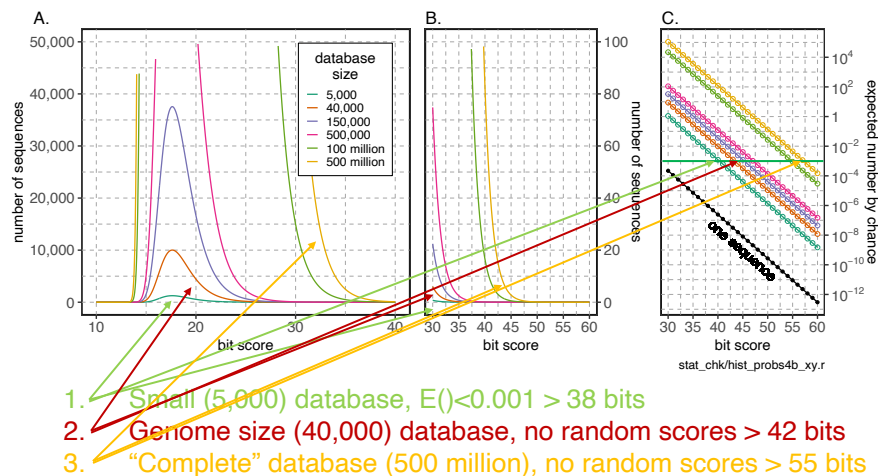
| Database | Size | Bits (0.001) |
|----------------|--------------|--------------|
| Landmark | 441 thousand | 47 |
| SwissProt | 480 thousand | 47 |
| Refseq_Select | 64 million | 53 |
| Refseq_Protein | 234 million | 55 |
| NR (clustered) | 242 million | 55 |
| NR | 510 million | 56 |

CSHL Programming for Biology

19

19

How many bits do I need?



(More sophisticated algorithms – PSIBLAST, JACKHMMR – also improve sensitivity, and they work better with large databases)

CSHL Programming for Biology

20

20

Effective Similarity Searching

1. Always search protein databases (possibly with translated DNA)
 2. Use E()-values, not percent identity, to infer homology
 - $E() < 0.001$ is significant in a single search
-
1. Search smaller (comprehensive) databases
 2. Change the scoring matrix for:
 - short sequences (exons, reads)
 - short evolutionary distances (mammals, vertebrates, a-proteobacteria)
 - high identity (>50% alignments) to reduce over-extension
 3. Is every aligned residue homologous?
 - alignment overextension

CSHL Programming for Biology

21

21

Scoring matrices – shifting lookback (where do those bits come from?)

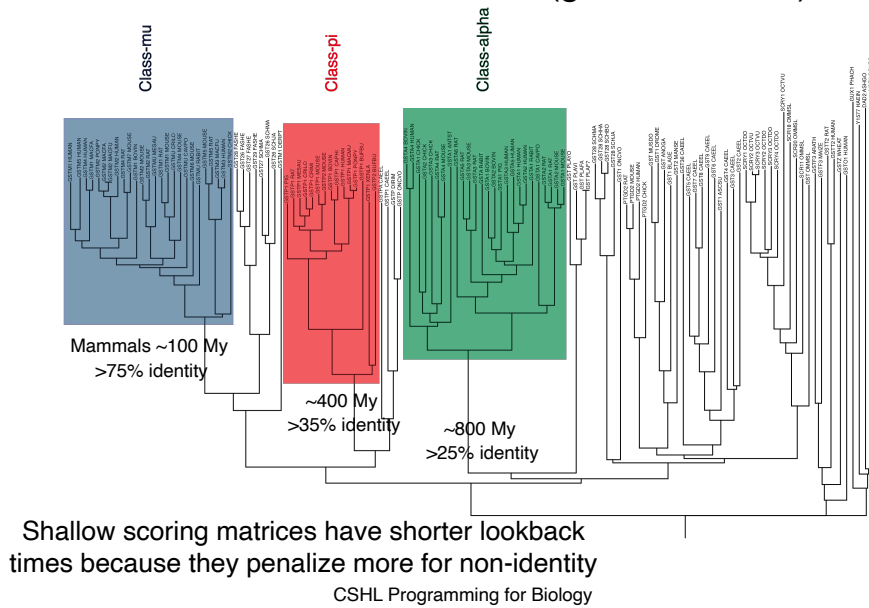
- Scoring matrices can set the evolutionary look-back time for a search
 - Lower PAM (PAM10/VT10 ... PAM/VT40) for closer (10% ... 50% identity)
 - Higher BLOSUM for higher conservation (BLOSUM50 distant, BLOSUM80 conserved)
- Shallow scoring matrices for short domains/short queries (metagenomics)
 - Matrices have “bits/position” (score/position), 40 aa at 0.45 bits/position (BLOSUM62) means 18 bit ave. score (50 bits significant)
- Deep scoring matrices allow alignments to continue, possibly outside the homologous region

CSHL Programming for Biology

22

22

Scoring matrices set look back time: Glutathione Transferases (gstm1_human)



23

23

Scoring matrices and alignment length

Pam40

| | A | R | N | D | E | I | L |
|---|----|-----|----|-----|-----|----|----|
| A | 8 | | | | | | |
| R | -9 | 12 | | | | | |
| N | -4 | -7 | 11 | | | | |
| D | -4 | -13 | 3 | 11 | | | |
| E | -3 | -11 | -2 | 4 | 11 | | |
| I | -6 | -7 | -7 | -10 | -7 | 12 | |
| L | -8 | -11 | -9 | -16 | -12 | -1 | 10 |

Pam250

| | A | R | N | D | E | I | L |
|---|----|----|----|----|----|---|---|
| A | 2 | | | | | | |
| R | -2 | 6 | | | | | |
| N | 0 | 0 | 2 | | | | |
| D | 0 | -1 | 2 | 4 | | | |
| E | 0 | -1 | 1 | 3 | 4 | | |
| I | -1 | -2 | -2 | -2 | -2 | 5 | |
| L | -2 | -3 | -3 | -4 | -3 | 2 | 6 |

$$\lambda S_{i,j} = \log_b \left(\frac{q_{i,j}}{p_i p_j} \right)$$

q_{ij} : homolog frequency w/ Pam40, 250

$$q_{R:N(40)} = 0.000435$$

$$p_R = 0.051$$

$$q_{R:N(250)} = 0.002193$$

$$p_N = 0.043$$

$$\lambda_2 S_{ij} = \lg_2 (q_{ij}/p_i p_j) \quad \lambda_e S_{ij} = \ln(q_{ij}/p_i p_j) \quad p_R p_N = 0.002193$$

$$\lambda_2 S_{R:N(40)} = \lg_2 (0.000435/0.002193) = -2.333$$

$$\lambda_2 = 1/3; S_{R:N(40)} = -2.333/l_2 = -7$$

$$\lambda S_{R:N(250)} = \lg_2 (0.002193/0.002193) = 0$$

CSHL Programming for Biology

24

24

Empirical matrix performance (median results from random alignments)

| Matrix | target % ident | bits/position | aln len (50 bits) |
|------------------|----------------|---------------|-------------------|
| VT160 -12/-2 | 23.8 | 0.26 | 192 |
| BLOSUM50 -10/-2 | 25.3 | 0.23 | 217 |
| BLOSUM62* -11/-1 | 28.9 | 0.45 | 111 |
| VT120 -11/-1 | 27.4 | 1.03 | 48 |
| VT80 -11/-1 | 51.9 | 1.55 | 32 |
| PAM70* -10/-1 | 33.8 | 0.64 | 78 |
| PAM30* -9/-1 | 45.5 | 1.06 | 47 |
| VT40 -12/-1 | 72.7 | 2.76 | 18 |
| VT20 -15/-2 | 84.6 | 3.62 | 13 |
| VT10 /16/-2 | 90.9 | 4.32 | 12 |

HMMs can be very "deep"

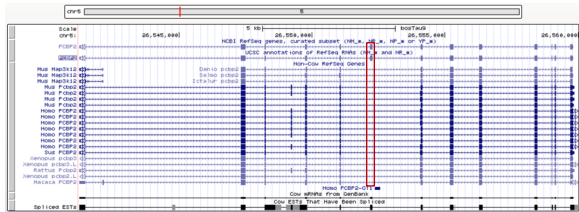
Pearson (2013) Curr. Prot.
Bioinformatics 3.5.1

CSHL Programming for Biology

25

25

Scoring matrices, alignment length, and exon detection – bovine PCBP2



```

human [ 10 20 30 40 50 60 70 80
MDTGVEGLNVLTLIRLLMHGKEVGSIIIGKKGESVKMRRESGARINISEGNCPEIITLAGPTNAIFKAFAMIIDKLE
bovin MDTGVEGLNVLTLIRLLMHGKEVGSIIIGKKGESVKMRRESGARINISEGNCPEIITLAGPTNAIFKAFAMIIDKLE
human [ 10 20 30 40 50 60 70 80
EDISSMTNSTAASRPFTLRLVVPASQCGSLIGKGGCKIKEIRESTGAQVQVAGDMLPSTERAITIAGIPQSIIECVK
bovin EDISSMTNSTAASRPFTLRLVVPASQCGSLIGKGGCKIKEIRESTGAQVQVAGDMLPSTERAITIAGIPQSIIECVK
human [ 90 100 110 120 130 140 150 160
QICVVMLETLSSQSPKGVITIPYRKPSSSPVIFAGGQDRYSTGSDSASFPHITPSMCLNPDLAGPPLAETIQQGYAIPQ
bovin QICVVMLETLSSQSPKGVITIPYRKPSSSPVIFAGGQDRYSTGSDSASFPHITPSMCLNPDLAGPPLAETIQQGYAIPQ
human [ 170 180 190 200 210 220 230 240
PDLTKLHQLAMQSHFPMTHGNTGFSIISSSPEVKGYELDASAQTTSHELTIPNDLIGCIIGRAKINEIROMSGAQ
bovin PDLTKLHQLAMQSHFPMTHGNTGFSIISSSPEVKGYELDASAQTTSHELTIPNDLIGCIIGRAKINEIROMSGAQ
human [ 250 260 270 280 290 300
IKIANPVGSTDRQVITIGSAASISLAQYLINVLRSSETGGMGSS
bovin IKIANPVGSTDRQVITIGSAASISLAQYLINVLRSSETGGMGSS

```

Is this exon really
missing from cow?

CSHL Programming for Biology

26

26

Scoring Matrices - Summary

- PAM and BLOSUM matrices greatly improve the sensitivity of protein sequence comparison – low identity with significant similarity
- PAM matrices have an evolutionary model - lower number, less divergence – lower=closer; higher=more distant
- BLOSUM matrices are sampled from conserved regions at different average identity – higher=more conservation
- Shallow matrices set maximum look-back time
- Short alignments (domains, exons, reads) require shallow (higher information content) matrices

CSHL Programming for Biology

29

29

Effective Similarity Searching

1. Always search protein databases (possibly with DNA – blastx, fastx)
 2. Use E()-values, not percent identity, to infer homology
 - $E() < 0.001$ is significant in a single search
-
1. Search smaller (comprehensive) databases
 - Less redundancy; better sensitivity
 2. Change the scoring matrix for:
 - short evolutionary distances (mammals, vertebrates, a-proteobacteria)
 - short sequences (exons, reads)
 - high identity (>50% alignments) to reduce over-extension

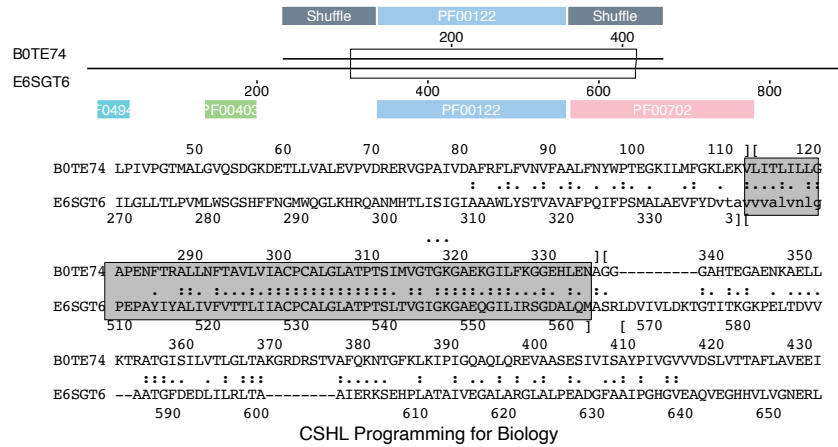
CSHL Programming for Biology

30

30

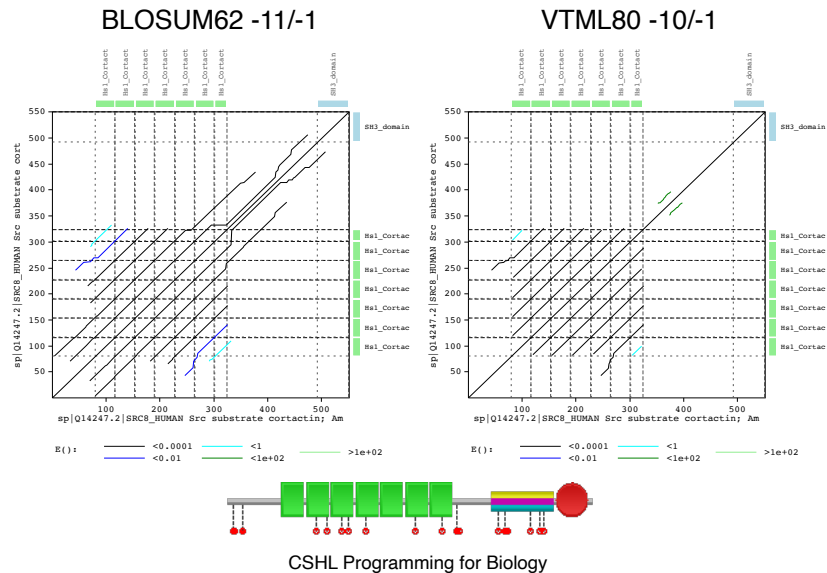
Sub-alignment scoring detects over-extension

```
>>sp|E6SGT6|E6SGT6 THEM7 Heavy metal translocating P-type ATPase EC=3.6.3.4 (888 aa)
qRegion: 81-112:309-340 : score=15; bits=12.3; Id=0.219; Q=0.0 : Shuffle
qRegion: 113-335:341-563 : score=736; bits=232.8; Id=0.641; Q=644.7 : PF00122
qRegion: 336-415:564-642 : score=14; bits=12.0; Id=0.236; Q=0.0 : Shuffle
Region: 81-111:309-339 : score=11; bits=11.1; Id=0.194; Q=0.0 : NODOM :0
Region: 112-334:340-562 : score=736; bits=232.8; Id=0.641; Q=644.7 : PF00122 Pfam
Region: 338-415:566-642 : score=16; bits=12.6; Id=0.244; Q=0.0 : PF00702 Pfam
s-w opt: 632 Z-score: 1048.6 bits: 204.2 E(274545): 3.7e-51
Smith-Waterman score: 765; 49.7% identity (73.3% similar) in 344 aa overlap (81-415:309-642)
```



31

Scoring matrices affect alignment boundaries (homologous over-extension)



32

Homology, non-homology, and over-extension

- Sequences that share statistically significant sequence similarity are homologous (simplest explanation)
- But not all regions of the alignment contribute uniformly to the score
 - lower identity/Q-value because of non-homology (over-extension) ?
 - lower identity/Q-value because more distant relationship (domains have different ages) ?
- Test by searching with isolated region
 - can the distant domain (?) find closer (significant) homologs?
- Similar (homology) or distinct (non-homology) structure is the gold standard
- Multiple sequence alignment can obscure over-extension
 - if the alignment is over-extended, part of the alignment is NOT homologous

CSHL Programming for Biology

33

33

Effective Similarity Searching

1. Always search protein databases (possibly with translated DNA)
 2. Use E()-values, not percent identity, to infer homology
 - $E() < 0.001$ is significant in a single search
-
1. Search smaller (comprehensive) databases
 2. Change the scoring matrix for:
 - short sequences (exons, reads)
 - short evolutionary distances (mammals, vertebrates, a-proteobacteria)
 - high identity (>50% alignments) to reduce over-extension
 3. Is every aligned residue homologous?
 - alignment overextension
 4. (Tomorrow) All methods (pairwise, HMM, PSSM) miss homologs, and find homologs the other methods miss

CSHL Programming for Biology

34

34

workshop II – parsing blast results

Goto:

fasta.bioch.virginia.edu/mol_evol/pfb_python_matrices.html

Your goal is to reproduce a version of this table:

| Matrix | target % ident | align_len | evalue |
|---------------------|-------------------|-----------|--------|
| VT160 | 29.7 | 67 | 2.1 |
| BLOSUM50 | 34.0 | 121 | 1.2 |
| BLOSUM62* -11/-1 | 31.2 | 90 | 0.37 |
| VT80 | 66.7 | 50 | 1.8 |
| VT40 | 72.7 | 11 | 1.3 |

CSHL Programming for Biology

35