

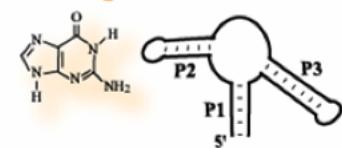
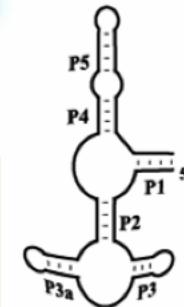
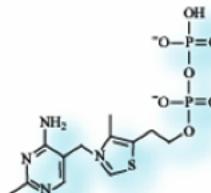
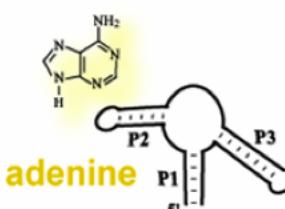
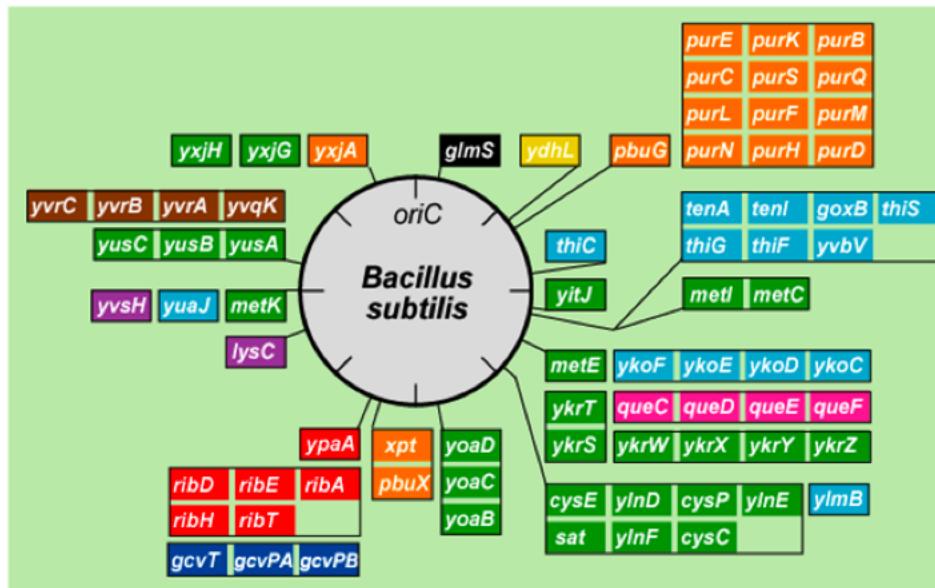
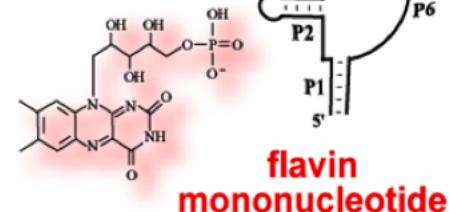
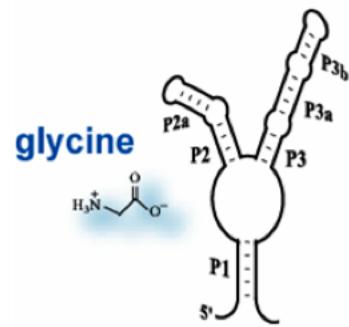
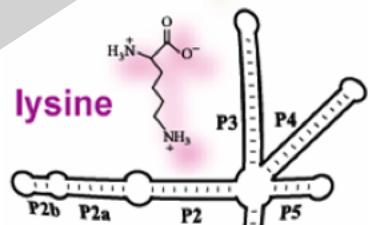
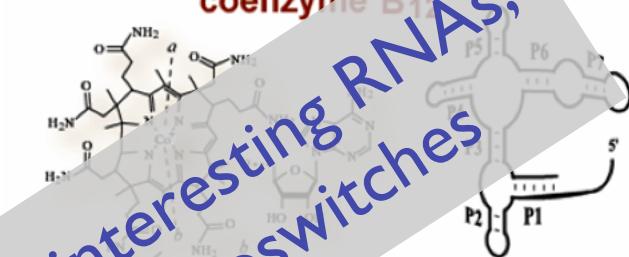
RNA Search and Motif Discovery

Genome 541
Intro to Computational
Molecular Biology

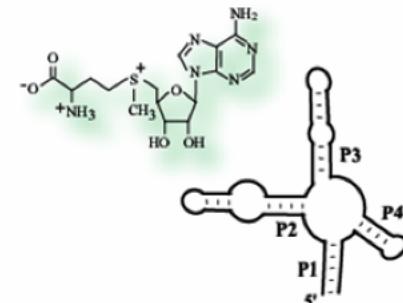
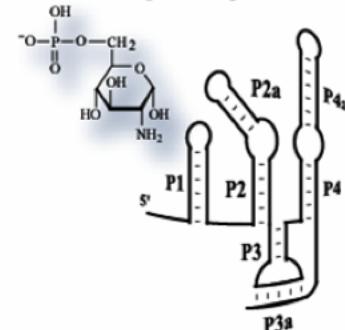
Day I

Many biologically interesting roles for RNA
RNA secondary structure prediction

Many interesting RNAs,
e.g. Riboswitches



glucosamine-6-phosphate



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: A Computation Order

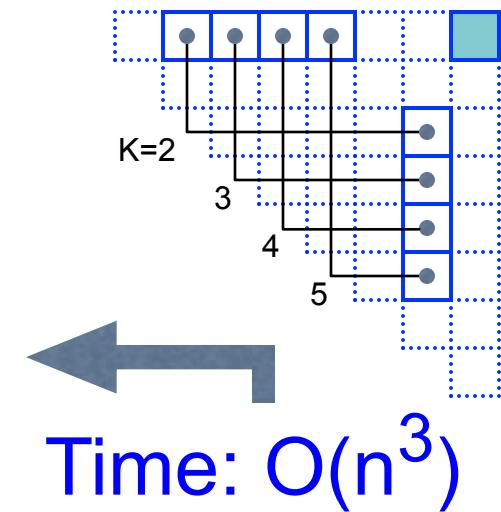
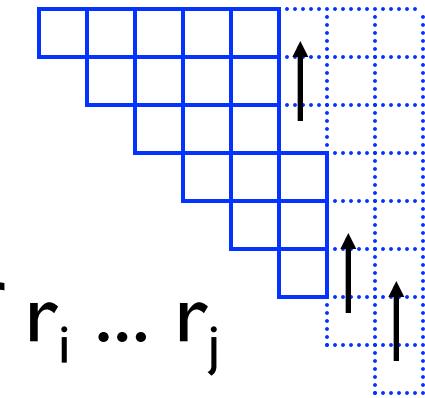
$B(i,j)$ = # pairs in optimal pairing of $r_i \dots r_j$

Or energy

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

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

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



Approaches, II

- + 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

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

Motif Description

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

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 Nusinov/Zucker RNA folding

Algorithms for searching

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

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”

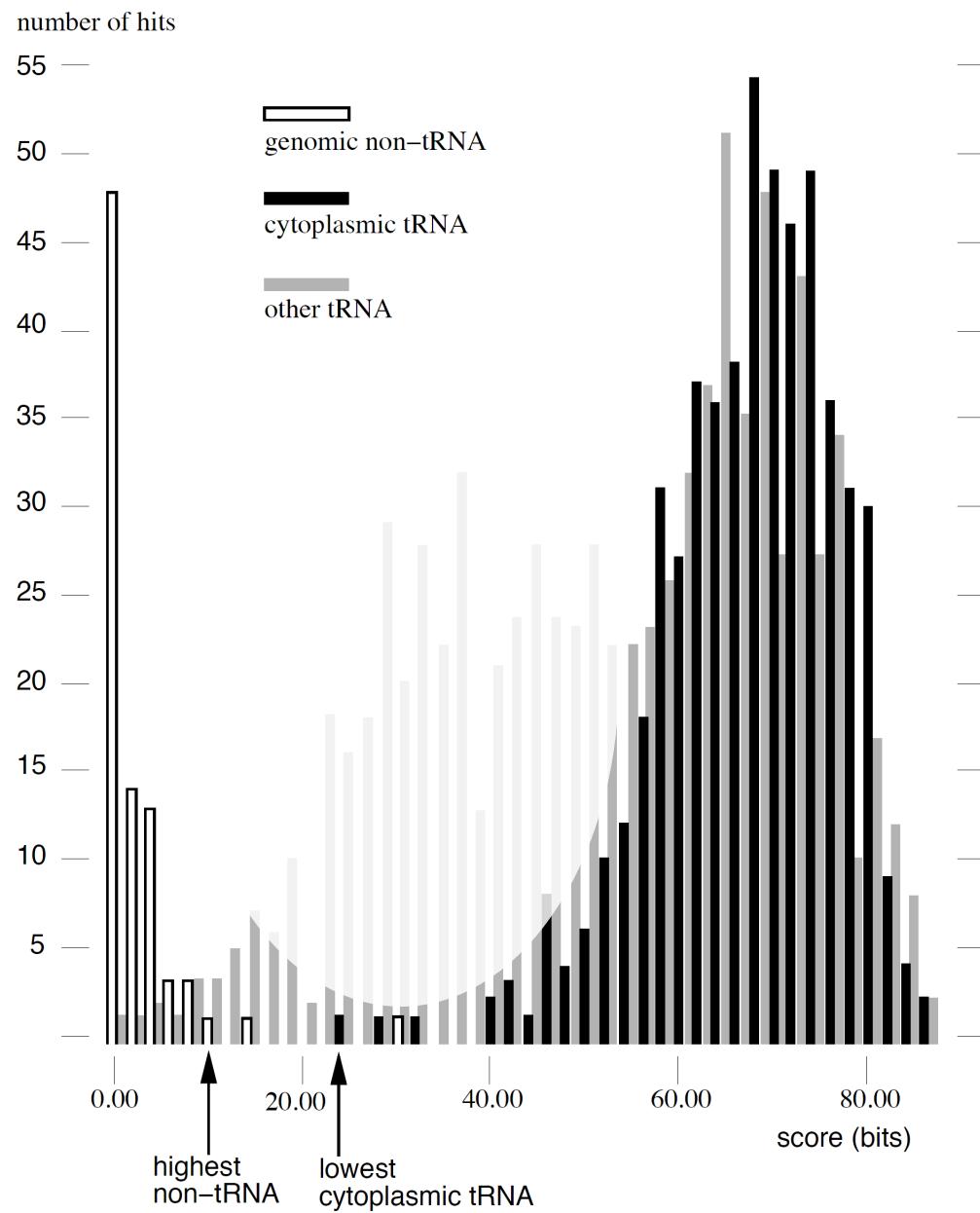
Scoring:

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

Viterbi approximation - find single best path

(Bonus: alignment & structure prediction)

Example: searching for tRNAs

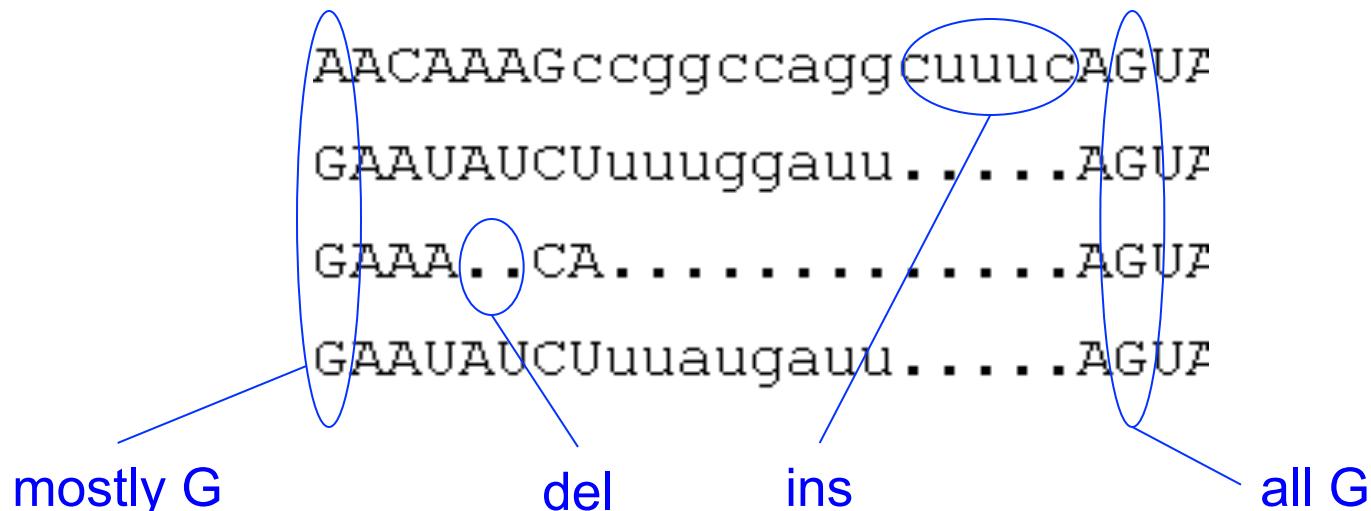


How to model an RNA “Motif”?

Conceptually, start with a profile HMM:

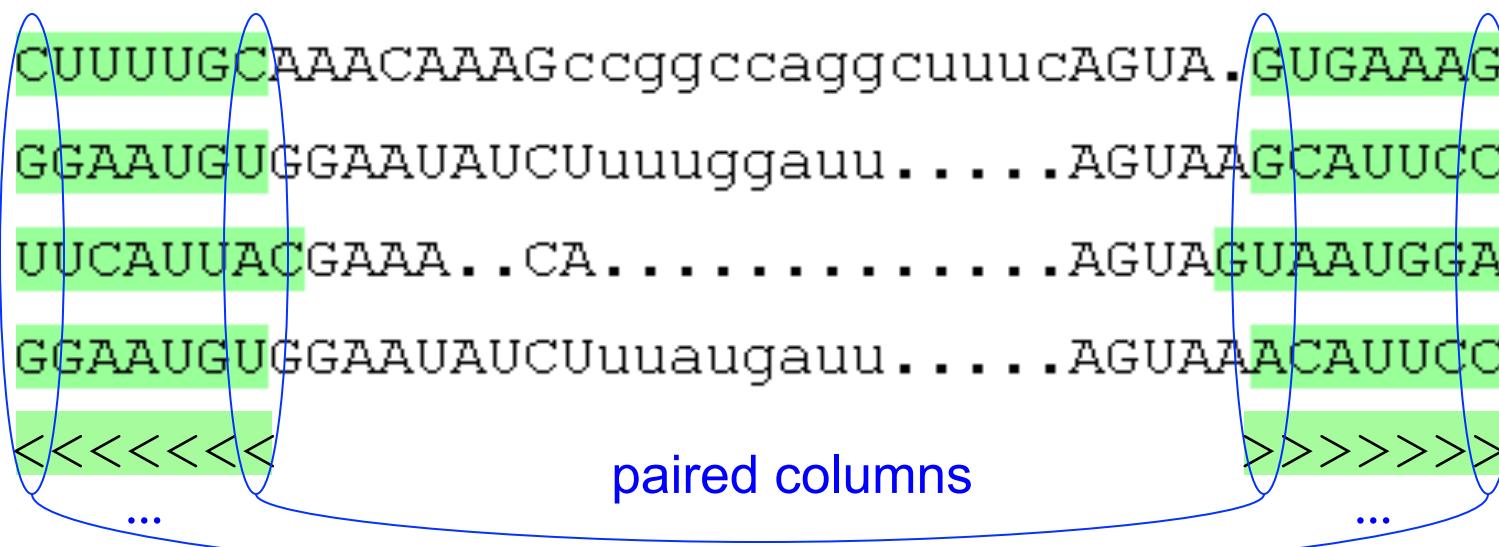
from a multiple alignment, estimate nucleotide/ insert/delete preferences for each position

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

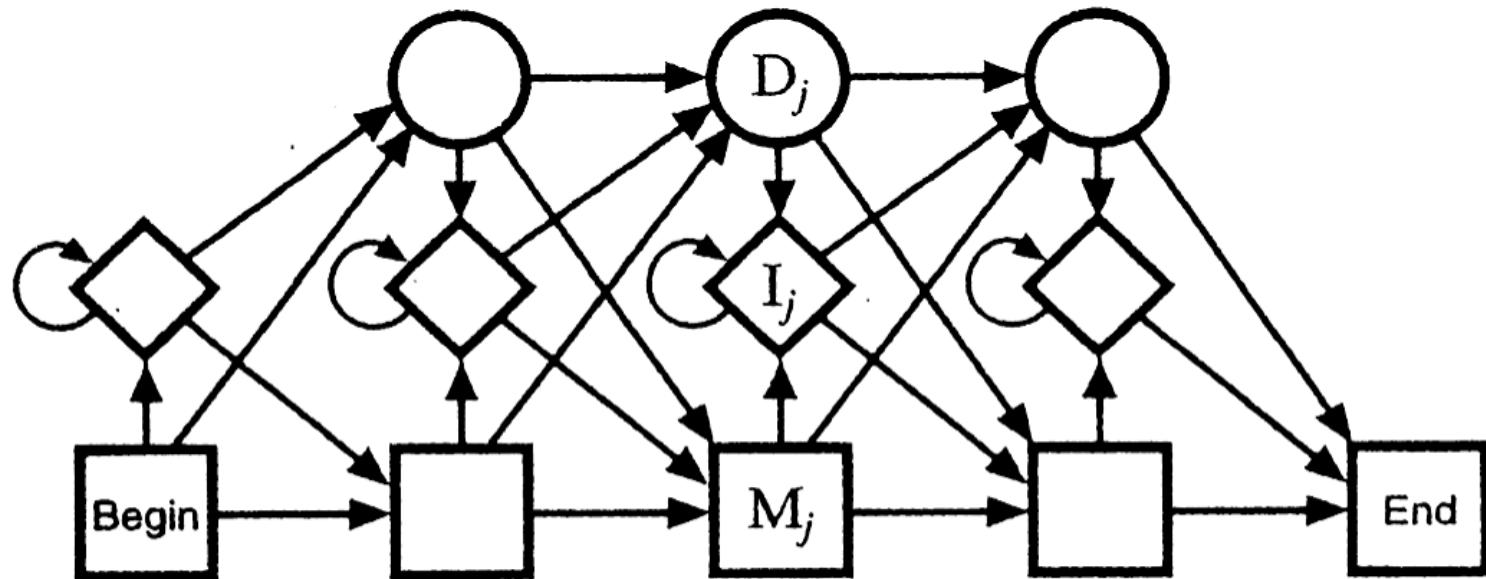


Figure 5.2 *The transition structure of a profile HMM.*

M_j : Match states (20 emission probabilities)

I_j : Insert states (Background emission probabilities)

D_j : Delete states (silent - no emission)

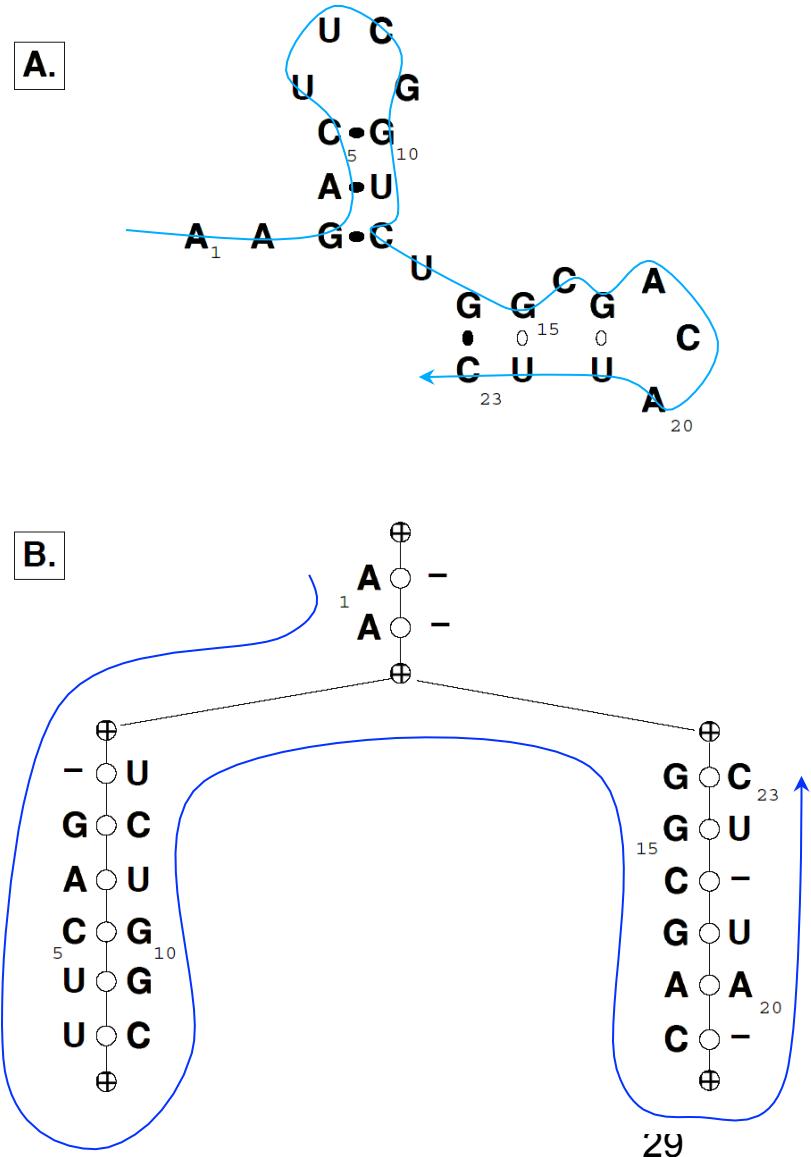
CM Structure

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)

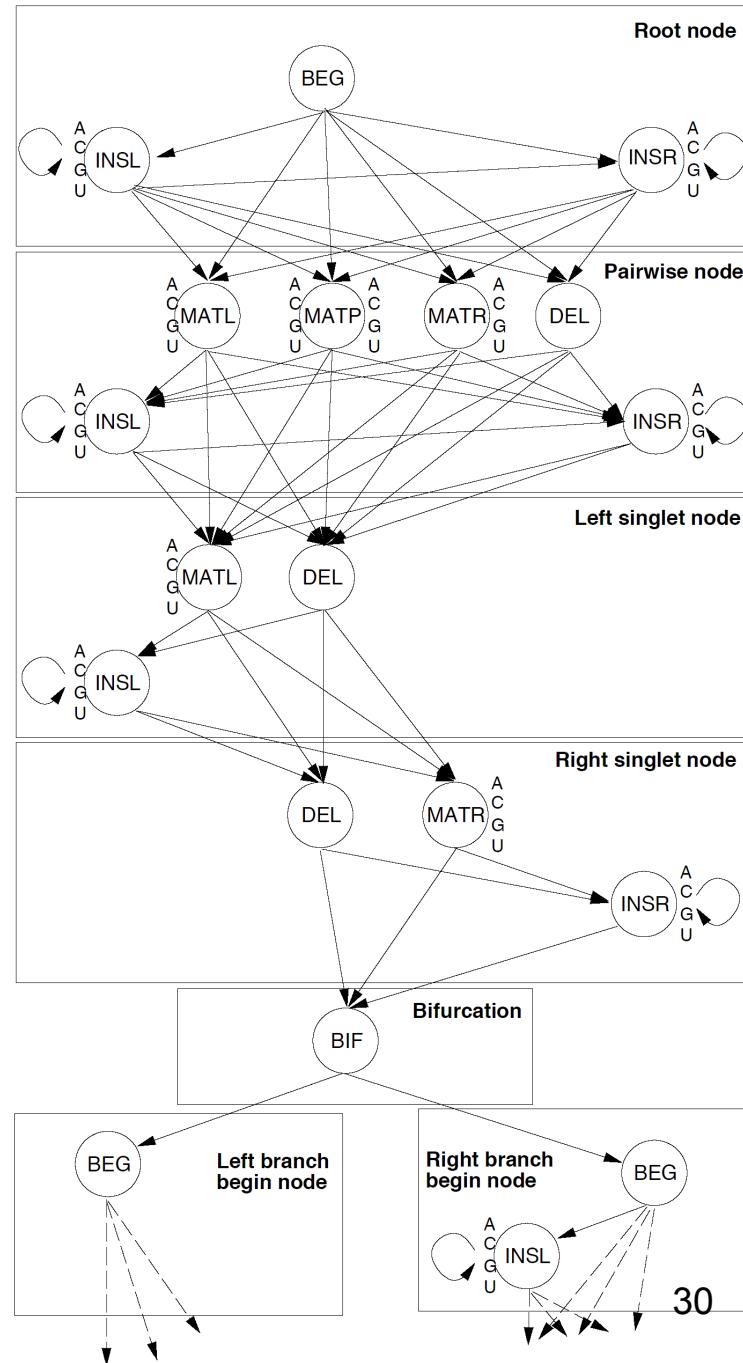


Overall CM Architecture

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)

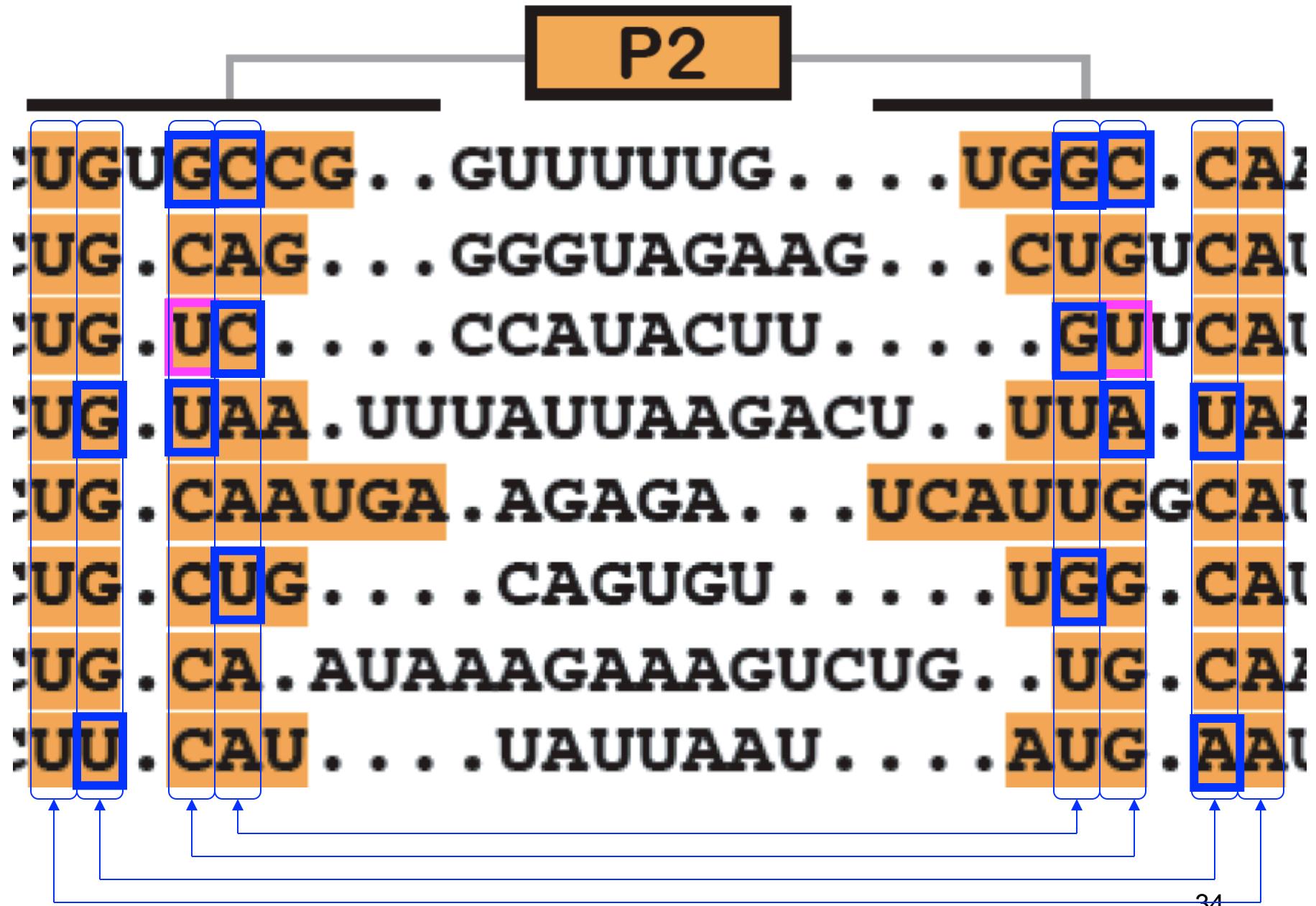
$S_{ij}^y = \max_{\pi} \log P(x_{ij} \text{ generated starting in state } y \text{ via path } \pi)$

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



Time $O(qn^3)$, q states, seq len n

compare: $O(qn)$ for profile HMM



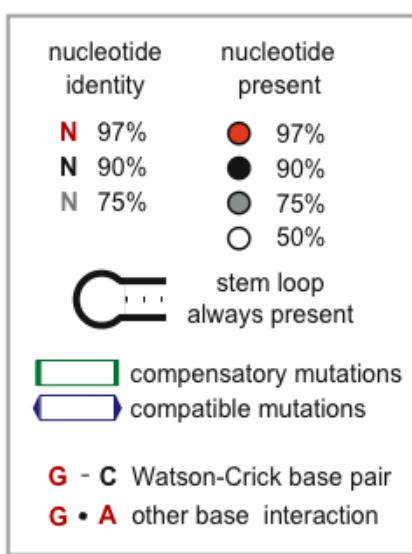
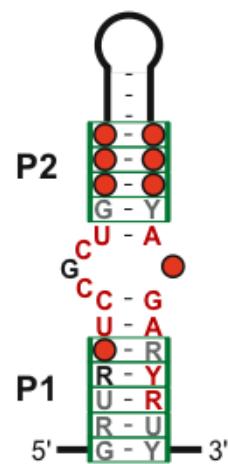
34

Covariation is strong evidence for base pairing

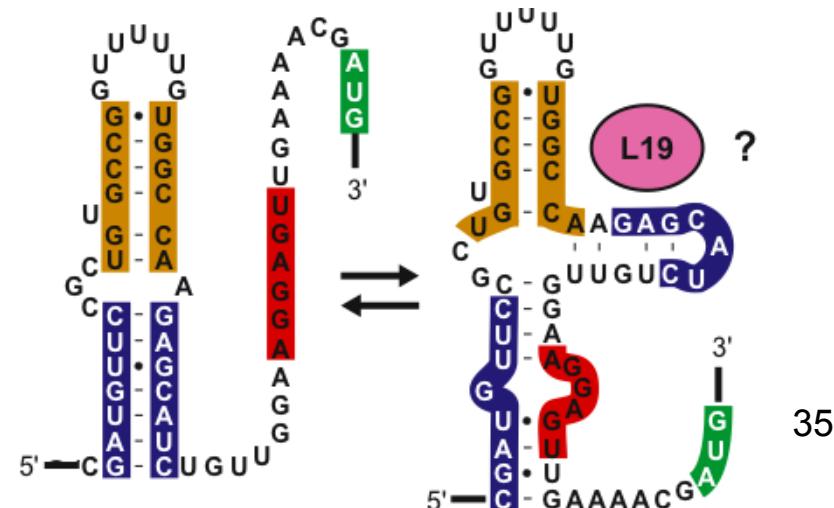
A mRNA leader

	-35	-10	TSS	P1	P2	RBS	Start	
<i>Bsu</i>	TTGCAT.	17.	TAAGAT.	40.	AAAACGAUGUUC	CGCUGUGCCG..	GUUUUUG... UGGC. CA	CACCAUCUG. 05. AGGAGU. 08. AUG
<i>Bha</i>	TTGTTC.	17.	TCTTCT.	17.	AUUACGAUGUUC	CGCUG. CAG..	GGGUAGAAG... CUGUCAU	GAGCAUCUG. 06. AGGAGG. 11. AUG
<i>Oih</i>	TTGAAAC.	17.	TATATT.	31.	UAAACGAUGUUC	CGCUG. UC..	CCAUAACUU.... GUUCAU	GAGCAUAG. 06. AGGAGU. 07. AUG
<i>Bce</i>	TTGCTA.	18.	TATGCT.	36.	UUAACGAUGUUC	CGCUG. UAA..	UUUAUUAAGACU.. UUA. UAA	GAGCAUCUG. 05. AGGAGA. 09. AUG
<i>Gka</i>	TTGCCT.	17.	TATCAT.	38.	AAAACGAUGUUC	CGCUG. CAAUGA.	AGAGA... UCAUUGGCAU	GAACAUUCUG. 04. AGGAGU. 08. AUG
<i>Bcl</i>	TTGTGC.	17.	TATGAT.	45.	AUUACGAUAUUC	CGCUG. CUG..	CAGUGU... UGG. CAU	GAAUGUCUG. 06. AGGAGG. 10. AUG
<i>Bac</i>	ATGACA.	17.	GATACT.	35.	AUAACGAUGUUC	CGCUG. CA.	AUAAAGAAAGUCUG.. UG. CA	CACCAUCUG. 05. AGGAGU. 08. AUG
<i>Lmo</i>	TTTACA.	17.	TAACCT.	28.	AUAACGAUAUUC	CGCUG. CAU..	UAUUAU... AUG. AAU	GAAUGUUG. 05. AGGAGA. 07. AUG
<i>Sau</i>	TTGAAA.	17.	TAACAT.	23.	AUCACUAUGAUC	CGCUG. CU..	AUAUAUUGUCG... AGGCAAG	PAACAUAGG. 04. AGAGGA. 09. AUG
<i>Cpe</i>	TTAAAG.	18.	TAACAT.	08.	GUACCGGCCGUC	CUCUGUCACA..	GAG..... UGUGUUAAGP	ACGUCAA. 17. AGGAGG. 08. AUG
<i>Chy</i>	TTGCAT.	17.	TATAAT.	09.	UACCAAACGUUC	CGCUG. GA..	CAGGGC... UC. CAU	GAACGUGC. 03. AGGAGG. 09. AUG
<i>Swo</i>	TTGAGA.	17.	TAAAAT.	16.	AAAAAGGUGGU	CGCUG. CAU..	AAACUAA... AAU	G. UAUGPACACC. 05. AGGAGG. 07. AUG
<i>Ame</i>	TTGCGG.	17.	TATAAT.	10.	UUACGGGCCGUC	CUCUA. UAC..	AGGA... GUA. UAAGPACGCUA	. 07. AGGAGG. 07. AUG
<i>Dre</i>	TTGCC.	17.	TATAAT.	16.	UUACGGACGGUC	CGCUG. CCU..	CUGGGAA... AGG. UAAGPACGCUA	. 04. AGGAAG. 12. GUG
<i>Spn</i>	TTTACT.	17.	TAACAT.	28.	AUACAGUUUAUC	CGCUG. AGGA..	AGAU... UCCU. CAACAUUGACAA	. 04. AGGAGA. 05. AUG
<i>Smu</i>	TTTACA.	17.	TACAAT.	26.	AAACGGCUAAUC	CGCUG. AG..	ACAGAGCA... CU. UAUGAUUAGUA	. 04. AGGAGA. 07. AUG
<i>Lpl</i>	TTGCGT.	18.	TATTCT.	21.	UUAACGAUGUUC	CGCUG. AC..	CAGGUU... GU. CACGAUGUCGG	. 04. AGGAAG. 09. AUG
<i>Efa</i>	TTTACA.	17.	TAACAT.	28.	AUUACAAUAUUC	CGCUG. UGG. CA..	GAAG... UGACCA. UAA	GAUUAUUG. 06. AGGAGA. 08. AUG
<i>Ljo</i>	TTTACA.	17.	TAACAT.	25.	UUAUGGGUAAUUC	CGCUG. GCAC..	AAG... GUGUUGAU	GAUUGC. 03. AGGAGA. 07. AUG
<i>sth</i>	TAGACA.	17.	TAAGAT.	29.	UAACGGCUAAUC	CGCUG. AGA.	CACAGAGGU. UGCCUCU. UAA	GAUUAGUA. 03. AGGAGU. 08. AUG
<i>Lac</i>	TTAAA.	17.	TTACTT.	39.	UUAUGGGUAAUUC	CGCUG. ACG..	CUGGUA... CGUUGAU	CAAUGC. 03. AGGAGA. 10. AUG
<i>Spy</i>	TTTACA.	17.	TAGAAT.	29.	UUACGGCUAAUC	CGCUG. AG..	ACAAGUA... CU. UAA	GAUUAGUA. 03. AGGAGA. 06. AUG
<i>Lsa</i>	TTTTAA.	17.	TAAAAT.	26.	ACAACGAUAUUC	CGCUG. GCG..	CAAGA... CGUUAU	AAUACUG. 06. AGGAGA. 07. AUG
<i>Lsl</i>	TTTACT.	17.	TATTTT.	24.	AUAACGAUAUUC	CGCUG. C..	AACUG... GACAU	GAUUGUCGG. 04. AGGAAA. 07. AUG
<i>Fnu</i>	TTGACA.	17.	TAAAAT.	12.	AAUUCGAUAUUC	CGCUG. UAA..	UAAA... UUA. AAU	GAUUAUCUU. 04. AGGAAG. 02. AUG

B



C mRNA leader switch?



Mutual Information

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

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

M.I. Example (Artificial)

* ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ *

	1	2	3	4	5	6	7	8	9	*
A	G	A	U	A	A	U	C	U		
A	G	A	U	C	A	U	C	U		
A	G	A	C	G	U	U	C	U		
A	G	A	U	U	U	U	C	U		
A	G	C	C	A	G	G	C	U		
A	G	C	G	C	G	G	C	U		
A	G	C	U	G	C	G	C	U		
A	G	C	A	U	C	G	C	U		
A	G	G	U	A	G	C	C	U		
A	G	G	G	C	G	C	C	U		
A	G	G	U	G	U	C	C	U		
A	G	G	C	U	U	C	C	U		
A	G	U	A	A	A	A	C	U		
A	G	U	C	C	A	A	C	U		
A	G	U	U	G	C	A	C	U		
A	G	U	U	U	C	A	C	U		
A	16	0	4	2	4	4	4	0	0	
C	0	0	4	4	4	4	4	16	0	
G	0	16	4	2	4	4	4	0	0	
U	0	0	4	8	4	4	4	0	16	

MI:	1	2	3	4	5	6	7	8	9
9	0	0	0	0	0	0	0	0	0
8	0	0	0	0	0	0	0	0	0
7	0	0	2	0.30	0	0	1		
6	0	0	1	0.55	1				
5	0	0	0	0.42					
4	0	0	0.30						
3	0	0							
2	0								
1									

Cols 1 & 9, 2 & 8: perfect conservation & *might* be base-paired, but unclear whether they are. M.I. = 0

Cols 3 & 7: No conservation, but always W-C pairs, so seems likely they do base-pair. M.I. = 2 bits.

Cols 7->6: unconserved, but each letter in 7 has only 2 possible mates in 6. M.I. = 1 bit.

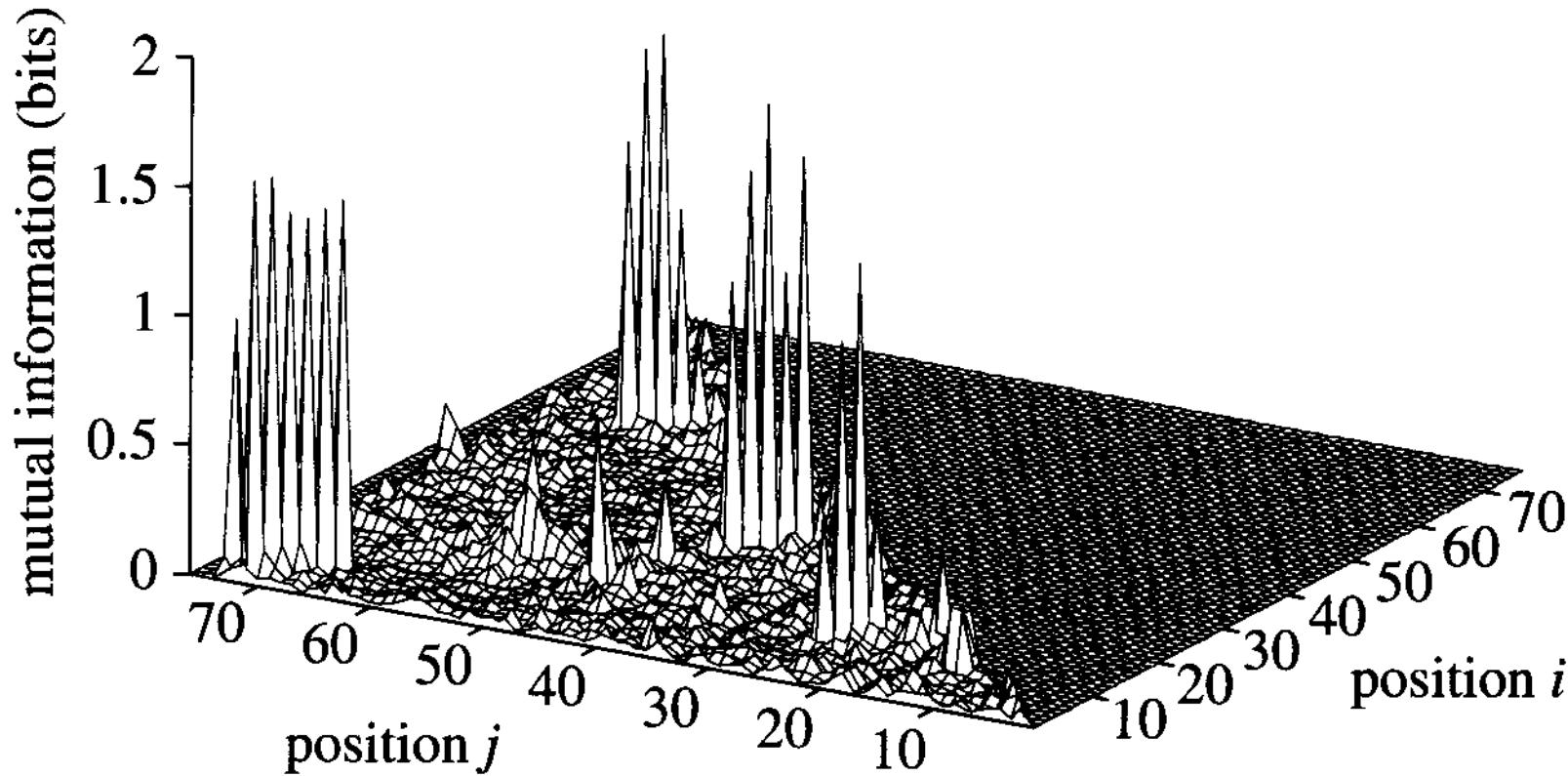
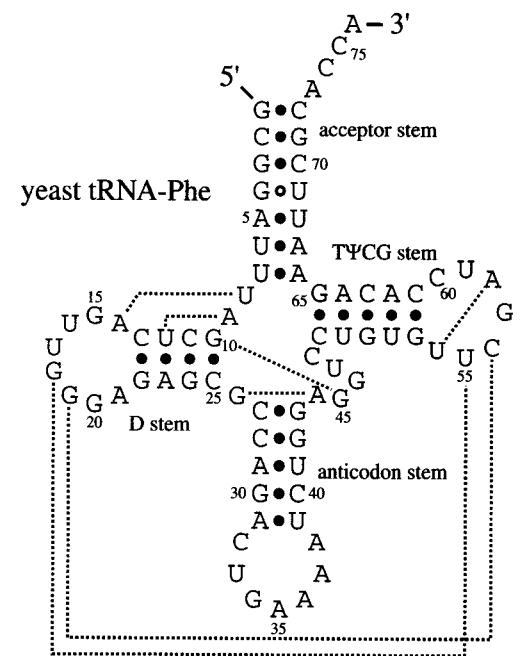


Figure 10.6 A mutual information plot of a tRNA alignment (top) shows four strong diagonals of covarying positions, corresponding to the four stems of the tRNA cloverleaf structure (bottom; the secondary structure of yeast phenylalanine tRNA is shown). Dashed lines indicate some of the additional tertiary contacts observed in the yeast tRNA-Phe crystal structure. Some of these tertiary contacts produce correlated pairs which can be seen weakly in the mutual information plot.



Primary vs Secondary Info

Dataset	Avg. id	Min id	Max id	ClustalV accuracy	1° info (bits)	2° info (bits)
TEST	.402	.144	1.00	64%	43.7	30.0-32.3
SIM100	.396	.131	.986	54%	39.7	30.5-32.7
SIM65	.362	.111	.685	37%	31.8	28.6-30.7

Disallowing / allowing
pseudoknots

$$\left(\sum_{i=1}^n \max_j M_{i,j} \right) / 2$$

Comparison to TRNASCAN

Fichant & Burks - best heuristic then

97.5% true positive

0.37 false positives per MB

CM AI415 (trained on trusted alignment)

> 99.98% true positives

< 0.2 false positives per MB

Current method-of-choice is “tRNAscanSE”, a CM-based scan with heuristic pre-filtering (including TRNASCAN?) for performance reasons.

Slightly different
evaluation criteria

tRNAscanSE

Uses 3 older heuristic tRNA finders as prefilter

Uses CM built as described for final scoring

Actually 3(?) different CMs

eukaryotic nuclear

prokaryotic

organellar

Used in all genome annotation projects

An Important Application: Rfam

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:

~8GB

Rel 7.0, 3/05: 503 families, 363k instances

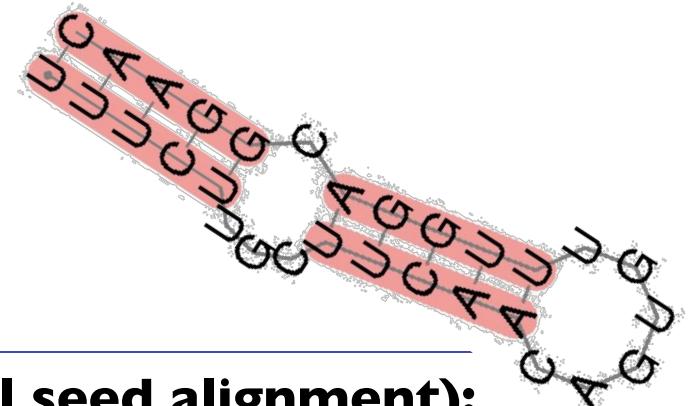
Rel 9.0, 7/08: 603 families, 636k instances

Rel 9.1, 1/09: 1372 families, 1148k instances

Rel 10.0, 1/10: 1446 families, 3193k instances

~160GB

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.	GUUCCUGCUUCAACAGUGUUUUGGAUGGAAAC
Hom. sap.	UUUCUUC. UUCAACAGUGUUUUGGAUGGAAAC
Hom. sap.	UUUCCUGUUUCAACAGUGCUUGGA. GGAAC
Hom. sap.	UUUAUC.. AGUGACAGAGUUUCACU. AUAAA
Hom. sap.	UCUCUUGCUUCAACAGUGUUUUGGAUGGAAAC
Hom. sap.	AUUAUC.. GGGAACAGUGUUUUCCC. AUAAU
Hom. sap.	UCUUGC.. UUCAACAGUGUUUUGGACGGAAG
Hom. sap.	UGUAUC.. GGAGACAGUGAUCUCC. AUAUG
Hom. sap.	AUUAUC.. GGAAGCAGUGCCUUCC. AUAAU
Cav. por.	UCUCCUGCUUCAACAGUGCUUGGACGGAGC
Mus. mus.	UAUAUC.. GGAGACAGUGAUCUCC. AUAUG
Mus. mus.	UUUCCUGCUUCAACAGUGCUUGAACCGGAAC
Mus. mus.	GUACUUGCUUCAACAGUGUUUUGAACCGGAAC
Rat. nor.	UAUAUC.. GGAGACAGUGACCUC. AUAUG
Rat. nor.	UAUCUUGCUUCAACAGUGUUUUGGACGGAAC
SS_cons	<<<<...<<<<.....>>>. >3>>