

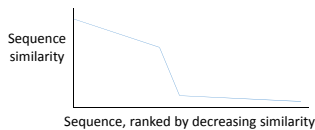
MCB 5472 Lecture #4:
Probabilistic models of homology:
Psi-BLAST and HMMs
February 17, 2014

From last week:

- BLASTp searches find homologs to a single sequence in a sequence database
 - Highest score to sequences best matching the query
 - Corollary: lower scores to distant sequences still matching the query

Finding all homologous sequences using BLASTp

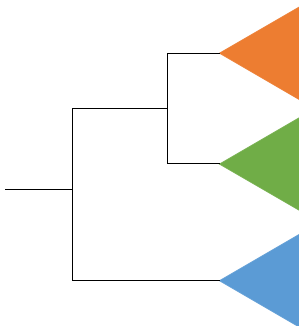
- In an ideal world (e.g., highly conserved sequences) a simple BLASTp search would recover all homologs of a single query sequence
 - Simply accept all sequences above some E-value cutoff
 - E-values can have a steep decline



Nice but...

- Assumes the query sequence perfectly represents all homologs
- False because:
 - Substitutions may be from lineage-specific biases and not conserved in homologs more generally, biasing search towards closer relatives
 - Conservation is always incomplete: the query may not contain conserved positions present in most other homologs
- Sequences close to a hard E-value cutoff can be easily excluded/included depending on search bias

Shown another way



Sensitivity vs specificity

- Trade off between minimizing false-negative detection (sensitivity) and false-positives (specificity)
 - A common trade-off in bioinformatics
- BLASTp is designed to maximize sensitivity
 - Specificity can be low – left to the user to cull out the false-positives

Net result:

- Divergent homologs are hard to detect
- Because they are close to typical E-value cutoffs, any bias can easily lead to them being excluded
- Discriminating between false- and true-positives can be problematic
 - Requires manual examination
- Finding deep homology is hard!

Better: use multiple queries representing all homologs

- Option 1: Run multiple individual BLASTps
 - Still easy to bias – “unknown unknowns”
- Option 2: Make a statistical model of sequence conservation amongst all homologs and use that to find different relatives
 - Diverse input sequences removes and averages out lineage-specific biases

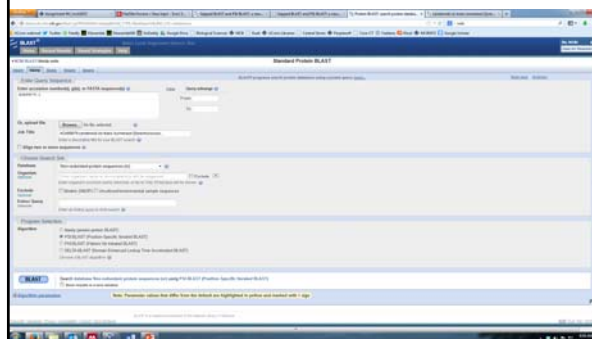
PSI-BLAST

- “Position-Specific Iterated BLAST”
- Works **only** for proteins
- Uses BLASTp to create a “position-specific score matrix” (PSSM)
 - Smith Waterman global alignments are an option for the command line but not the web (slower, more accurate)
- Uses matrix for subsequent database searches
- Matrix updated on each iteration
 - Bias reduced each time
 - Sensitivity increased towards distant homologs
 - False-positives reduced by model refinement

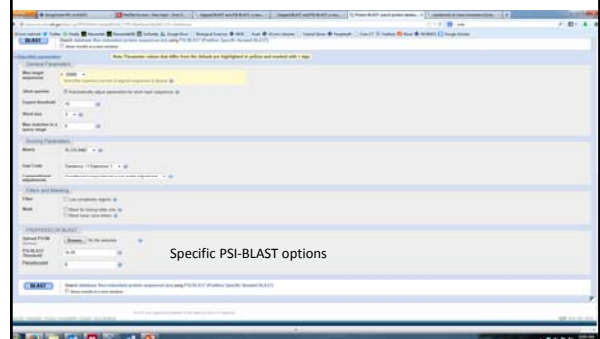
PSI-BLAST: step #1

- First iteration: standard BLASTp using a single sequence
- All homologs above a specified E-value threshold kept to make PSSM
 - Can be specified via parameters, manually edited on NCBI website implementation

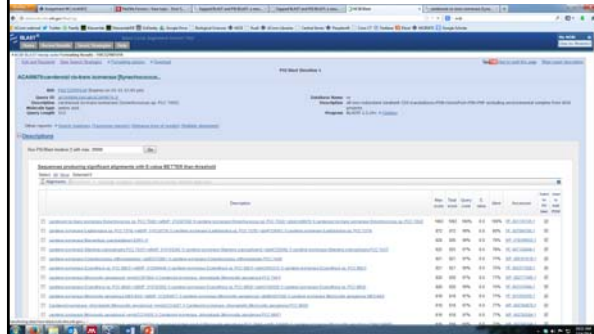
NCBI web implimentation



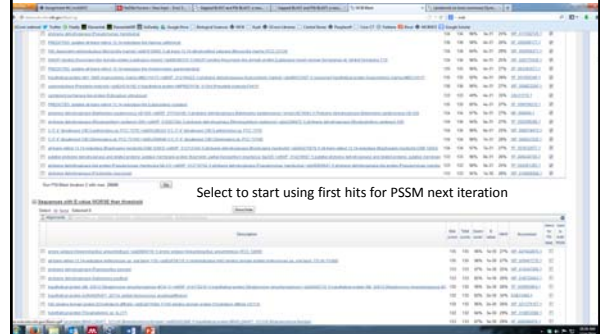
NCBI web implementation



NCBI web implementation



NCBI web implementation



PSI-BLAST: step #2

- Makes (rough) multiple sequence alignment for the selected BLASTp results
- All hits aligned to the query
 - Not a true multiple sequence alignment
 - Possible to input an externally generated alignment via terminal version (but not web)
- Alternative at terminal: Smith Waterman global pairwise alignments
 - Not available for web
 - Slower but more accurate

PSI-BLAST: step #3

- Use sequence alignment to create Position-Specific Scoring Matrix (PSSM)
- PSSM:
 - Unique substitution matrix for each sequence alignment column
 - Extra column for gap penalty
 - Matrix is 21 x [query length] vs. 20 x 20 for normal matrix
 - Scores merge standard distance matrix with position-specific frequencies from 1st iteration, weighted by sequence similarity

Example PSSM

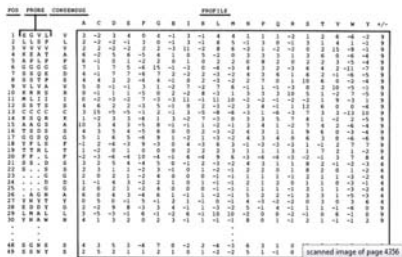


Fig. 1. The concept of a profile. Left: A few degrees of profile entropy. Right: An example of a profile for the immunoglobulin variable region domain, generated from the four protein sequences shown at the left (see Fig. 2B for details). The profile is shown in the box. The rightmost column of the profile gives the penalty for insertions (between 1-17). Positions 18-47 of the profile are omitted from the figure for clarity. Notice that where gaps appear in some of the profile sequences, the insertion/indletion penalty is lower than elsewhere.

Gribkov et al. (1987) PNAS 84: 4355-4358

PSI-BLAST: step #4

- Query reference database using the PSSM
 - Recall: BLASTp looks for 2-3 amino acid words similar to the query sequence above some threshold score calculated from the distance matrix
 - An equivalent calculation can be performed using the PSSM; find possible words having a score > the same threshold
 - Subsequent BLAST steps are the same: extend matching words, recalculate with gaps, calculate statistics
 - E-values now reflect similarity to the query profile, not any individual sequence

PSI-BLAST: step #5

- Perform as many iterations as you like
- PSSM updated each time based on hits passing E-value threshold on the previous iteration
- Sequence-specific bias reduced each time as the PSSM is adjusted to reflect homolog in the entire input set

CrtR protein 1e-50 threshold

Iteration	Hits > 1e-50	Notes
1	151	
2	215	
3	258	Query not top hit, top E-value != 0.0
4	271	
5	271	
6	271	

Model corruption

- If a non-homologous sequence is included during model construction, can bias the model away from true homologs
- With subsequent iteration, model can be made completely useless
- Using a higher E-value cutoff can ameliorate
- On web can examine results and limit selection
 - Can't do this in high-throughput at terminal

Command line PSI-BLAST

- Part of BLAST+ package so same basic parameters apply
 - Additional flags:


```
-num_iterations [number]
-out_pssm [filename]
-out_ascii_pssm [filename]
-comp_based_stats 0 # required
```
- e.g., [jlklassen@bbcsrv3 ~]\$ psiblast -query test.faa -db all.faa -num_iterations 3 -out test_vs_all.psiblast -out_ascii_pssm pssm.out -comp_based_stats 0

Starting PSI-BLAST with pre-computed PSSM

- Create PSSM using PSI-BLAST with


```
-out_pssm flag (not -out_ascii_psm)
```
- Use `-in_pssm` flag instead of `-query`
- e.g., [jlklassen@bbcsrv3 ~]\$ psiblast


```
-in_pssm pssm.out -db all.faa
-num_iterations 3 -out
test_vs_all.psiblast
-comp_based_stats 0
```

RPSBLAST

- PSI-BLAST queries a sequence database with an individual PSSM
- RPSBLAST does the opposite: queries an individual sequence with a database of PSSMs
 - e.g., From NCBI's Conserved Domain Database (CDD) to annotate sequences according to NCBI's ortholog family description
- In BLAST+, command is `rpsblast+` and works similarly to other BLAST+ commands except `-db` is now PSSMs, not sequences
- Possible to make your own PSSM database but complicated (most people use HMMer instead)

Beyond BLAST

- Recall that all BLAST programs are local alignments
 - Trade-off between speed and accuracy
- FASTA: alternative package for database
 - <http://www.ebi.ac.uk/Tools/sss/fasta/>
 - Heuristics like BLAST using word matching for initial sequence matching
 - Final alignments use Smith-Waterman global pairwise alignment method
- Advances in computer science and statistics on both fronts (i.e., better accuracy & better approximations)

HMMER

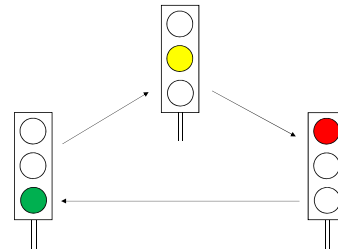
- HMMER is a software package similar to PSI-BLAST, i.e., searching databases with homology models
- Uses HMMs instead of PSSMs
- Advantages:
 - More statistically explicit models
 - HMMER3 as ~fast as BLAST
 - Easy to use at command line
 - Can make models for DNA, RNA, protein
- Disadvantages:
 - Initial alignment is always a second step
 - No NCBI interface (database-specific instead)

The purpose of HMMs

- To evaluate the probability of a sequence matching a model
 - Assumes preexisting model
- Essentially a classification problem
 - Given data, how well does it fit a model?
 - Given data and multiple models, which fits best?
 - e.g., Does a gene belong to a gene family?

Markov Chains

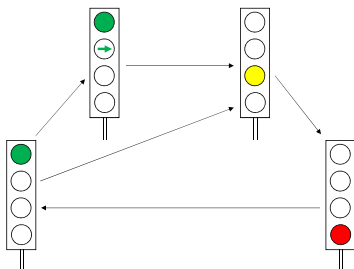
- “A finite number of states connected by transitions”



http://www.ch.emblnet.org/CourseMBnet/Base03/slides/PSSM_HMM.pdf

Markov Chains

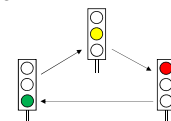
- “A finite number of states connected by transitions”



http://www.ch.emblnet.org/CourseMBnet/Base03/slides/PSSM_HMM.pdf

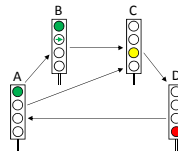
Transition probabilities

- The probability of moving from one state to another
- $P(\text{Yellow}|\text{Green}) = 1$
 - “The probability of transitioning to yellow given green”
- $P(\text{Red}|\text{Yellow}) = 1$
- $P(\text{Green}|\text{Red}) = 1$
- $P(\text{Green}|\text{Yellow}) = 0$
- $P(\text{Yellow}|\text{Red}) = 0$
- $P(\text{Red}|\text{Green}) = 0$



Transmission probabilities

- The probability of moving from one state to another
- $P(C|B) = 1$
- $P(D|C) = 1$
- $P(A|D) = 1$
- $P(A|B) = P(\text{vehicles turning left})$
- $P(A|B) = P(\text{no vehicles turning left})$
- $P(B|A) = 0$ etc.



Hidden Markov Models

- HMMs are like Markov Chains in that they comprise states connected by transitions
- Difference: each state does not comprise a single symbol but rather a distribution of them
 - e.g., a column of a sequence alignment will contain some frequency of A, C, G and T
- Each state can “emit” a symbol with some probability

Hidden Markov Models

- Known:
 - The number of states
 - The transition probabilities
 - The emission probabilities
- Question: how well does a sequence match the model?
 - Evaluate the global probability by multiplying the probability of each step through the graph

A simplified model: identifying a 5' splice site

Eddy 2004 Nat. Biotechnol. 22: 1315-1316

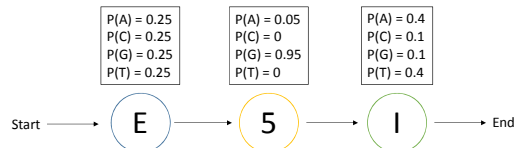
Three possible states: exon, 5' splice site, intron



CTTCATGTGAAAGCAGACGTAAGTCA

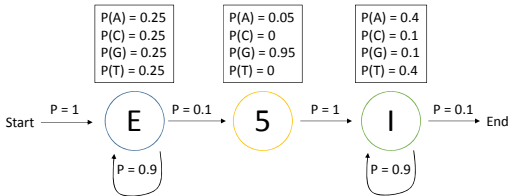
Emission probabilities

- 5' splice site nearly always G, occasionally A
- Exon sequence distributed uniformly
- Intron sequence is AT rich

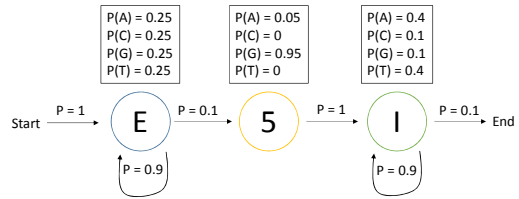


Transition probabilities

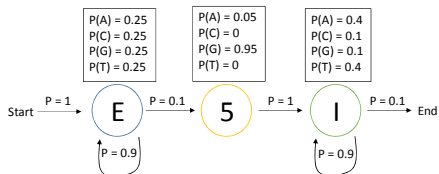
- Exon & intron can be multiple bases long
- 5' splice site only one base long
- Exact probabilities are flexible



Discuss: given this sequence, where is the splice site?

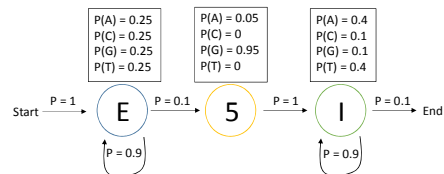


CTTTCATGTGAAAAGCAGACGTAAGTCA

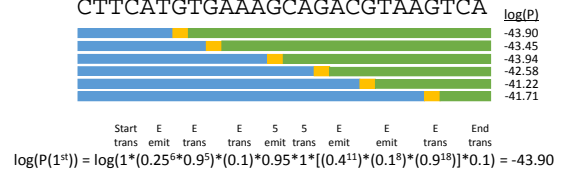


CTTTCATGTGAAAAGCAGACGTAAGTCA

- 14 possibilities (number of "A"s and "G"s)
- "A"s are sufficiently improbable that we will not work them out here



CTTTCATGTGAAAAGCAGACGTAAGTCA



HMMs

- Given a model and input data, we can calculate the likelihood of any given classification
- Because model is fully parameterized, significance of each path can be determined in a Bayesian statistical framework
 - "Posterior decoding"

Posterior decoding

- Definition: probability of chosen path divided by sum of the probability of all other paths
- e.g., $-41.22 / ((-43.90) + (-43.45) + (-43.94) + \dots)$

