

RNA Search and Motif Discovery

Genome 54I
Intro to Computational
Molecular Biology

Many biologically interesting roles for RNA
RNA secondary structure prediction

Many interesting RNAs, e.g. Riboswitches

Approaches to Structure Prediction

- Maximum Pairing**
 - + works on single sequences
 - + simple
 - too inaccurate
- Minimum Energy**
 - + works on single sequences
 - ignores pseudoknots
 - only finds "optimal" fold
- Partition Function**
 - + finds all folds
 - ignores pseudoknots

Nussinov's Algorithm

Computation Order

Or energy

$B(i,j) = \# \text{ pairs in optimal pairing of } r_i \dots r_j$

$B(i,j) = 0$ for all i, j with $i \geq j-4$; otherwise

$B(i,j) = \max$ of:

$$\begin{cases} B(i,j-1) \\ \max \{ B(i,k-1) + 1 + B(k+1,j-1) \mid i \leq k < j-4 \text{ and } r_k-r_j \text{ may pair} \} \end{cases}$$

Time: $O(n^3)$

Approaches, II

- Comparative sequence analysis**
 - + handles all pairings (potentially incl. pseudoknots)
 - requires several (many?) aligned, appropriately diverged sequences
- Stochastic Context-free Grammars**
 - Roughly combines min energy & comparative, but no pseudoknots
- Physical experiments (x-ray crystallography, NMR)**

Day 2

Day 1:

Many biologically interesting roles for RNA
RNA secondary structure prediction

Today:

Covariance Models (CMs) represent
RNA sequence/structure motifs
Fast CM search

8

Motif Description

What

A probabilistic model for RNA families

The “Covariance Model”

≈ A Stochastic Context-Free Grammar

A generalization of a profile HMM

Algorithms for Training

From aligned or unaligned sequences

Automates “comparative analysis”

Complements Nussinov/Zucker RNA folding

Algorithms for searching

16

Computational Problems

~~How to predict secondary structure~~

How to model an RNA “motif”
(i.e., sequence/structure pattern)

Given a motif, how to search for instances

Given (unaligned) sequences, find motifs

How to score discovered motifs

How to leverage prior knowledge

9

RNA Motif Models

“Covariance Models” (Eddy & Durbin 1994)

aka profile stochastic context-free grammars

aka hidden Markov models on steroids

Model position-specific nucleotide
preferences *and* base-pair preferences

Pro: accurate

Con: model building hard, search slow

15

Main Results

Very accurate search for tRNA

(Precursor to tRNAscanSE - current favorite)

Given sufficient data, model construction
comparable to, but not quite as good as,
human experts

Some quantitative info on importance of
pseudoknots and other tertiary features

17

Probabilistic Model Search

As with HMMs, given a sequence, you calculate likelihood ratio that the model could generate the sequence, vs a background model

You set a score threshold

Anything above threshold → a “hit”

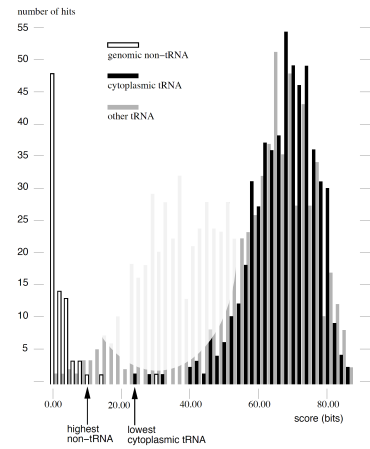
Scoring:

“Forward” / “Inside” algorithm - sum over all paths

Viterbi approximation - find single best path

(Bonus: alignment & structure prediction)

Example: searching for tRNAs

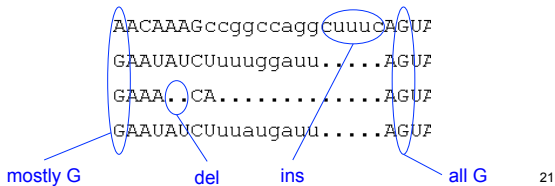


How to model an RNA “Motif”?

Conceptually, start with a profile HMM:

from a multiple alignment, estimate nucleotide/ insert/delete preferences for each position

given a new seq, estimate likelihood that it could be generated by the model, & align it to the model



How to model an RNA “Motif”?

Add “column pairs” and pair emission probabilities for base-paired regions



Profile HMM Structure

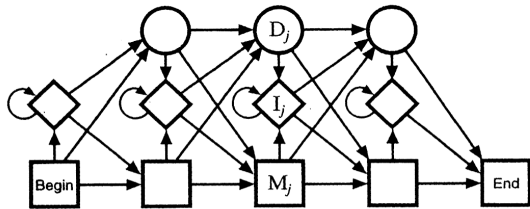


Figure 5.2 The transition structure of a profile HMM.

- Mj: Match states (20 emission probabilities)
- Ij: Insert states (Background emission probabilities)
- Dj: Delete states (silent - no emission)

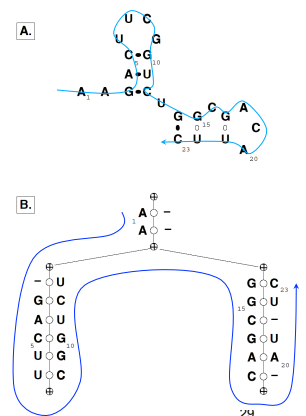
CM Structure

A: Sequence + structure

B: the CM “guide tree”

C: probabilities of letters/ pairs & of indels

Think of each branch being an HMM emitting both sides of a helix (but 3’ side emitted in reverse order)

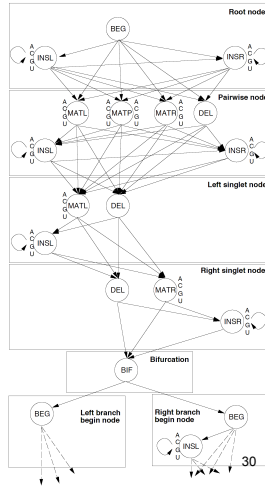


Overall CM Architecture

One box ("node") per node of guide tree

BEG/MATL/INS/DEL just like an HMM

MATP & BIF are the key additions: MATP emits pairs of symbols, modeling base-pairs; BIF allows multiple helices



CM Viterbi Alignment (the "inside" algorithm)

x_i = i^{th} letter of input

x_{ij} = substring i, \dots, j of input

T_{yz} = $P(\text{transition } y \rightarrow z)$

E_{x_i, x_j}^y = $P(\text{emission of } x_i, x_j \text{ from state } y)$

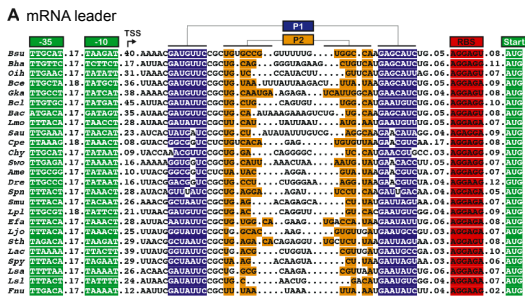
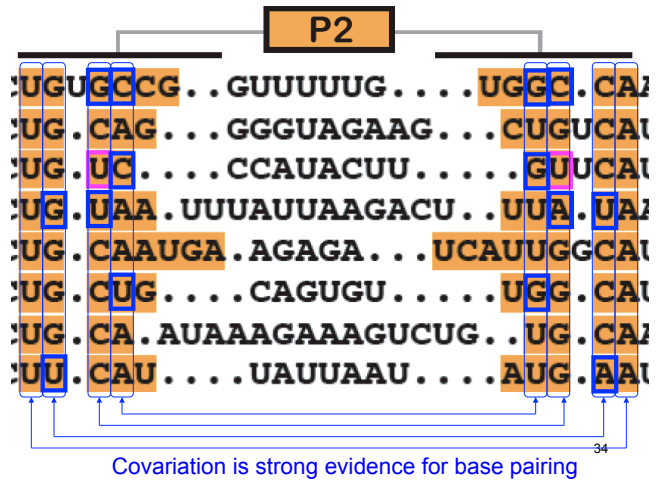
$S_{ij}^y = \max_{\pi} \log P(x_{ij} \text{ gen'd starting in state } y \text{ via path } \pi)$

CM Viterbi Alignment (the "inside" algorithm)

$S_{ij}^y = \max_{\pi} \log P(x_{ij} \text{ generated starting in state } y \text{ via path } \pi)$

$$S_{ij}^y = \begin{cases} \max_z [S_{i+1, j-1}^z + \log T_{yz} + \log E_{x_i, x_j}^y] & \text{match pair} \\ \max_z [S_{i+1, j}^z + \log T_{yz} + \log E_{x_i}^y] & \text{match/insert left} \\ \max_z [S_{i, j-1}^z + \log T_{yz} + \log E_{x_j}^y] & \text{match/insert right} \\ \max_z [S_{i, j}^z + \log T_{yz}] & \text{delete} \\ \max_{l < k \leq j} [S_{i, k}^{y_{left}} + S_{k+1, j}^{y_{right}}] & \text{bifurcation} \end{cases}$$

Time $O(qn^3)$, q states, seq len n
compare: $O(qn)$ for profile HMM



Mutual Information

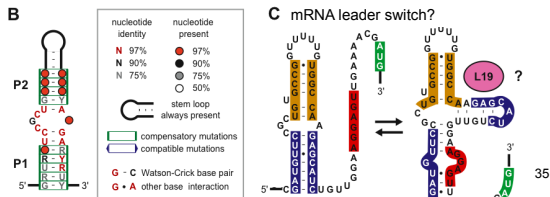
$$M_{ij} = \sum_{x_i, x_j} f_{x_i, x_j} \log_2 \frac{f_{x_i, x_j}}{f_{x_i} f_{x_j}}; \quad 0 \leq M_{ij} \leq 2$$

Max when no seq conservation but perfect pairing

MI = expected score gain from using a pair state

Finding optimal MI, (i.e. opt pairing of cols) is hard(?)

Finding optimal MI without pseudoknots can be done by dynamic programming



Rfam – an RNA family DB

Griffiths-Jones, et al., NAR '03, '05, '08

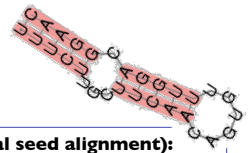
Biggest scientific computing user in Europe -
1000 cpu cluster for a month per release

Rapidly growing:

Rel 1.0, 1/03: 25 families, 55k instances	DB size:
Rel 7.0, 3/05: 503 families, 363k instances	~8GB
Rel 9.0, 7/08: 603 families, 636k instances	
Rel 9.1, 1/09: 1372 families, 1148k instances	
Rel 10.0, 1/10: 1446 families, 3193k instances	~160GB

48

Example Rfam Family



Input (hand-curated):

MSA "seed alignment"

SS_cons

Score Thresh T

Window Len W

Output:

CM

scan results & "full alignment"

phylogeny, etc.

IRE (partial seed alignment):

Hom. sap.	GUUCCUGCUUCAACAGUGUUGGAUGGAAC
Hom. sap.	UUUCUUC . UUCAACAGUGUUGGAUGGAAC
Hom. sap.	UUUCCUGUUCAACAGUGCUUGGA . GGAAC
Hom. sap.	UUUAUC . .AGUGACAGAUUCACU . AUAAA
Hom. sap.	UCUCUUGCUUCAACAGUGUUGGAUGGAAC
Hom. sap.	AUAUAC . .GGAAACAGUGUUCUCC . AUAAU
Hom. sap.	UCUUGC . .UUCAACAGUGUUGGACGGAAG
Hom. sap.	UGUAUC . .GGAGACAGUAUCUCC . AUAUG
Hom. sap.	AUAUAC . .GGAACAGUCCUCC . AUAAU
Cav. por.	UCUCCUGCUUCAACAGUGUUGGACGGAAC
Mus. mus.	UAUAUC . .GGAGACAGUAUCUCC . AUAUG
Mus. mus.	UUUCCUGCUUCAACAGUGUUGGACGGAAC
Mus. mus.	GUACUUGCUUCAACAGUGUUGGACGGAAC
Rat. nor.	UAUAUC . .GGAGACAGUACUCC . AUAUG
Rat. nor.	UAUCUUGCUUCAACAGUGUUGGACGGAAC
SS_cons	<<<< . . <<<< >>>> . >>>>