

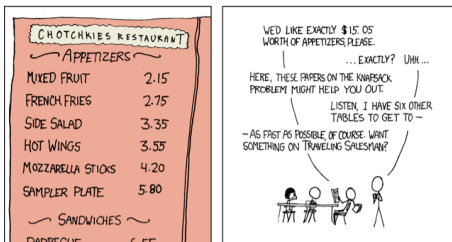
Heuristic Approaches

Mark Voorhies

5/5/2017

MY HOBBY:

EMBEDDING NP-COMPLETE PROBLEMS IN RESTAURANT ORDERS



[HTTP://WWW.XKCD.COM/287/](http://www.xkcd.com/287/)

PAM (Dayhoff) and BLOSUM matrices

- PAM1 matrix originally calculated from manual alignments of highly conserved sequences (myoglobin, cytochrome C, etc.)

PAM (Dayhoff) and BLOSUM matrices

- PAM1 matrix originally calculated from manual alignments of highly conserved sequences (myoglobin, cytochrome C, etc.)
- We can think of a PAM matrix as evolving a sequence by one unit of time.

PAM (Dayhoff) and BLOSUM matrices

- PAM1 matrix originally calculated from manual alignments of highly conserved sequences (myoglobin, cytochrome C, etc.)
- We can think of a PAM matrix as evolving a sequence by one unit of time.
- If evolution is uniform over time, then PAM matrices for larger evolutionary steps can be generated by multiplying PAM1 by itself (so, higher numbered PAM matrices represent greater evolutionary distances).

PAM (Dayhoff) and BLOSUM matrices

- PAM1 matrix originally calculated from manual alignments of highly conserved sequences (myoglobin, cytochrome C, etc.)
- We can think of a PAM matrix as evolving a sequence by one unit of time.
- If evolution is uniform over time, then PAM matrices for larger evolutionary steps can be generated by multiplying PAM1 by itself (so, higher numbered PAM matrices represent greater evolutionary distances).
- The BLOSUM matrices were determined from automatically generated ungapped alignments. Higher numbered BLOSUM matrices correspond to *smaller* evolutionary distances. BLOSUM62 is the default matrix for BLAST.

Motivation for scoring matrices

Frequency of residue i :

p_i

Motivation for scoring matrices

Frequency of residue i :

$$p_i$$

Frequency of residue i aligned to residue j :

$$q_{ij}$$

Motivation for scoring matrices

Frequency of residue i :

$$p_i$$

Frequency of residue i aligned to residue j :

$$q_{ij}$$

Expected frequency if i and j are independent:

$$p_i p_j$$

Motivation for scoring matrices

Frequency of residue i :

$$p_i$$

Frequency of residue i aligned to residue j :

$$q_{ij}$$

Expected frequency if i and j are independent:

$$p_i p_j$$

Ratio of observed to expected frequency:

$$\frac{q_{ij}}{p_i p_j}$$

Motivation for scoring matrices

Frequency of residue i :

$$p_i$$

Frequency of residue i aligned to residue j :

$$q_{ij}$$

Expected frequency if i and j are independent:

$$p_i p_j$$

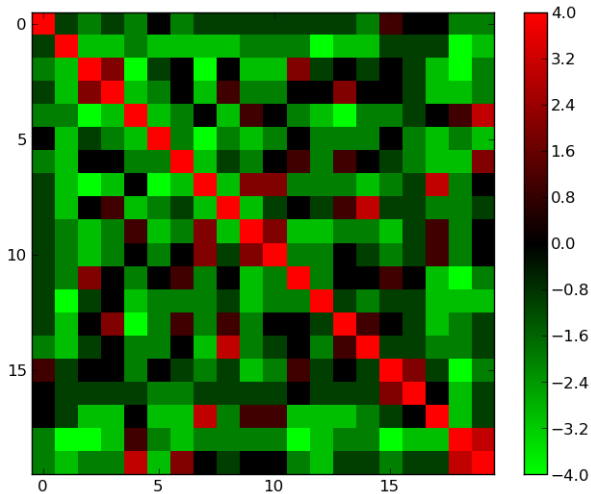
Ratio of observed to expected frequency:

$$\frac{q_{ij}}{p_i p_j}$$

Log odds (LOD) score:

$$s(i, j) = \log \frac{q_{ij}}{p_i p_j}$$

BLOSUM45 in alphabetical order



Clustering amino acids on log odds scores

```
import networkx as nx
try:
    import Pycluster
except ImportError:
    import Bio.Cluster as Pycluster

class ScoreCluster:
    def __init__(self, S, alpha_aa = "ACDEFGHIKLMNPQRSTVWY"):
        """ Initialize from numpy array of scaled log odds scores. """
        (x,y) = S.shape
        assert(x == y == len(alpha_aa))

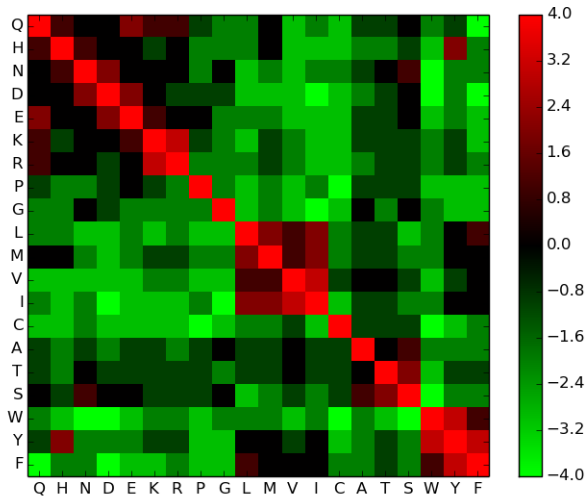
        # Interpret the largest score as a distance of zero
        D = max(S.reshape(x**2))-S
        # Maximum-linkage clustering, with a user-supplied distance matrix
        tree = Pycluster.treecluster(distancematrix = D, method = "m")

        # Use NetworkX to read out the amino-acids in clustered order
        G = nx.DiGraph()
        for (n,i) in enumerate(tree):
            for j in (i.left, i.right):
                G.add_edge(-(n+1),j)

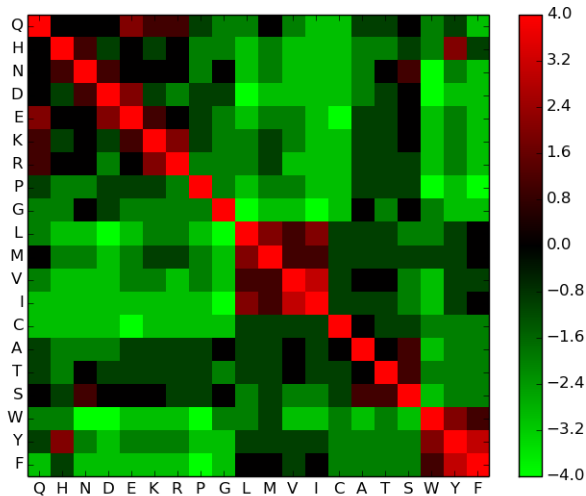
        self.ordering = [i for i in nx.dfs_preorder(G, -len(tree)) if(i >= 0)]
        self.names = "".join(alpha_aa[i] for i in self.ordering)
        self.C = self.permute(S)

    def permute(self, S):
        """ Given square matrix S in alphabetical order, return rows and columns
        of S permuted to match the clustered order. """
        return array([[S[i][j] for j in self.ordering] for i in self.ordering])
```

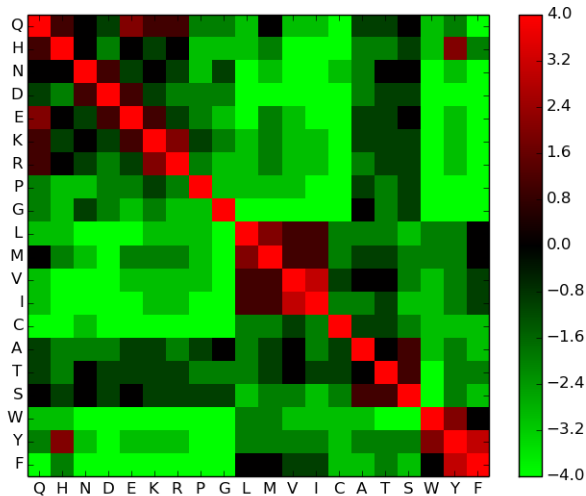
BLOSUM45 – maximum linkage clustering



BLOSUM62 with BLOSUM45 ordering



BLOSUM80 with BLOSUM45 ordering



The implementation of local alignment is the same as for global alignment, with a few changes to the rules:

- Initialize edges to 0 (*no penalty for starting in the middle of a sequence*)
- The maximum score is never less than 0, and no pointer is recorded unless the score is greater than 0 (*note that this implies negative scores for gaps and bad matches*)
- The trace-back starts from the highest score in the matrix and ends at a score of 0 (*local, rather than global, alignment*)

Because the naive implementation is essentially the same, the time and space requirements are also the same.

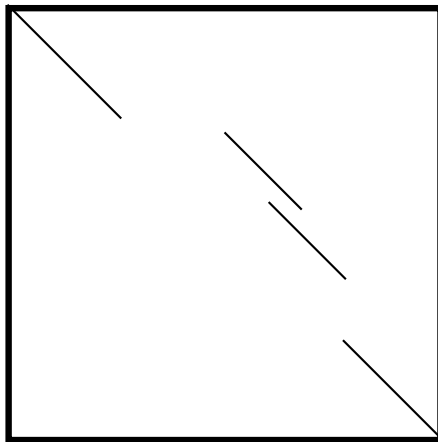
Smith-Waterman

	A	G	C	G	G	T	A	
	0	0	0	0	0	0	0	
G	0	0	1	0	0	1	0	
A	0	1	0	0	0	0	1	
G	0	0	2	1	1	1	0	
C	0	0	1	3	2	1	0	
G	0	0	0	2	4	3	2	
G	0	0	1	1	3	5	4	
A	0	1	0	0	2	4	4	5

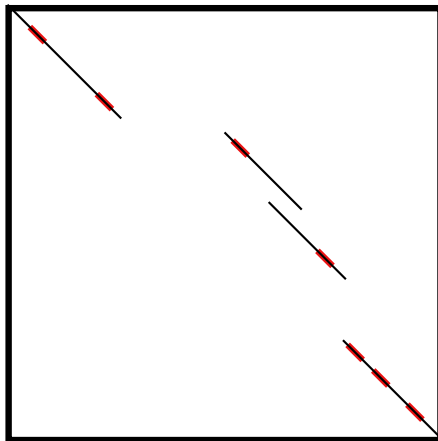
Why BLAST?

- Fast, heuristic approximation to a full Smith-Waterman local alignment
- Developed with a statistical framework to calculate expected number of false positive hits.
- Heuristics biased towards “biologically relevant” hits.

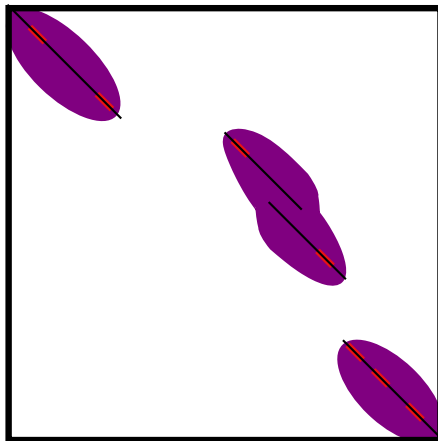
BLAST: A quick overview



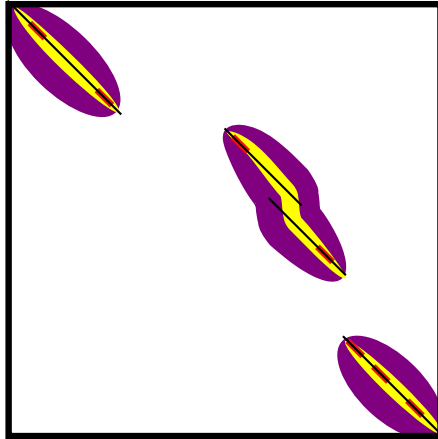
BLAST: Seed from exact word hits



BLAST: Myers and Miller local alignment around seed pairs



BLAST: High Scoring Pairs (HSPs)



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

- E : Expected number of “random” hits in a database of this size scoring *at least* S .
- S : HSP score
- m : Query length
- n : Database size
- k : Correction for similar, overlapping hits
- λ : normalization factor for scoring matrix

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

- E : Expected number of “random” hits in a database of this size scoring *at least* S .
- S : HSP score
- m : Query length
- n : Database size
- k : Correction for similar, overlapping hits
- λ : normalization factor for scoring matrix

A variant of this formula is used to generate sum probabilities for combined HSPs.

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

- E : Expected number of “random” hits in a database of this size scoring *at least* S .
- S : HSP score
- m : Query length
- n : Database size
- k : Correction for similar, overlapping hits
- λ : normalization factor for scoring matrix

A variant of this formula is used to generate sum probabilities for combined HSPs.

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

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

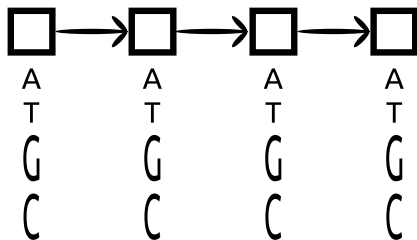
- E : Expected number of “random” hits in a database of this size scoring *at least* S .
- S : HSP score
- m : Query length
- n : Database size
- k : Correction for similar, overlapping hits
- λ : normalization factor for scoring matrix

A variant of this formula is used to generate sum probabilities for combined HSPs.

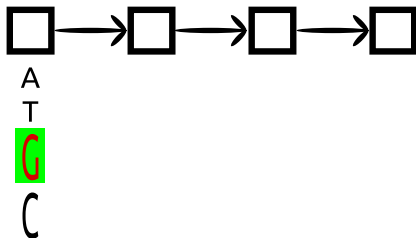
$$p = 1 - e^{-E}$$

(If you care about the difference between E and p , you're already in trouble)

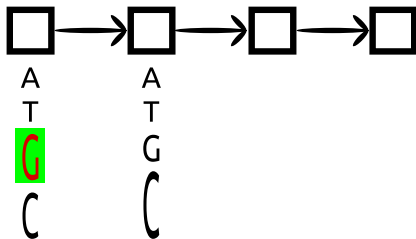
0th order Markov Model



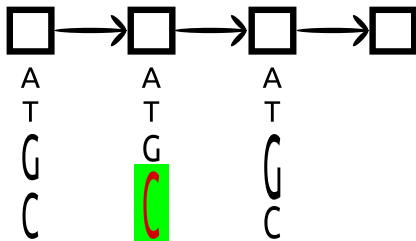
1st order Markov Model



1st order Markov Model



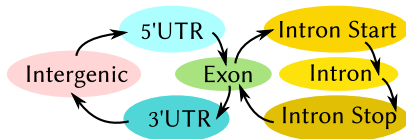
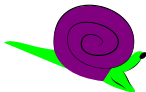
1st order Markov Model



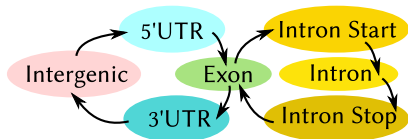
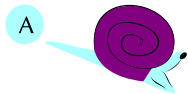
What are Markov Models good for?

- Background sequence composition
- Spam

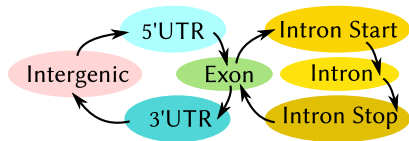
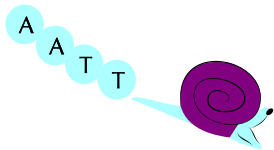
Hidden Markov Models



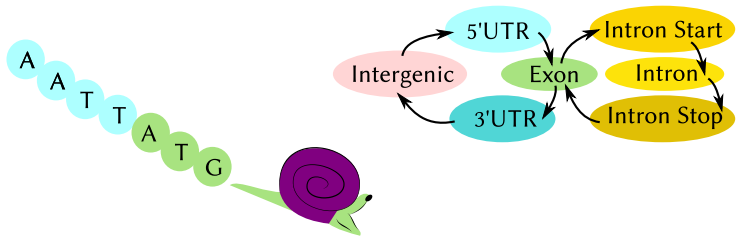
Hidden Markov Models



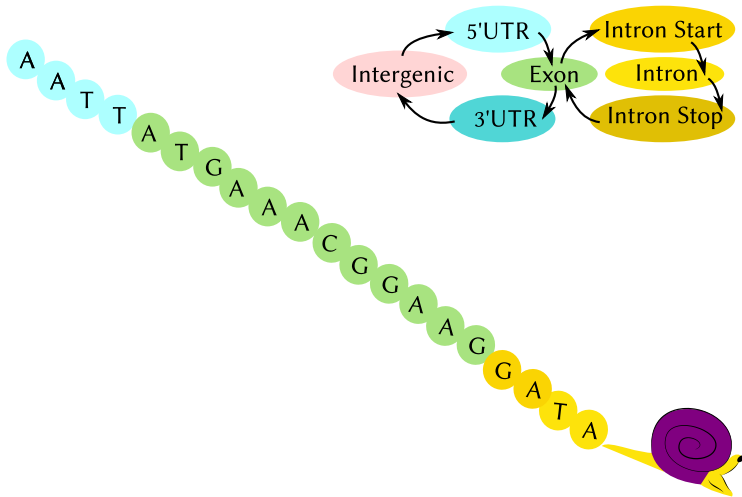
Hidden Markov Models



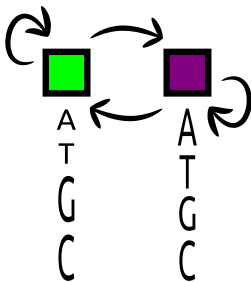
Hidden Markov Models



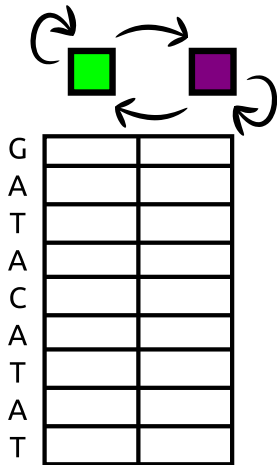
Hidden Markov Models



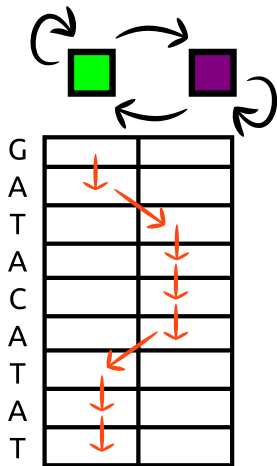
Hidden Markov Model



The Viterbi algorithm: Alignment

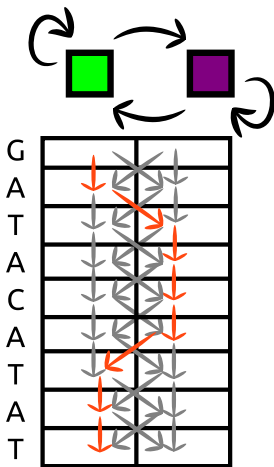


The Viterbi algorithm: Alignment



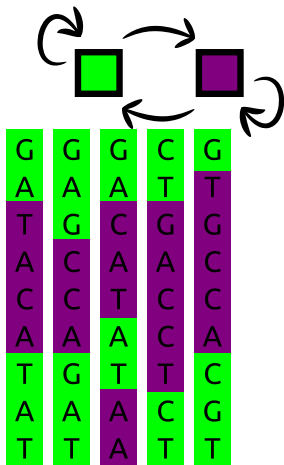
- Dynamic programming, like Smith-Waterman
- Sums *best* log probabilities of emissions and transitions (*i.e.*, multiplying independent probabilities)
- Result is most likely annotation of the target with hidden states

The Forward algorithm: Net probability



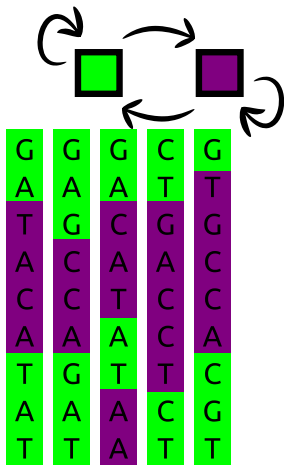
- Probability-weighted sum over all possible paths
- Simple modification of Viterbi (although *summing* probabilities means we have to be more careful about rounding error)
- Result is the probability that the observed sequence is explained by the model
- In practice, this probability is compared to that of a null model (e.g., random genomic sequence)

Training an HMM



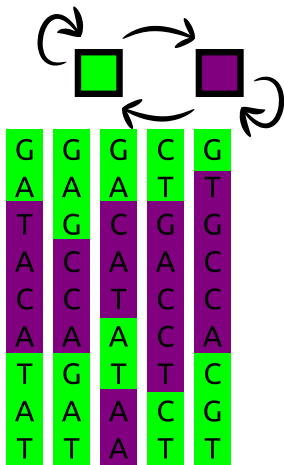
- If we have a set of sequences with known hidden states (e.g., from experiment), then we can calculate the emission and transition probabilities directly

Training an HMM



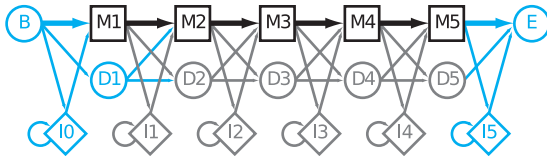
- If we have a set of sequences with known hidden states (e.g., from experiment), then we can calculate the emission and transition probabilities directly
- Otherwise, they can be iteratively fit to a set of unlabeled sequences that are known to be true matches to the model

Training an HMM



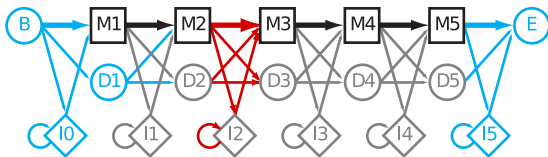
- If we have a set of sequences with known hidden states (e.g., from experiment), then we can calculate the emission and transition probabilities directly
- Otherwise, they can be iteratively fit to a set of unlabeled sequences that are known to be true matches to the model
- The most common fitting procedure is the Baum-Welch algorithm, a special case of expectation maximization (EM)

Profile Alignments: Plan 7



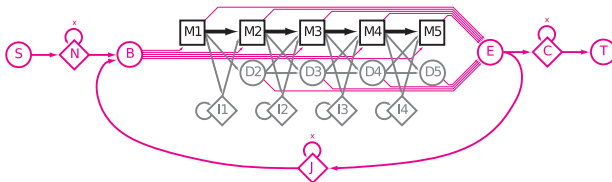
(Image from Sean Eddy, PLoS Comp. Biol. 4:e1000069)

Profile Alignments: Plan 7 (from Outer Space)



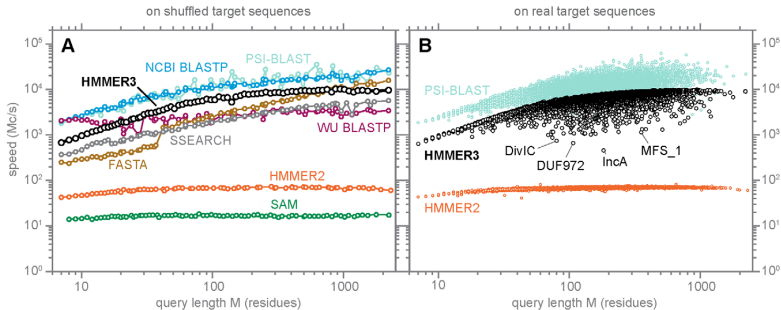
(Image from Sean Eddy, PLoS Comp. Biol. 4:e1000069)

Rigging Plan 7 for Multi-Hit Alignment



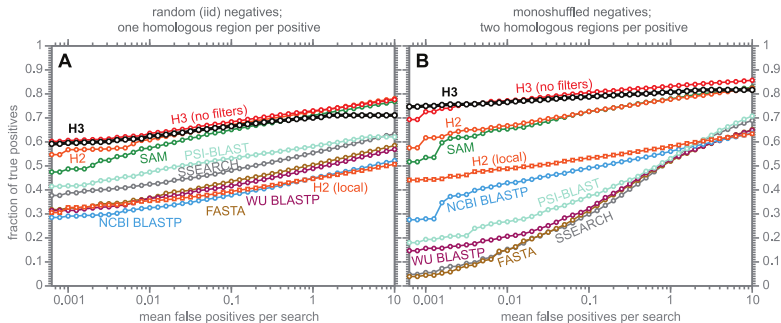
(Image from Sean Eddy, PLoS Comp. Biol. 4:e1000069)

HMMer3 speeds



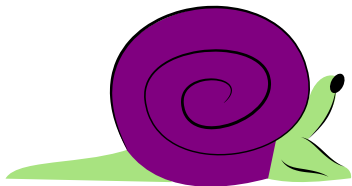
Eddy, PLoS Comp. Biol. 7:e1002195

HMMer3 sensitivity and specificity



Eddy, PLoS Comp. Biol. 7:e1002195

Stochastic Context Free Grammars



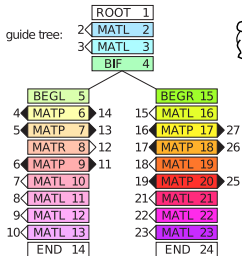
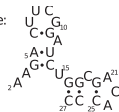
- Can emit from both sides \rightarrow base pairs
- Can duplicate emitter \rightarrow bifurcations

INFERNAL/Rfam

input multiple alignment:

[structure]	. . . <<<	> - >>	: << - >	. . . >>>	.		
human	. AAGACUUCG	GAUCUGGCG	. ĀĀĀ .	CCC .			
mouse	a UACACUUCG	GAUG - CACC	. AAA .	GUG a			
orc	. AGGUCUUC -	GCACGGGCA	gCCA c	UUC .			
	1	5	10	15	20	25	28

example structure:

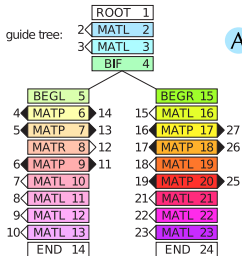
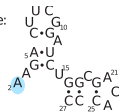


INFERNAL/Rfam

input multiple alignment:

[structure]	. . . <<<	> - >>	<< - <	. . . >>>	.		
human	. AAGACUUCG	GAUCUGGCG	. ĀĀĀ	. CCC	.		
mouse	a UACACUUCG	AUG - CACC	. AAA	. GUG	a		
orc	. AGGUCUUC -	GCACGGGCA	gCCA	cUUC	.		
	1	5	10	15	20	25	28

example structure:

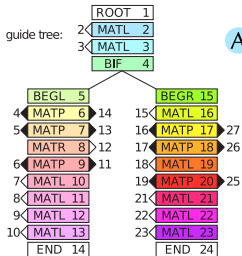
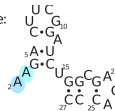


INFERNAL/Rfam

input multiple alignment:

[structure]	. . . <<<<	>>>>	<<<<	>>>>	. . . >>>>	.	
human	. AAGACUUCGG	GAUCUGGCG	. ĀĀĀ .	CCC .			
mouse	a UACACUUCGGAUG	-CACC	. AAA .	GUG a			
orc	. AGGUCUUC -	GCACGGGCA	g CCA c	UUC .			
	1	5	10	15	20	25	28

example structure:

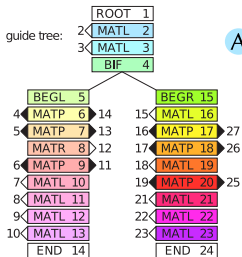


INFERNAL/Rfam

input multiple alignment:

[structure]	. . . <<< <<< >> >> >> >> >> . . . >>> .
human	. AAGACUUCGGGAUCUGGCG . ĀĀĀ . CCC .
mouse	a UACACUUCGGAUG - CACC . AAA . GUG a
orc	. AGGUCUUC - GCACGGGCA g CCA c UUC .
	1 5 10 15 20 25 28

example structure:

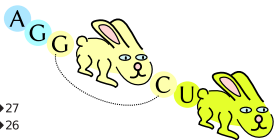
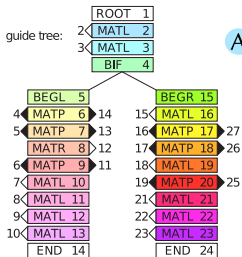
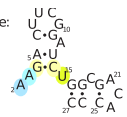


INFERNAL/Rfam

input multiple alignment:

[structure]	. . . <<<<	>>>>	<<<<	>>>>	. . . >>>>	.	
human	. AAGACUUCG	GAUCUGGCG	. ĀĀĀ .	CCC .			
mouse	aUACACUUCG	AUG - CACC	. AAA .	GUG a			
orc	. AGGUCUUC -	GCACGGGCA	gCCA c	UUC .			
	1	5	10	15	20	25	28

example structure:

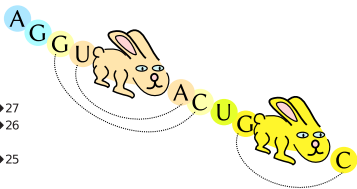
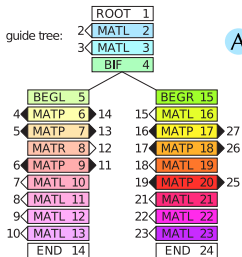
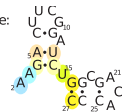


INFERNAL/Rfam

input multiple alignment:

[structure]	. . . <<<<	>> >>	<<<<	>>>>	. . . >>>>	.	
human	. AAGACUUCG	GAUCUGGCG	. ĀĀĀ .	CCC .			
mouse	a UACACUUCG	AUG - CACC	. AAA .	GUG a			
orc	. AGGUCUUC -	GCACGGGCA	g CCA c	UUC .			
	1	5	10	15	20	25	28

example structure:

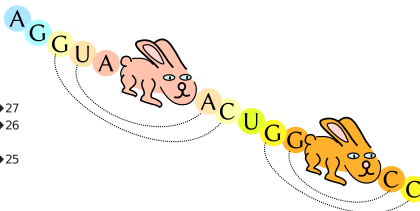
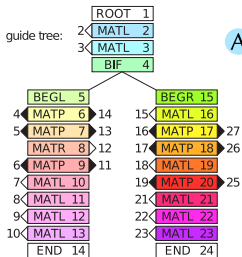
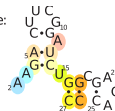


INFERNAL/Rfam

input multiple alignment:

[structure]	. . . <<<<	>>>>	<<<<	>>>>	. . . >>>>	.	
human	. AAGACUUCG	GAUCUGGCG	. ĀĀĀ .	CCCC .			
mouse	aUACACUUCG	GAUG-CACC	. AAA .	GUG a			
orc	. AGGUCUUC-	GCACGGGCA	gCCA c	UUC .			
	1	5	10	15	20	25	28

example structure:



- Keep working on your dynamic programming code.