

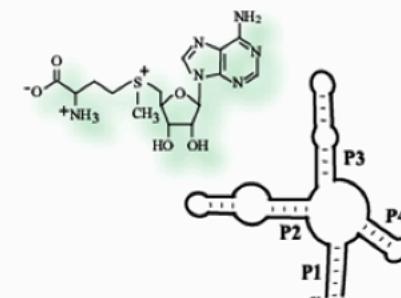
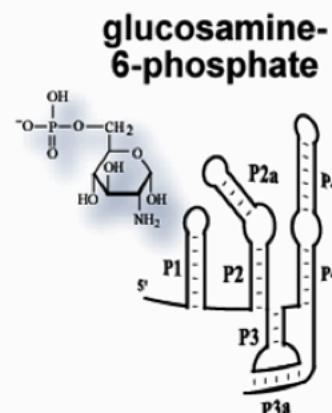
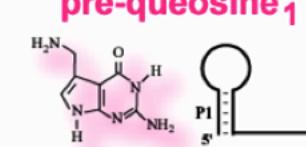
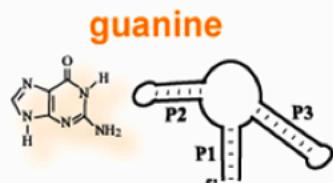
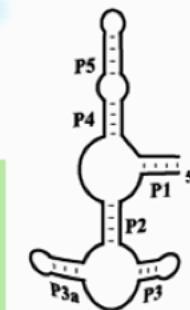
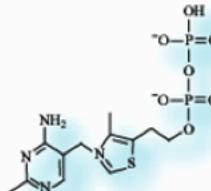
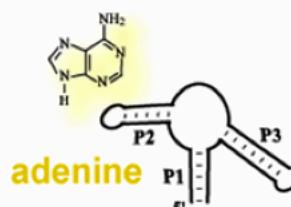
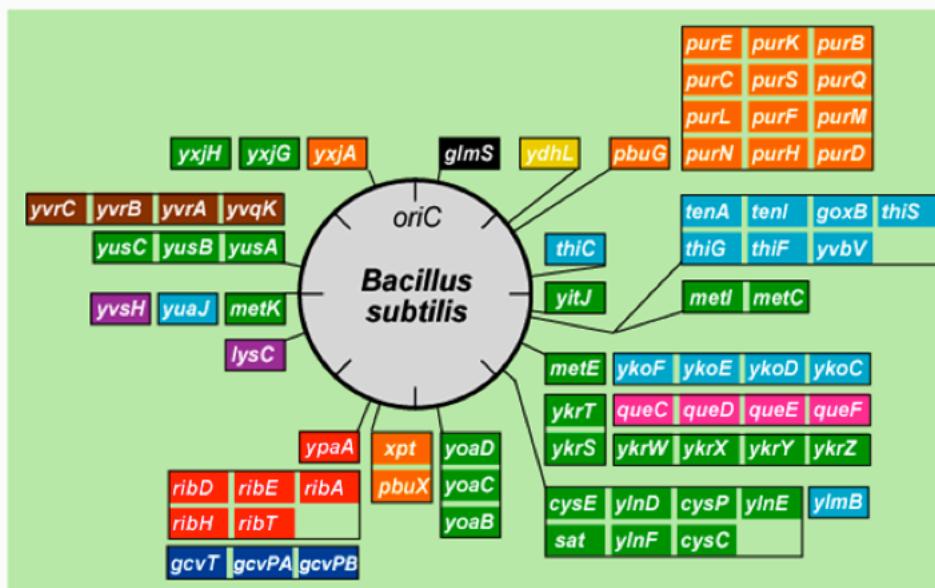
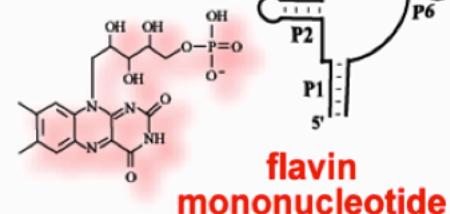
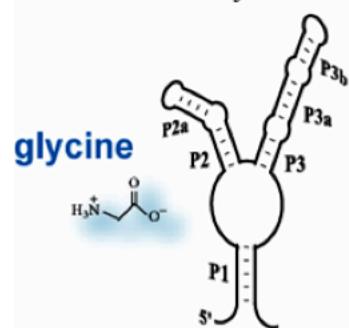
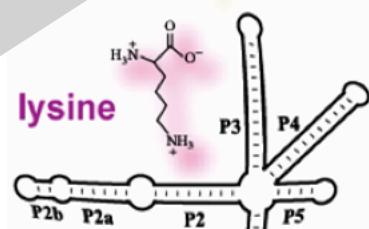
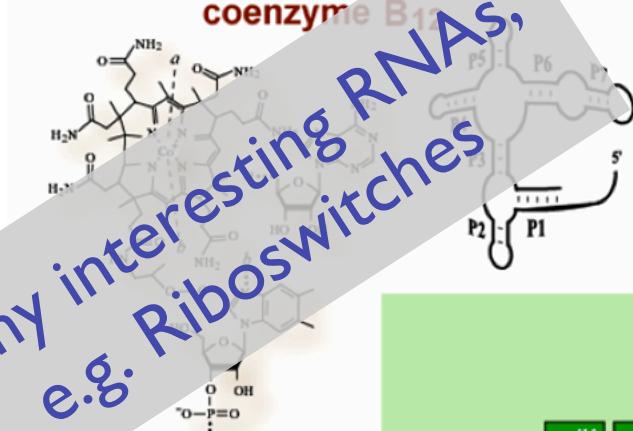
RNA Search and Motif Discovery

**CSEP 590 B
Computational Biology**

Previous Lecture

Many biologically interesting roles for RNA
RNA secondary structure prediction

Many interesting RNAs,
e.g. Riboswitches



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: Structure Prediction

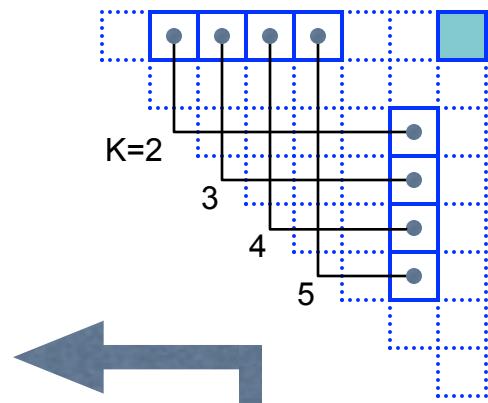
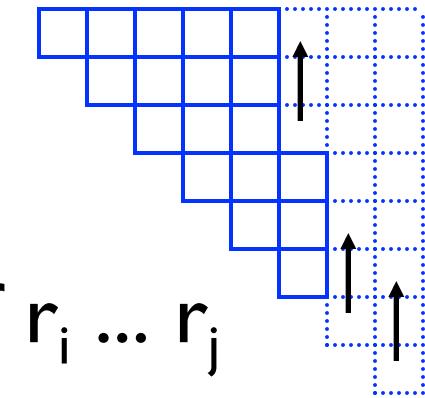
Computation Order

$B(i,j)$ = **# 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:

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



Time: $O(n^3)$

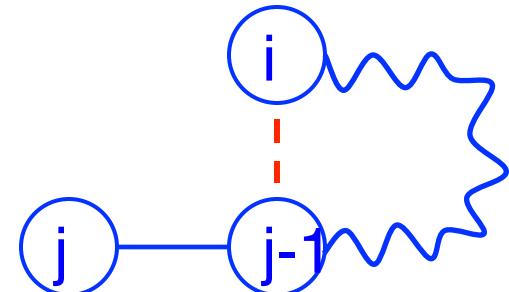
Loop-based energy version is better; recurrences similar, slightly messier

Optimal pairing of $r_i \dots r_j$

Two possibilities

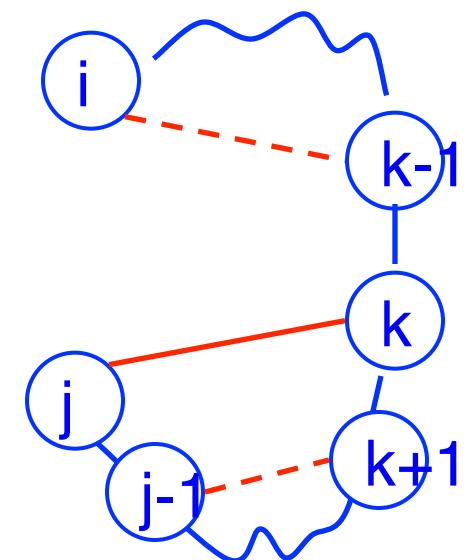
j Unpaired:

Find best pairing of $r_i \dots r_{j-1}$



j Paired (with some k):

Find best $r_i \dots r_{k-1}$ +
best $r_{k+1} \dots r_{j-1}$ plus 1



Why is it slow?

Why do pseudoknots matter?

Today

Structure prediction via comparative analysis

Covariance Models (CMs) represent
RNA sequence/structure motifs

Fast CM search

Motif Discovery

Applications in prokaryotes & vertebrates

Approaches, II

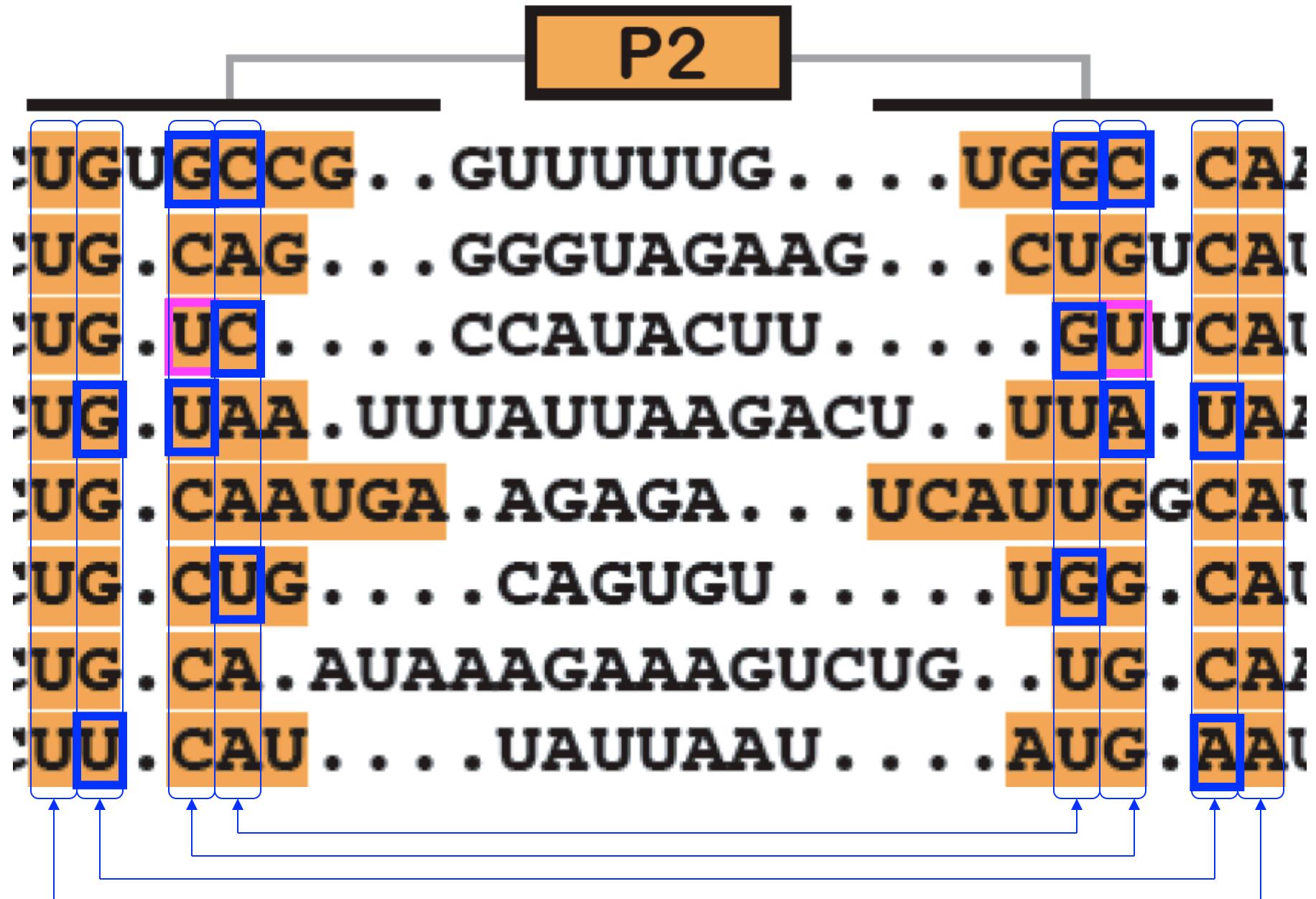
Comparative sequence analysis

- + 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)



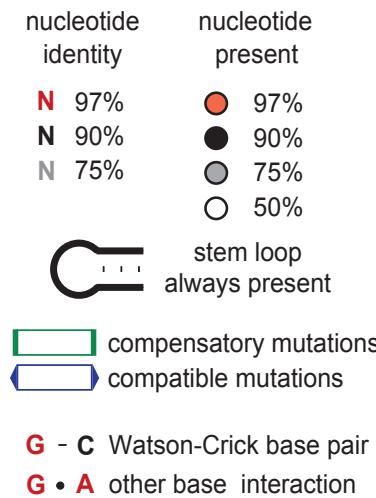
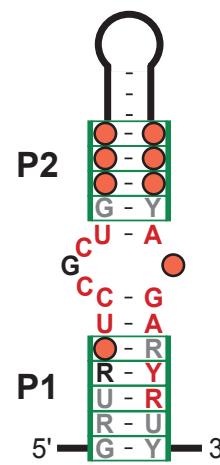
Covariation is strong evidence for base pairing

Example: Ribosomal Autoregulation:
Excess L19 represses L19 (RF00556; 555-559 similar)

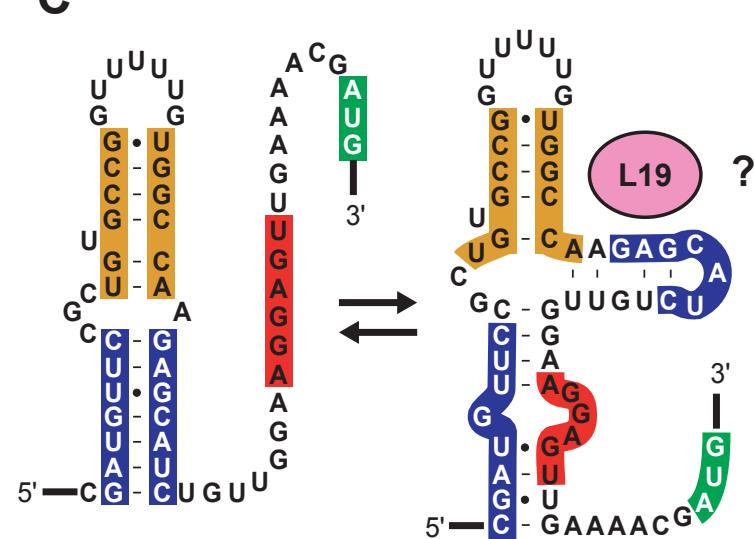
A L19 (*rplS*) mRNA leader

	-35	-10	TSS →	P1	P2	RBS	Start			
<i>Bsu</i>	TTGCAT.	17.	TAAGAT.	40. AAAAC	GAUGUUC	CGCUGUGCCG.. GUUUUUG...	GGC. CAAAGAGCAUC	UG. 05. AGGAGU. 08. AUG		
<i>Bha</i>	TTGTTTC.	17.	TCTTCT.	17. AUUAC	GAUGUUC	CGCUG. CAG... GGGUAGAAG...	CUGUCAUGAGCAUC	UG. 06. AGGAGG. 11. AUG		
<i>Oih</i>	TTGAAC.	17.	TATATT.	31. UAAAC	GAUGUUC	CGCUG. UC... CCAUACUU...	GUUCAUGAGCAUU	AG. 06. AGGAGU. 07. AUG		
<i>Bce</i>	TTGCTA.	18.	TATGCT.	36. UUAAC	GAUGUUC	CGCUG. UAA. UUUUAUAAGACU...	UUA. UAAGAGCAUC	UG. 05. AGGAGA. 09. AUG		
<i>Gka</i>	TTGCCT.	17.	TATCAT.	38. AAAAC	GAUGUUC	CGCUG. CAAUGA. AGAGA...	UCAUUGGCAUGAACAU	UG. 04. AGGAGU. 08. AUG		
<i>Bcl</i>	TTGTGC.	17.	TATGAT.	45. AUUAC	GAUAUUC	CGCUG. CUG... CAGUGU...	UGG. CAUGAAUGUC	UG. 06. AGGAGG. 10. AUG		
<i>Bac</i>	ATGACA.	17.	GATACT.	35. AUUAC	GAUGUUC	CGCUG. CA. AUAAAAGAAAGUCUG...	UG. CAAGAGCAUC	UG. 05. AGGAGU. 08. AUG		
<i>Lmo</i>	TTTACA.	17.	TAACCT.	28. AUUAC	GAUAUUC	CGCUU. CAU... UAUUAAU...	AUG. AAUGAAUGUU	UG. 05. AGGAGA. 07. AUG		
<i>Sau</i>	TTGAAA.	17.	TAACAT.	23. AUCAC	UAUCAUCC	CGCUG. CU... AUUAUAUUGUCG...	AGGCAAGAACAU	AG. 04. AGAGGA. 09. AUG		
<i>Cpe</i>	TTAAAG.	18.	TAACAT.	08. GUACC	GGCGGUCC	CUCUGUCACA...	UGUGUUAAGAACGUCA	AA. 17. AGGAGG. 08. AUG		
<i>Chy</i>	TTGCAT.	17.	TATAAT.	09. UACCAA	ACGUUC	CGCUG. GA... CAGGGC...	UC. CAUGAACGUGCC	03. AGGAGG. 09. AUG		
<i>Swo</i>	TTGAGA.	17.	TAAAAT.	16. AAAAA	GGUGGUCC	CGCUG. CAUU...	AAUG. UAUGAACACC	UU. 05. AGGAGG. 07. AUG		
<i>Ame</i>	TTGCGG.	17.	TATAAT.	10. UUACG	GGCGGUCC	CUCUA. UAC...	GU. UAAGAACGUCA	UA. 07. AGGAGG. 07. AUG		
<i>Dre</i>	TTGCC.	17.	TATAAT.	16. UUACG	GACGGUCC	CGCUG. CCU...	CUGGAA...	AGG. UAAGAACGUCA. 04. AGGAAG. 12. GUG		
<i>Spn</i>	TTTACT.	17.	TAAACT.	28. AUAC	GUUAUCC	CGCUG. AGGA...	AGAU...	UCCU. CAAGAUUGACAA. 04. AGGAGA. 05. AUG		
<i>Smu</i>	TTTACA.	17.	TACAAT.	26. AAACG	GCUAUAC	CGCUG. AG...	ACAGAGCA...	CU. UAUGAUUAAGUA. 04. AGGAGA. 07. AUG		
<i>Lpl</i>	TTGCGT.	18.	TATTCT.	21. UUAC	GAUGUUC	CGCUG. AC...	CAGGUU...	GU. CACGAAUGUC	GG. 04. AGGAAG. 09. AUG	
<i>Efa</i>	TTTACA.	17.	TAAACT.	28. AUUAC	AAUAUUC	CGCUG. UGG. CA...	GAAG...	UGACCA. UAAGAUAU	UG. 06. AGGAGA. 08. AUG	
<i>Ljo</i>	TTTACA.	17.	TAAACT.	25. UUAUG	GGUAUUC	CGCUG. GCAC...	AAG...	GUGUUGAU	GAAUGCC	GU. 03. AGGAGA. 07. AUG
<i>Sth</i>	TAGACA.	17.	TAAGAT.	29. UUACG	GGCUAAUC	CGCUG. AGA. CACAGAGGU...	UGCUCU...	UAAGAUUA	GUAA. 03. AGGAGU. 08. AUG	
<i>Lac</i>	TTAAAA.	17.	TTACTT.	39. UUAUG	GGGUAUUC	CGCUG. ACG...	CUGGU...	CGUUGAU	GAAUGCC	GA. 03. AGGAGA. 10. AUG
<i>Spy</i>	TTTACA.	17.	TAGAAAT.	29. UUACG	GGCUAAUC	CGCUA. AG...	ACAAGUA...	CU. UAAGAUUA	GUAA. 03. AGGAGA. 06. AUG	
<i>Lsa</i>	TTTTAA.	17.	TAAAAT.	26. ACAAC	GAUAUUC	CGCUG. GCG...	CAAGA...	CGUUAU	GAAUAUC	UG. 06. AGGAGA. 07. AUG
<i>Lsl</i>	TTTACT.	17.	TATTCT.	24. AUUAC	GAUAUUC	CGCUG. C...	AACUG...	GACAU	GAAUGUC	GG. 04. AGGAAA. 07. AUG
<i>Fnu</i>	TTGACA.	17.	TTAAAT.	12. AAUUC	GAUAUUC	CGCUU. UAA...	UAAA...	UUA. AAU	GAUAUAC	UU. 04. AGGAAG. 02. AUG

B



C



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 (*below*)

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.

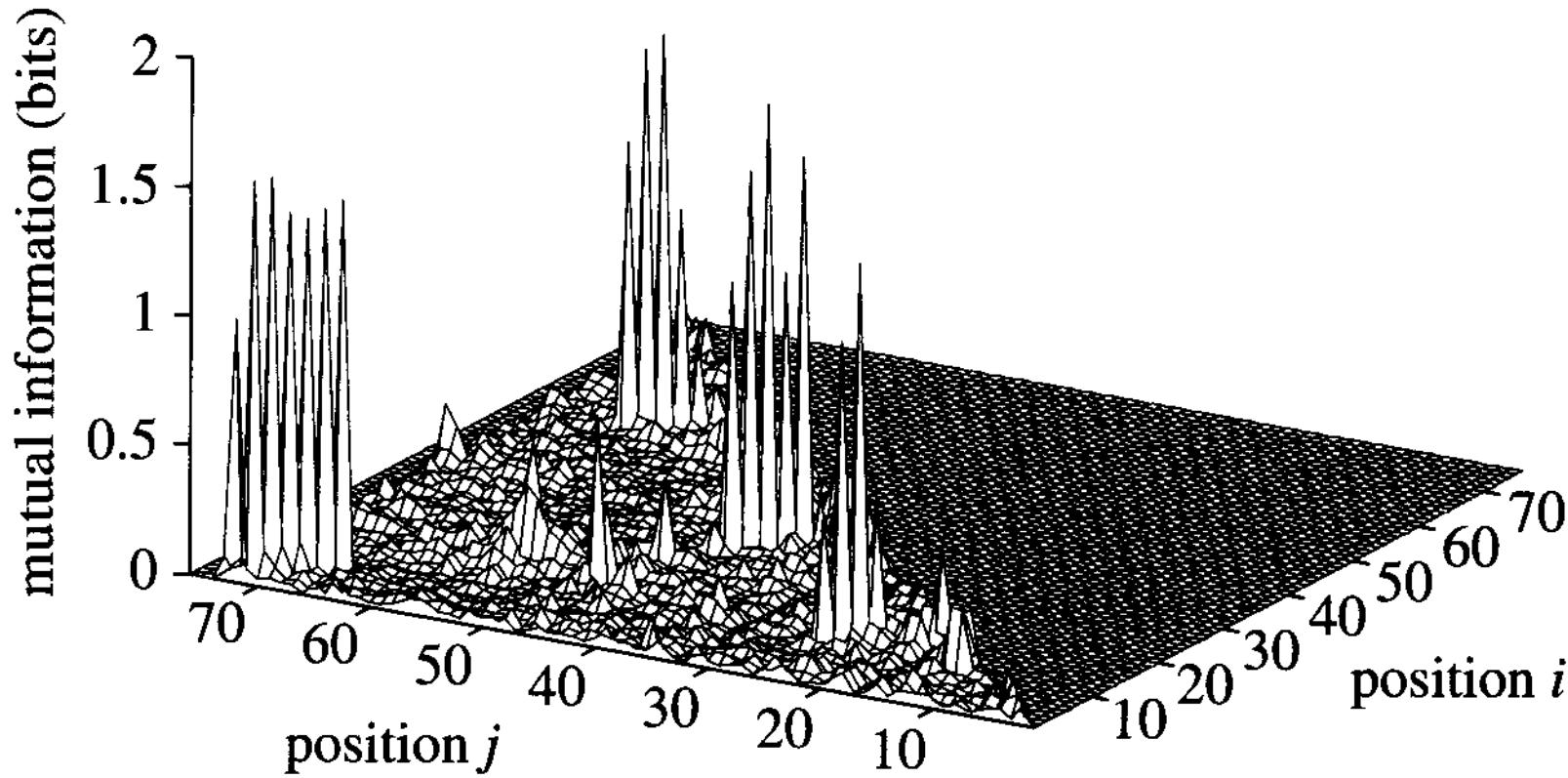
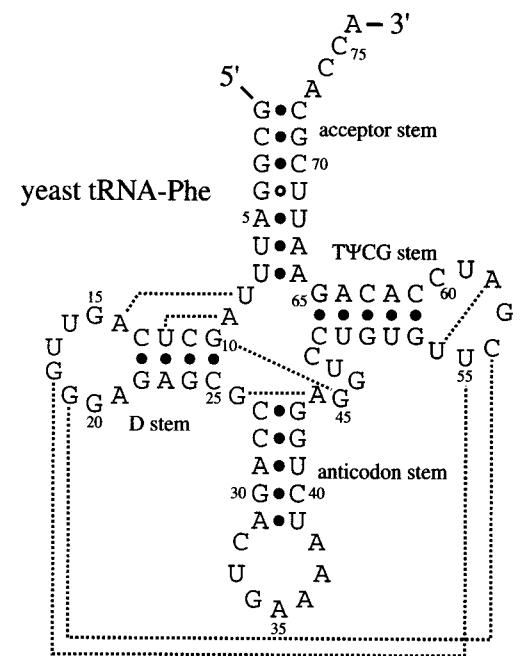


Figure 10.6 A mutual information plot of a tRNA alignment (top) shows four strong diagonals of covarying positions, corresponding to the four stems of the tRNA cloverleaf structure (bottom; the secondary structure of yeast phenylalanine tRNA is shown). Dashed lines indicate some of the additional tertiary 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.



MI-Based Structure-Learning

Find best (max total MI) subset of column pairs among $i \dots j$, subject to absence of pseudo-knots

$$S_{i,j} = \max \begin{cases} S_{i,j-1} & j \text{ unpaired} \\ \max_{i \leq k < j-4} S_{i,k-1} + M_{k,j} + S_{k+1,j-1} & j \text{ paired} \end{cases}$$

“Just like Nussinov/Zucker folding”

BUT, need enough data---enough sequences at right phylogenetic distance

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

Eddy & Durbin 1994: 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”

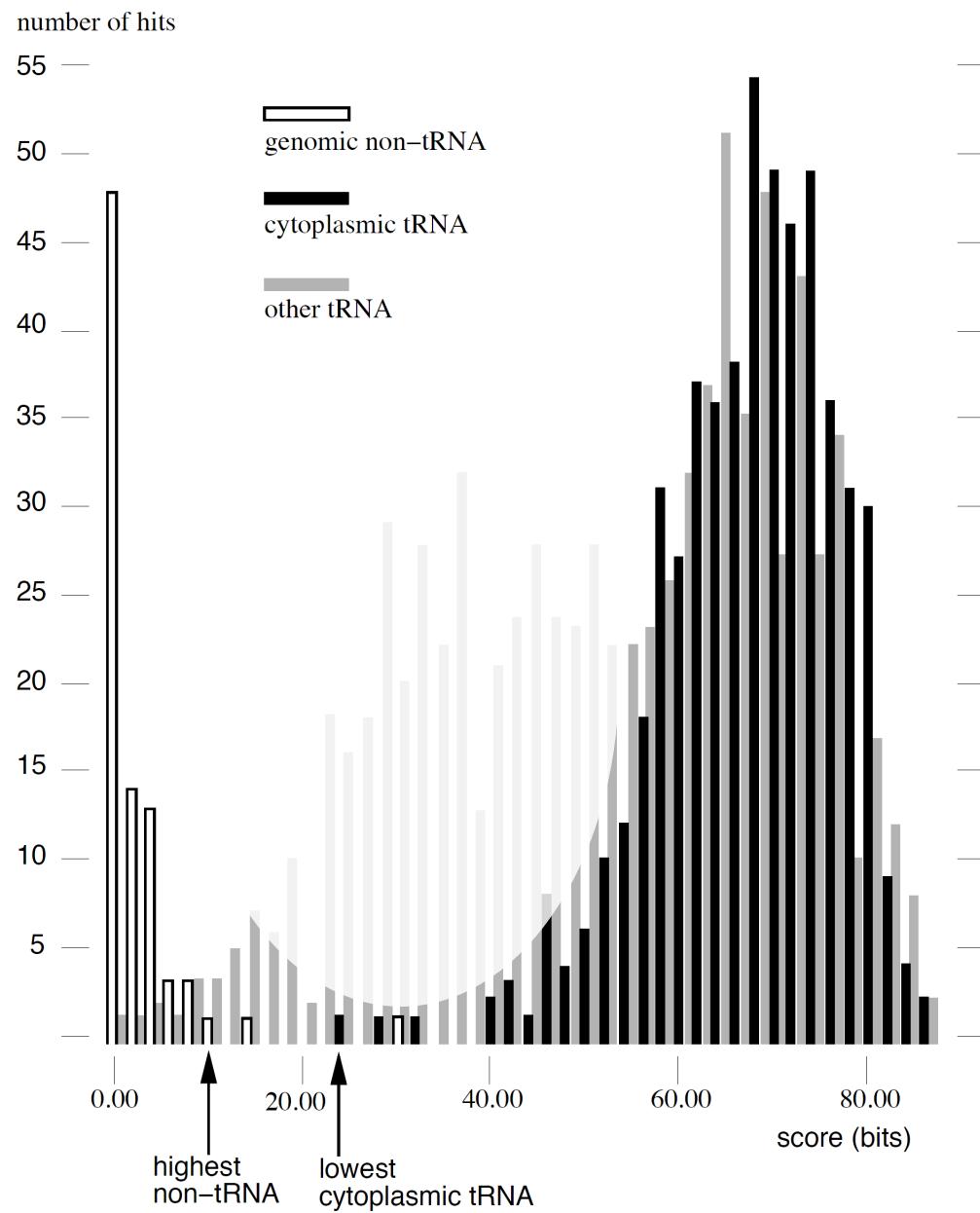
Scoring:

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

Viterbi approximation - find single best path

(Bonus: alignment & structure prediction)

Example: searching for tRNAs



Profile Hmm Structure

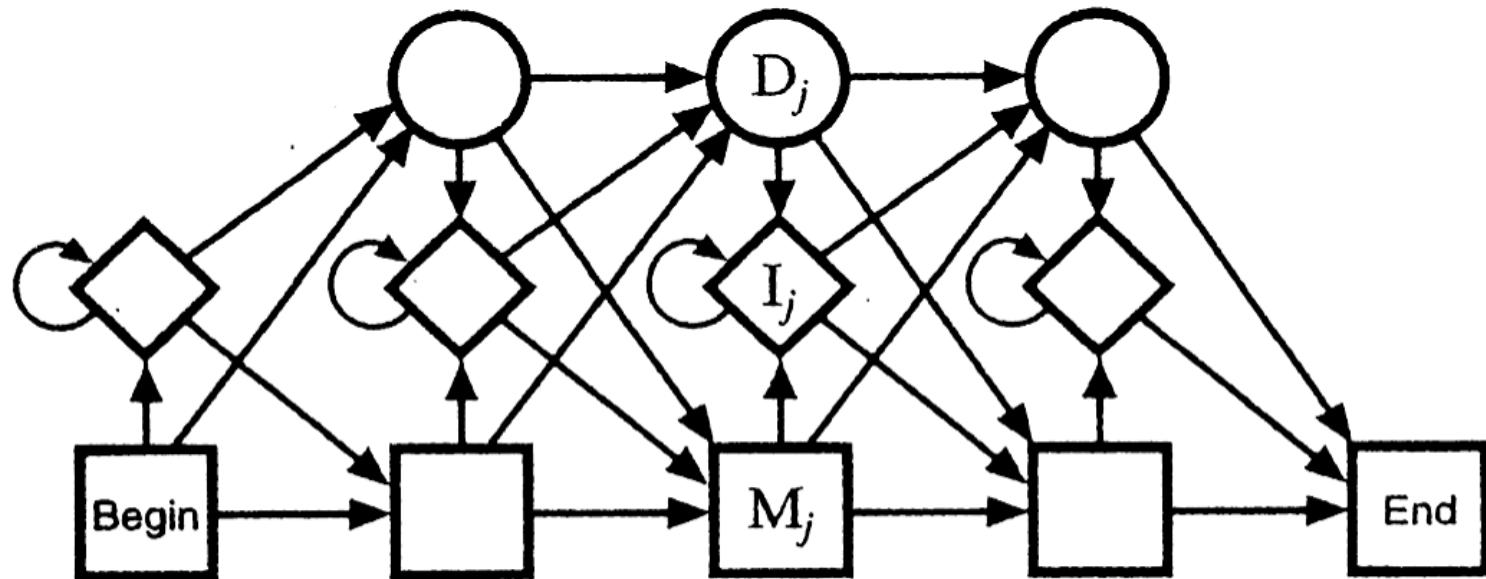


Figure 5.2 *The transition structure of a profile HMM.*

M_j : Match states (20 emission probabilities)

I_j : Insert states (Background emission probabilities)

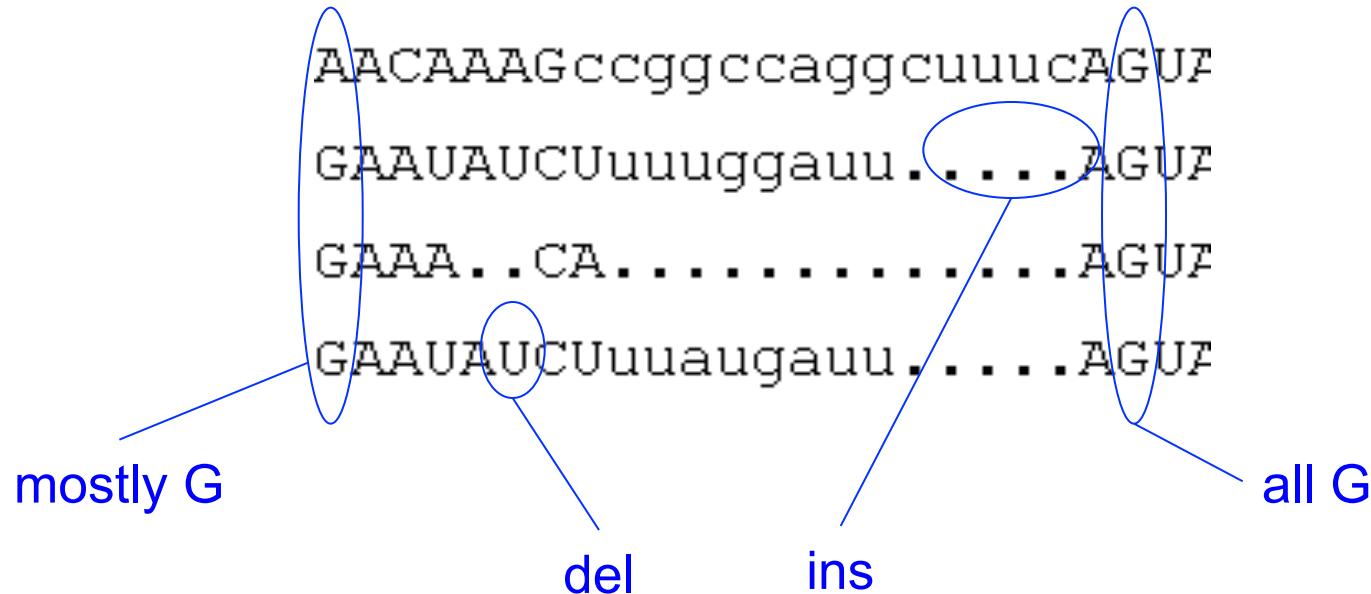
D_j : Delete states (silent - no emission)

How to model an RNA “Motif”?

Conceptually, start with a profile HMM:

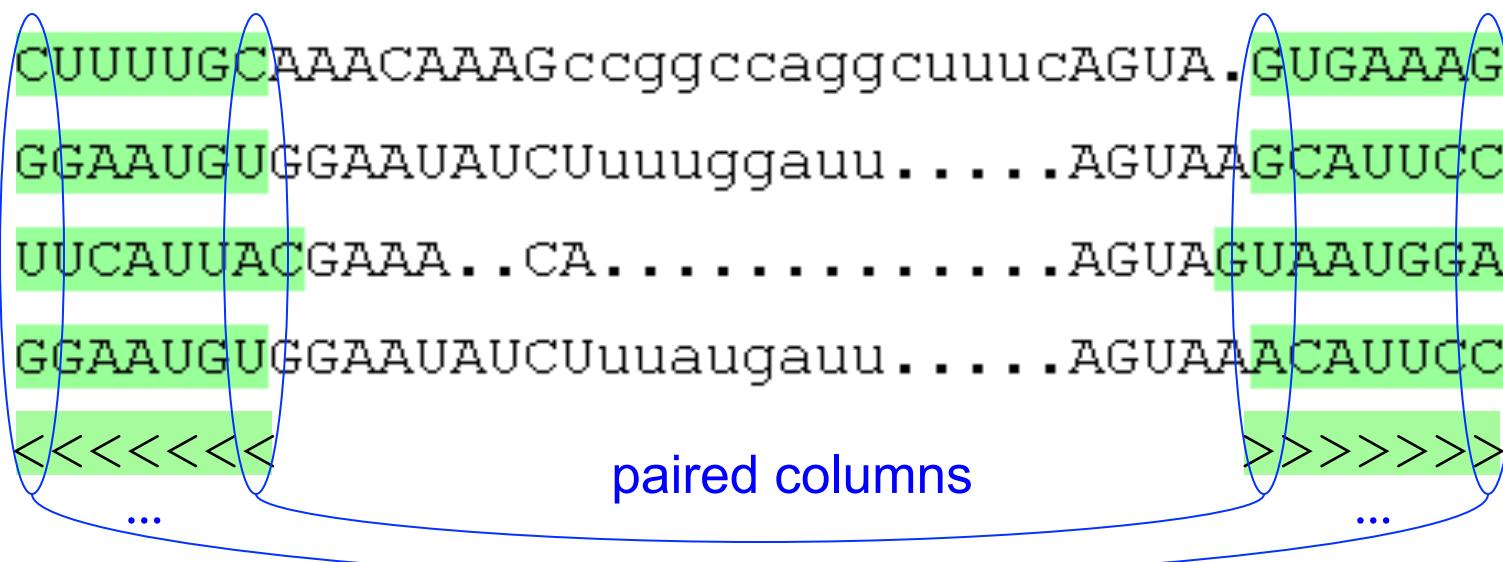
from a multiple alignment, estimate nucleotide/ insert/delete preferences for each position

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



Does not handle “paired
columns” above

Profile Hmm Structure

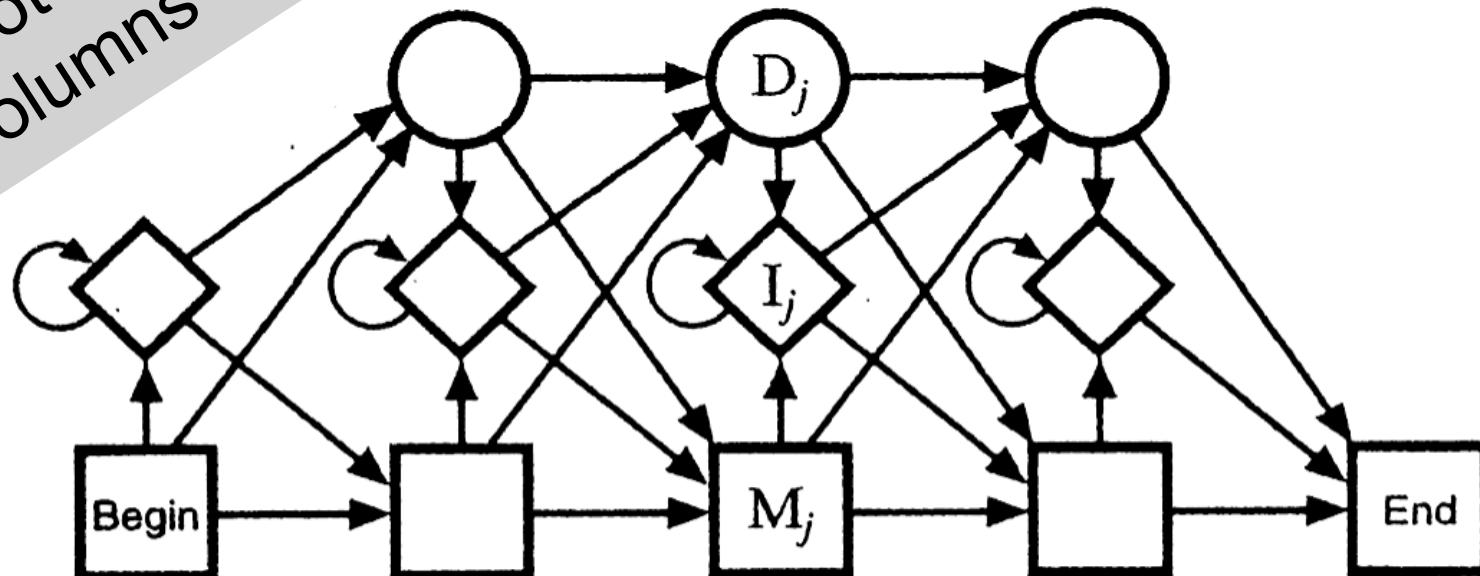


Figure 5.2 The transition structure of a profile HMM.

M_j : Match states (20 emission probabilities)

I_j : Insert states (Background emission probabilities)

D_j : Delete states (silent - no emission)

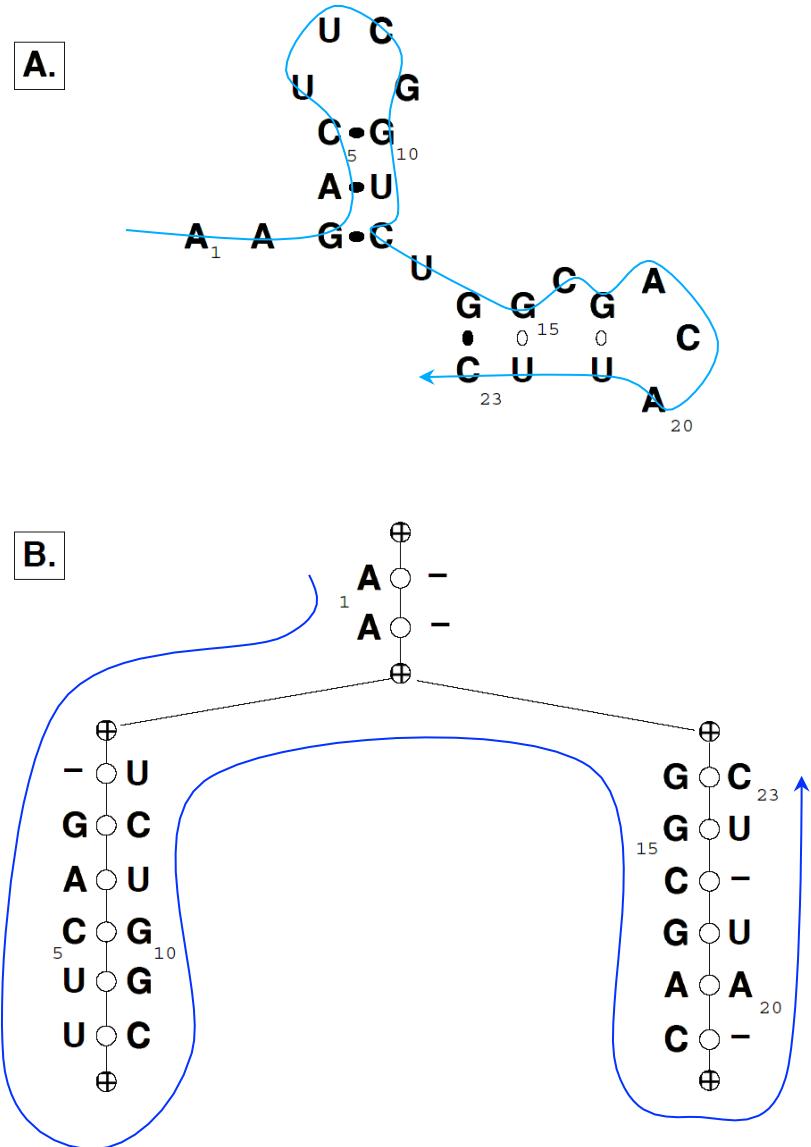
CM Structure

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)

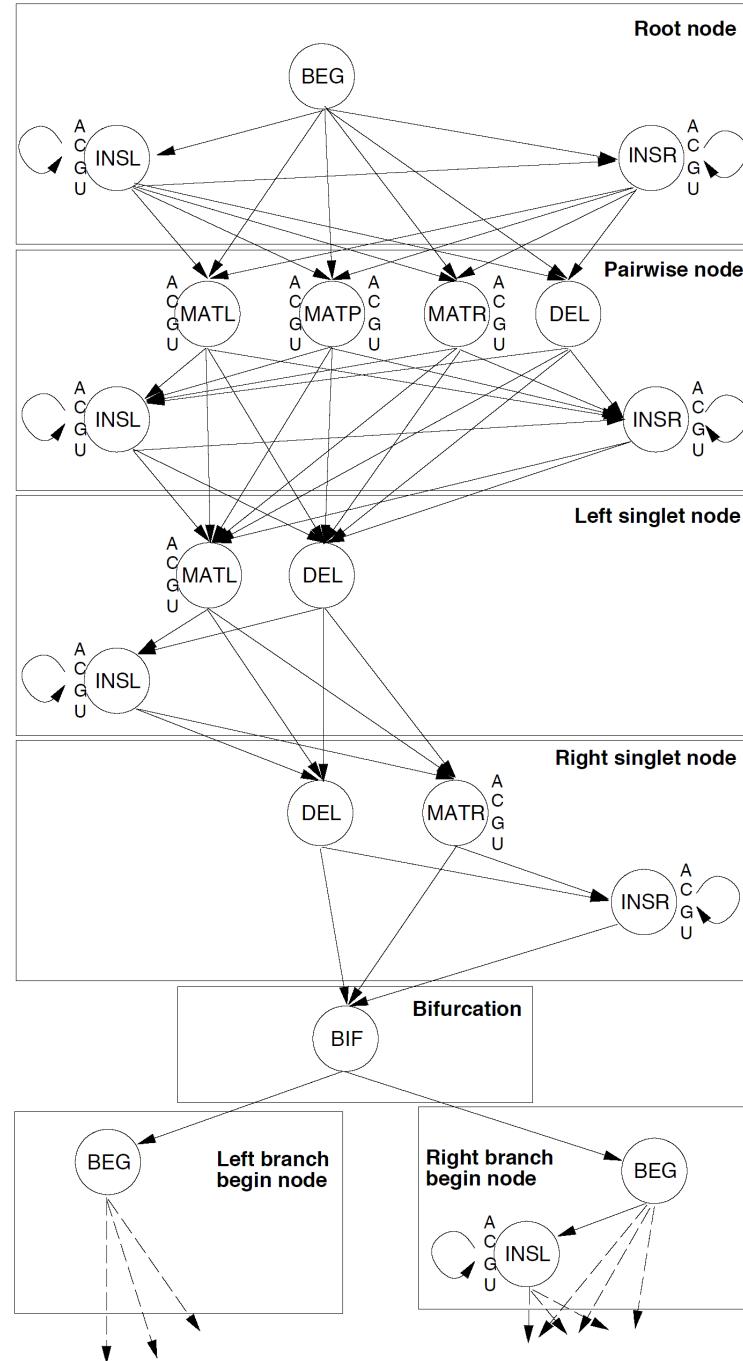


Overall CM Architecture

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

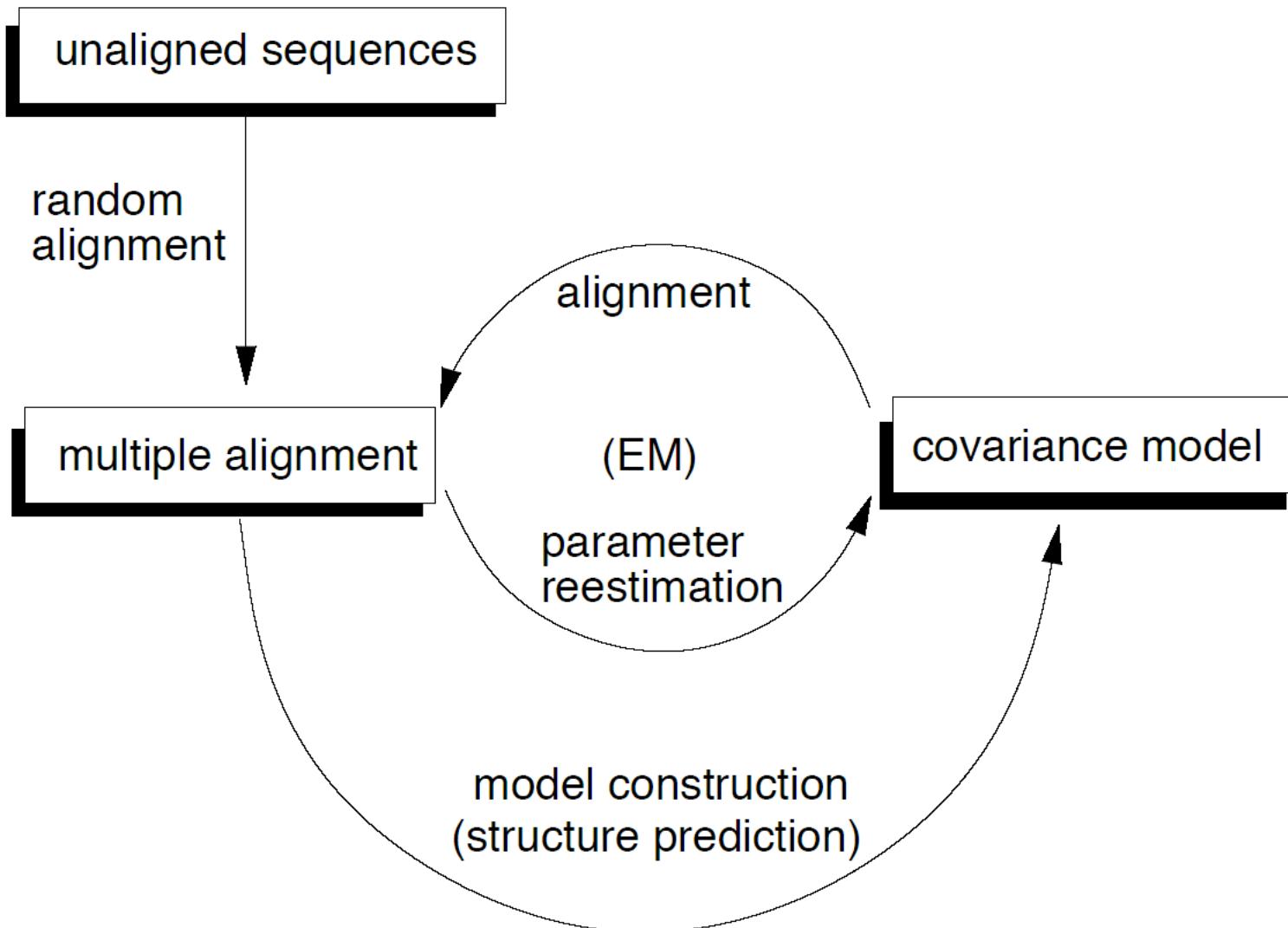
Primary vs Secondary Info

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

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

Model Training



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

Was biggest scientific comp 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

RF00037: Example Rfam Family

Input (hand-curated):

MSA “seed alignment”

SS_cons

Score Thresh T

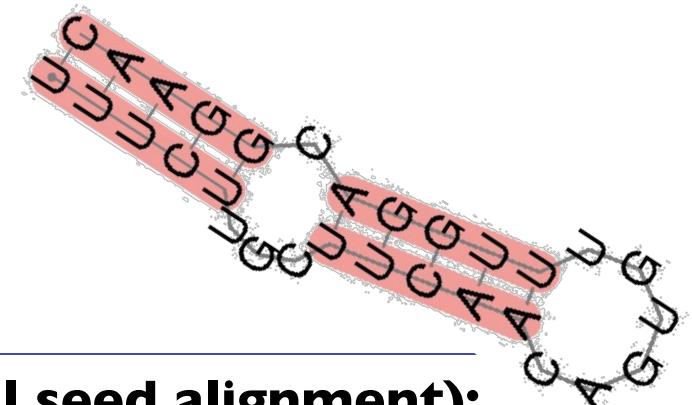
Window Len W

Output:

CM

scan results & “full
alignment”

phylogeny, etc.



IRE (partial seed alignment):

Hom. sap.	GUUCCUGCUUCAACAGUGUUUUGGAUGGAAC
Hom. sap.	UUUCUUC . UUCAACAGUGUUUUGGAUGGAAC
Hom. sap.	UUUCCUGUUUCAACAGUGCUUGGA . GGAAC
Hom. sap.	UUUAUC .. AGUGACAGAGUUUCACU . AUAAA
Hom. sap.	UCUCUUGCUUCAACAGUGUUUUGGAUGGAAC
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	<<<<...<<<<.....>>>>.>>>>

Rfam – key issues

Overly narrow families

Variant structures/unstructured RNAs

Spliced RNAs

RNA pseudogenes

Human ALU is SRP related w/ 1.1m copies

Mouse B2 repeat (350k copies) tRNA related

Speed & sensitivity

Motif discovery/hand-made models

Homology search

“Homolog” – similar by descent from common ancestor

Sequence-based

Smith-Waterman

FASTA

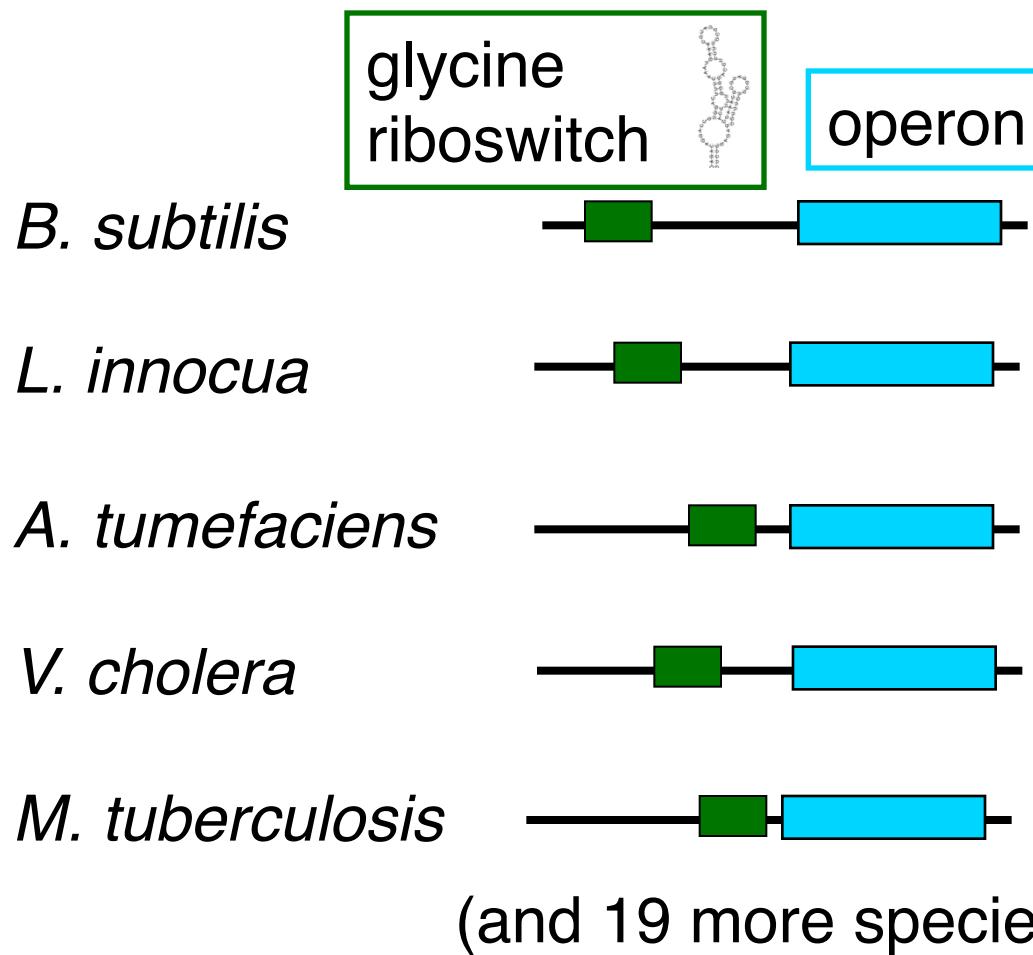
BLAST

For RNA, sharp decline in sensitivity at ~60-70% identity

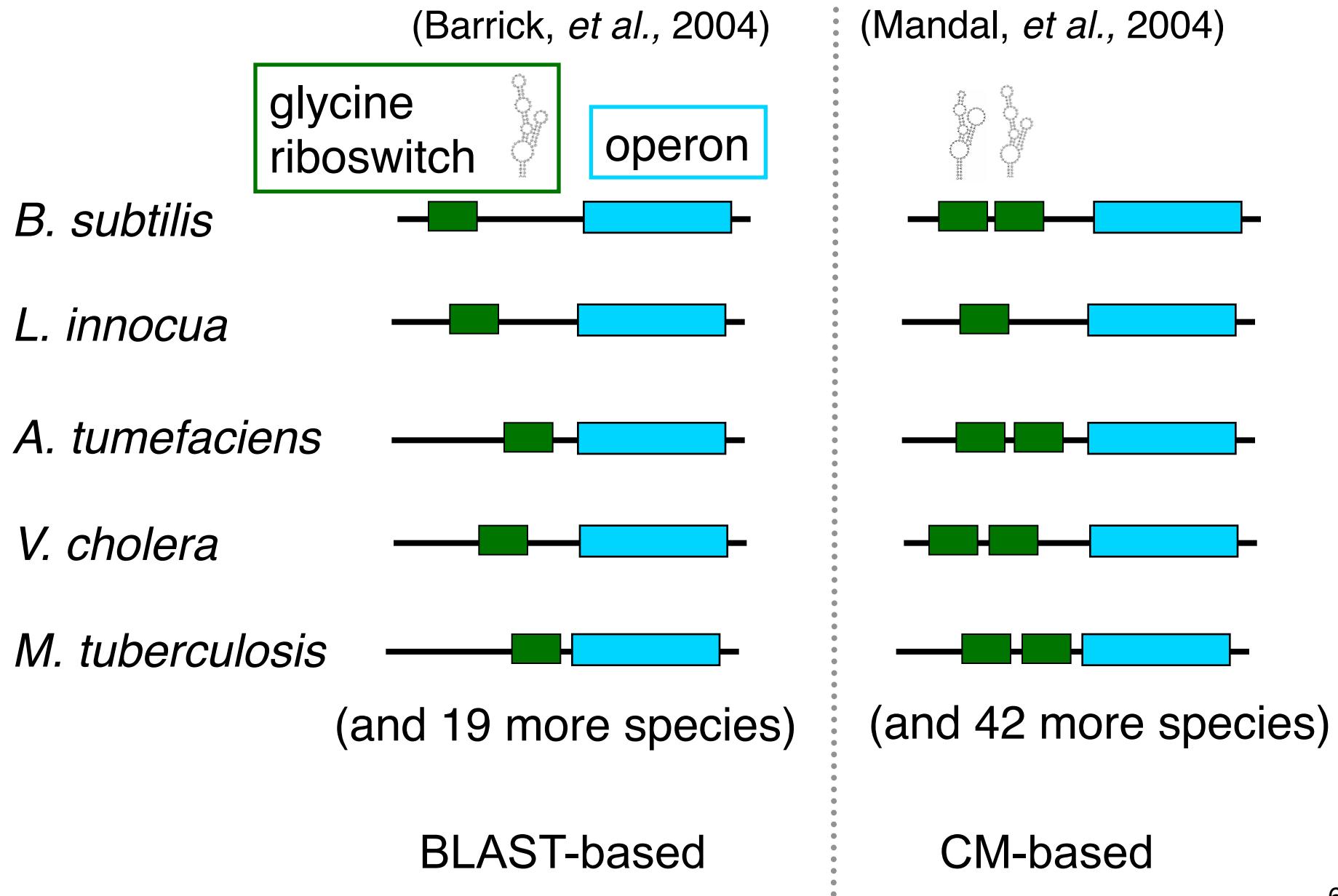
So, use structure, too

Impact of RNA homology search

(Barrick, *et al.*, 2004)



Impact of RNA homology search



Faster Genome Annotation of Non-coding RNAs Without Loss of Accuracy

Zasha Weinberg

& W.L. Ruzzo

Recomb '04, ISMB '04, Bioinfo '06

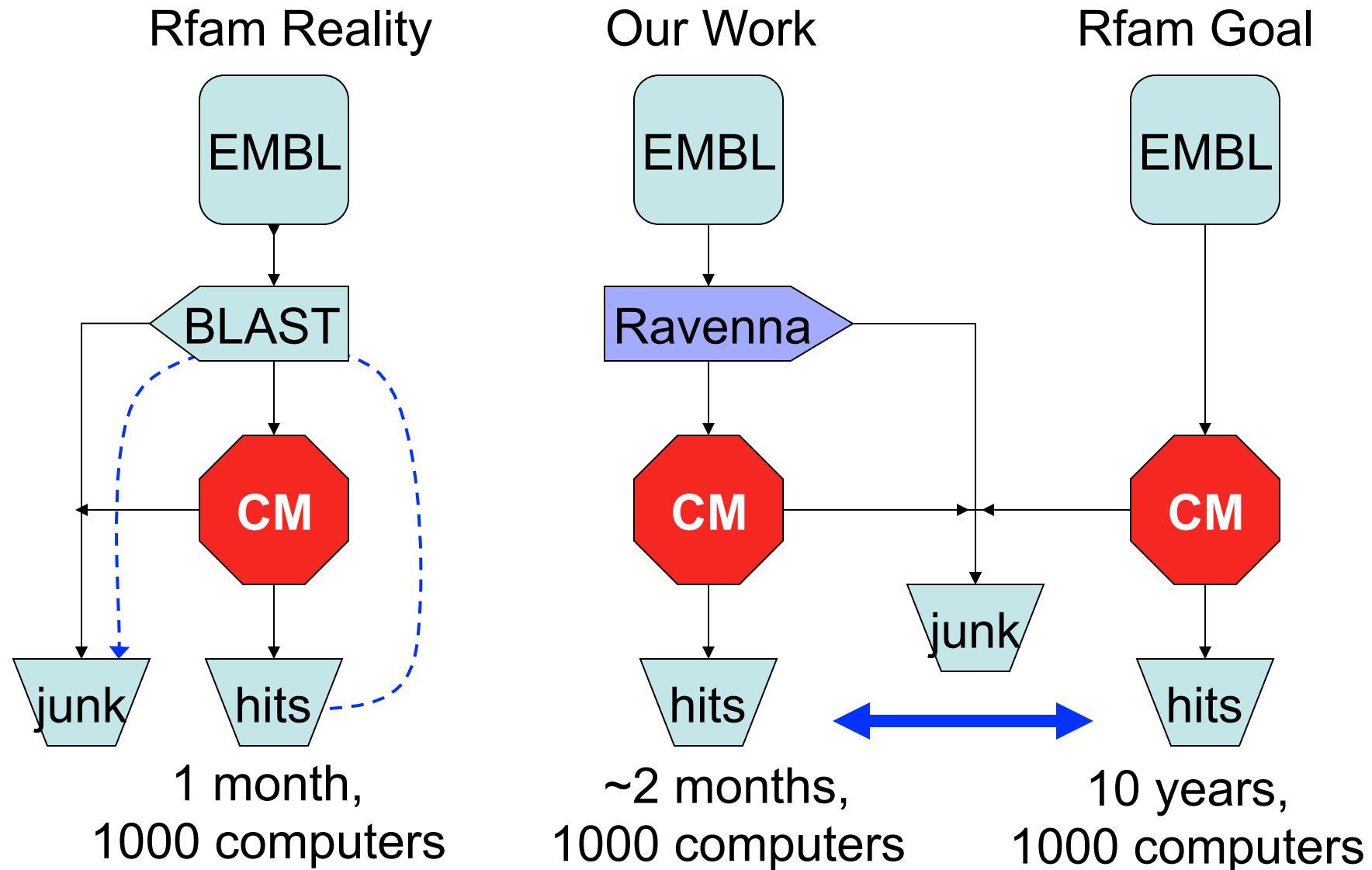
RaveNnA: Genome Scale RNA Search

Typically 100x speedup over raw CM, w/ no loss in accuracy:

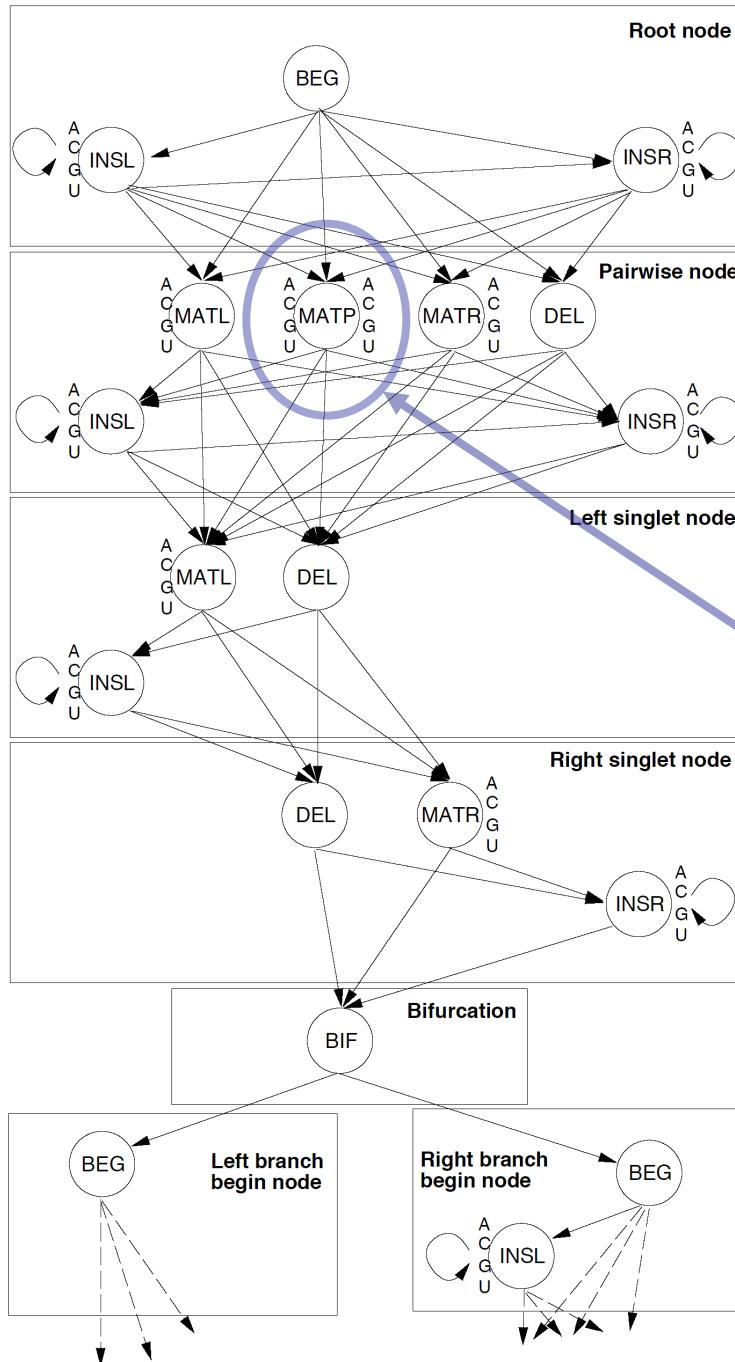
- Drop structure from CM to create a (faster) HMM
- Use that to pre-filter sequence;
- Discard parts where, provably, CM score < threshold;
- Actually run CM on the rest (the promising parts)
- Assignment of HMM transition/emission scores is key
 - (a large convex optimization problem)

Weinberg & Ruzzo, *Bioinformatics*, 2004, 2006

CM's are good, but slow

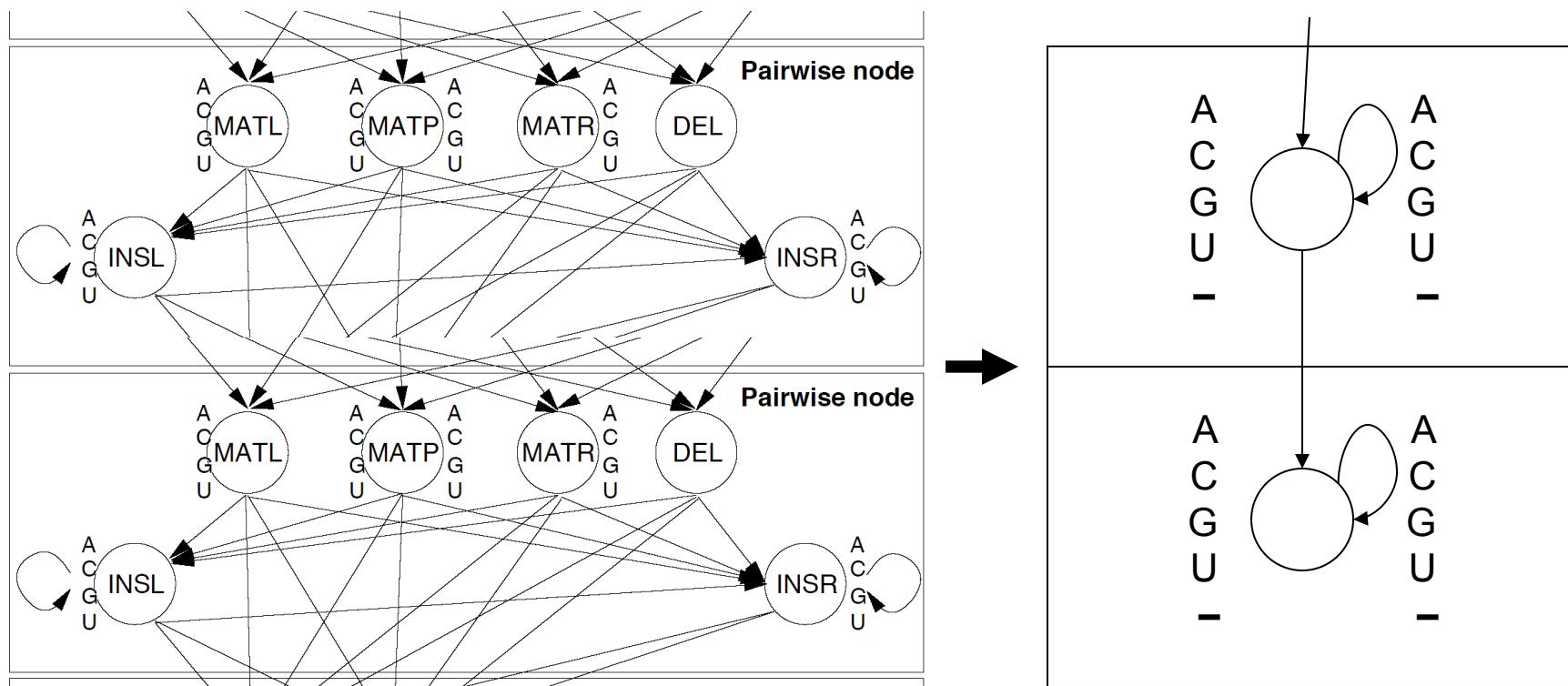


Covariance Model

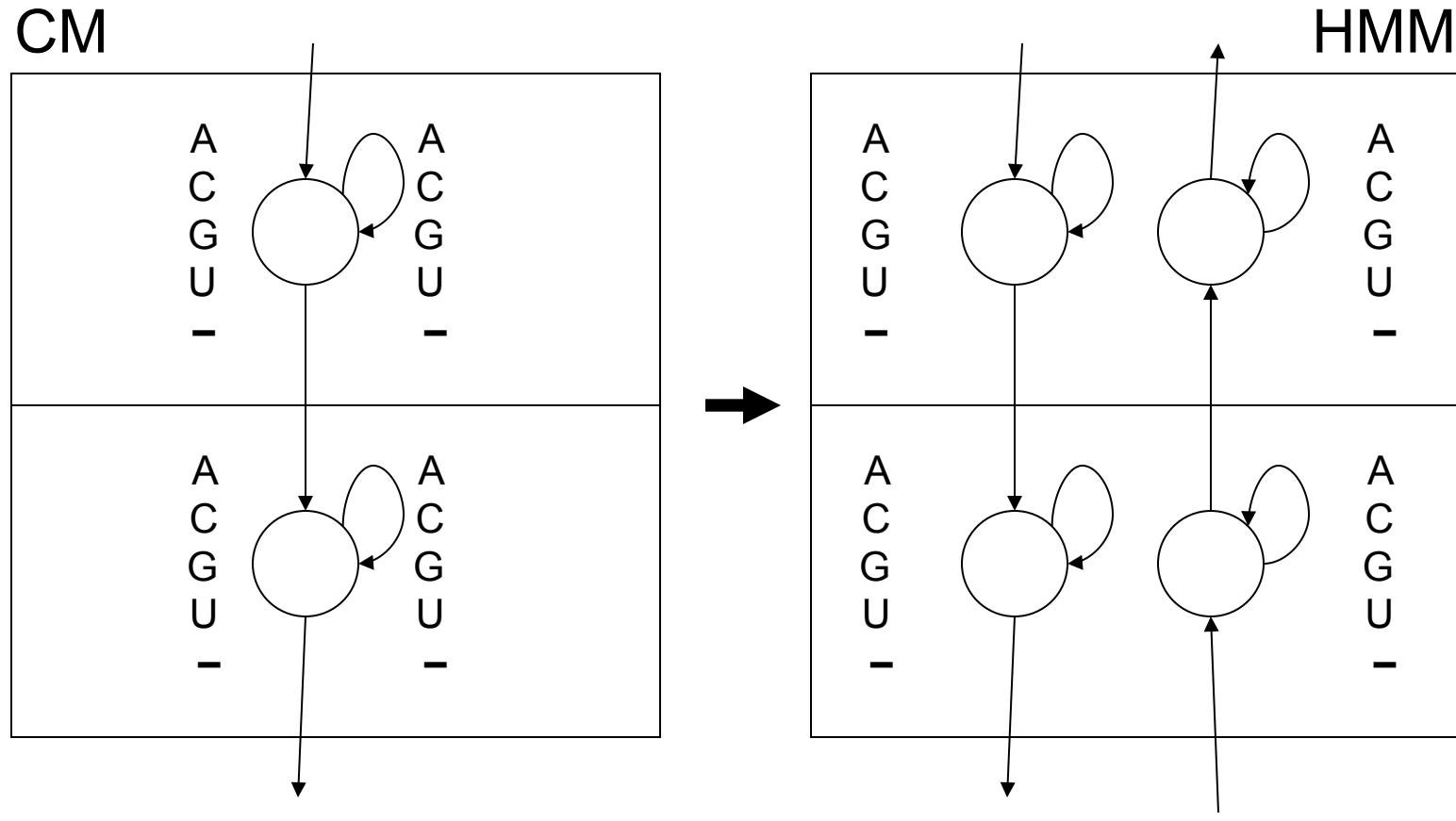


Key difference of CM vs HMM:
 Pair states emit paired symbols, corresponding to base-paired nucleotides; 16 emission probabilities here.

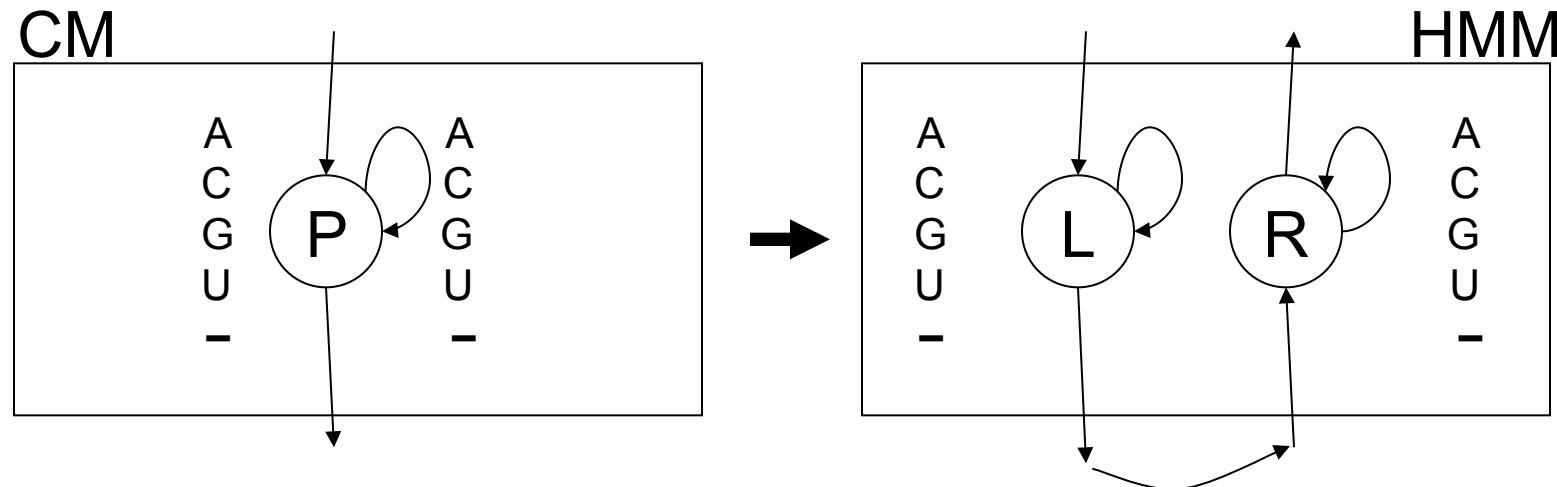
Oversimplified CM (for pedagogical purposes only)



CM to HMM



Key Issue: 25 scores → 10



Need: $\log \text{Viterbi scores } \text{CM} \leq \text{HMM}$

Viterbi/Forward Scoring

Path π defines transitions/emissions

$\text{Score}(\pi) = \text{product of "probabilities" on } \pi$

NB: ok if “probs” aren’t, e.g. $\sum \neq 1$

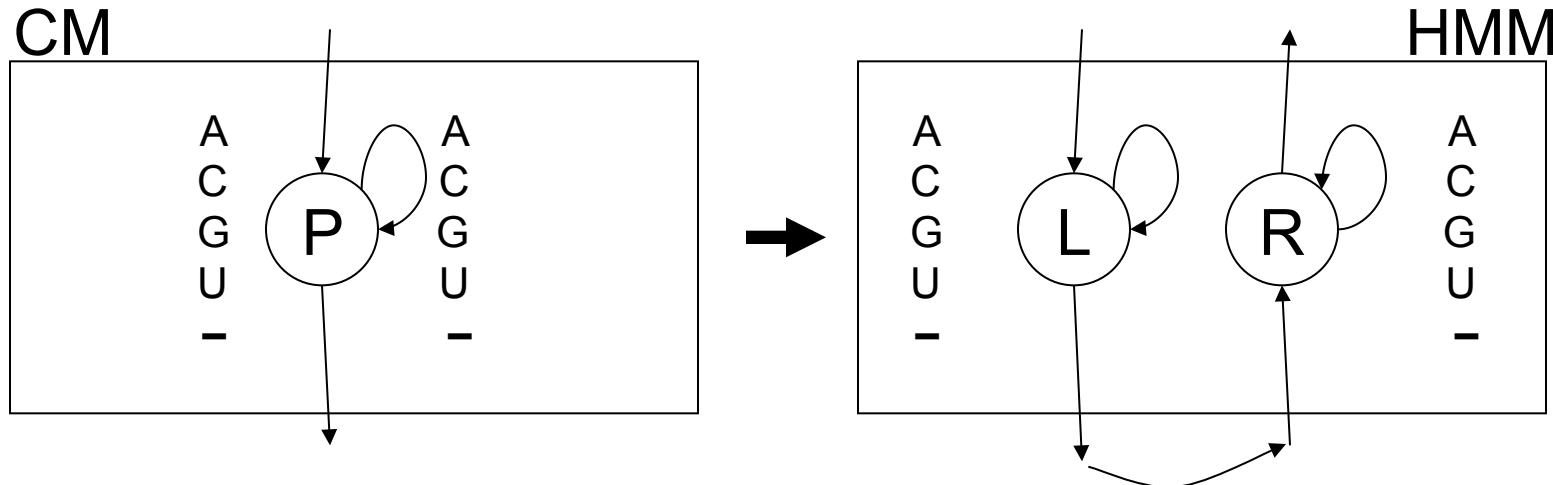
(e.g. in CM, emissions are odds ratios vs
0th-order background)

For any nucleotide sequence x :

$\text{Viterbi-score}(x) = \max\{ \text{score}(\pi) \mid \pi \text{ emits } x\}$

$\text{Forward-score}(x) = \sum\{ \text{score}(\pi) \mid \pi \text{ emits } x\}$

Key Issue: 25 scores → 10



Need: $\log \text{Viterbi scores } \text{CM} \leq \text{HMM}$

$$P_{AA} \leq L_A + R_A$$

$$P_{AC} \leq L_A + R_C$$

$$P_{AG} \leq L_A + R_G$$

$$P_{AU} \leq L_A + R_U$$

$$P_{A-} \leq L_A + R_-$$

$$P_{CA} \leq L_C + R_A \quad \dots$$

$$P_{CC} \leq L_C + R_C \quad \dots$$

$$P_{CG} \leq L_C + R_G \quad \dots$$

$$P_{CU} \leq L_C + R_U \quad \dots$$

$$P_{C-} \leq L_C + R_- \quad \dots$$

NB: HMM not a prob. model

$$\begin{aligned}P_{AA} &\leq L_A + R_A \\P_{AC} &\leq L_A + R_C \\P_{AG} &\leq L_A + R_G \\P_{AU} &\leq L_A + R_U \\P_{A-} &\leq L_A + R_- \\&\dots\end{aligned}$$

Rigorous Filtering

Any scores satisfying the linear inequalities give rigorous filtering

Proof:

CM Viterbi path score

\leq “corresponding” HMM path score

\leq Viterbi HMM path score

(even if it does not correspond to *any* CM path)

Some scores filter better

$$P_{UA} = I \leq L_U + R_A$$

$$P_{UG} = 4 \leq L_U + R_G$$

Option 1:

$$L_U = R_A = R_G = 2$$

Option 2:

$$L_U = 0, R_A = 1, R_G = 4$$

Assuming ACGU ≈ 25%

Opt 1:

$$L_U + (R_A + R_G)/2 = 4$$

Opt 2:

$$L_U + (R_A + R_G)/2 = 2.5$$

Optimizing filtering

For any nucleotide sequence x :

$$\text{Viterbi-score}(x) = \max\{ \text{score}(\pi) \mid \pi \text{ emits } x \}$$

$$\text{Forward-score}(x) = \sum\{ \text{score}(\pi) \mid \pi \text{ emits } x \}$$

Expected Forward Score

$$E(L_i, R_i) = \sum_{\text{all sequences } x} \text{Forward-score}(x) * \Pr(x)$$

NB: E is a function of L_i, R_i only

Under 0th-order background model

Optimization:

Minimize $E(L_i, R_i)$ subject to score Lin.Ineq.s

This is heuristic (“forward $\downarrow \Rightarrow$ Viterbi $\downarrow \Rightarrow$ filter \downarrow ”)

But still rigorous because “subject to score Lin.Ineq.s”

Calculating $E(L_i, R_i)$

$$E(L_i, R_i) = \sum_x \text{Forward-score}(x) * \text{Pr}(x)$$

Forward-like: for every state, calculate expected score for all paths ending there; easily calculated from expected scores of predecessors & transition/emission probabilities/scores

Minimizing $E(L_i, R_i)$

Calculate $E(L_i, R_i)$
symbolically, in terms of
emission scores, so we
can do partial derivatives
for numerical convex
optimization algorithm

Forward:

$$\begin{aligned} f_k(i) &= P(x_1 \dots x_i, \pi_i = k) \\ f_l(i+1) &= e_l(x_{i+1}) \sum_k f_k(i) a_{k,l} \end{aligned}$$

Viterbi:

$$v_l(i+1) = e_l(x_{i+1}) \cdot \max_k (v_k(i) a_{k,l})$$

$$\frac{\partial E(L_1, L_2, \dots)}{\partial L_i}$$

Assignment of scores/ “probabilities”

Convex optimization problem

Constraints: enforce rigorous property

Objective function: filter as aggressively as possible

Problem sizes:

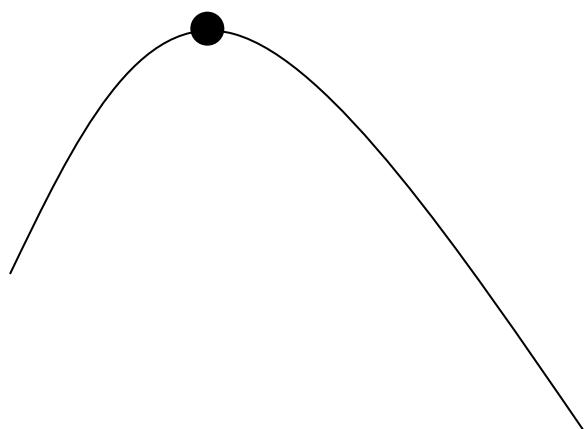
1000-10000 variables

10000-100000 inequality constraints

“Convex” Optimization

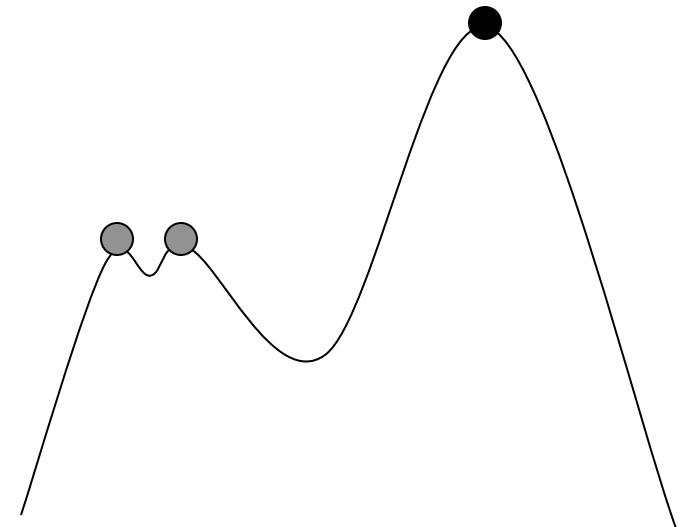
Convex:

local max = global max;
simple “hill climbing” works



Nonconvex:

can be many local maxima,
≪ global max;
“hill-climbing” fails



Estimated Filtering Efficiency

(139 Rfam 4.0 families)

Filtering fraction	# families (compact)	# families (expanded)
$< 10^{-4}$	105	110
$10^{-4} - 10^{-2}$	8	17
.01 - .10	11	3
.10 - .25	2	2
.25 - .99	6	4
.99 - 1.0	7	3

break even → ~100x speedup

Averages 283 times faster than CM

Results: new ncRNAs (?)

Name	# Known (BLAST + CM)	# New (rigorous filter + CM)
<i>Pyrococcus</i> snoRNA	57	123
Iron response element	201	121
Histone 3' element	1004	102*
Retron msr	11	48
Hammerhead I	167	26
Hammerhead III	251	13
U6 snRNA	1462	2
U7 snRNA	312	1
cobalamin riboswitch	170	7
13 other families	5-1107	0

Results: With additional work

	# with BLAST+CM	# with rigorous filter series + CM	# new
Rfam tRNA	58609	63767	5158
Group II intron	5708	6039	331
tRNAscan-SE (human)	608	729	121
tmRNA	226	247	21
Lysine riboswitch	60	71	11
And more...			

“Additional work”

Profile HMM filters use no 2^{ary} structure info

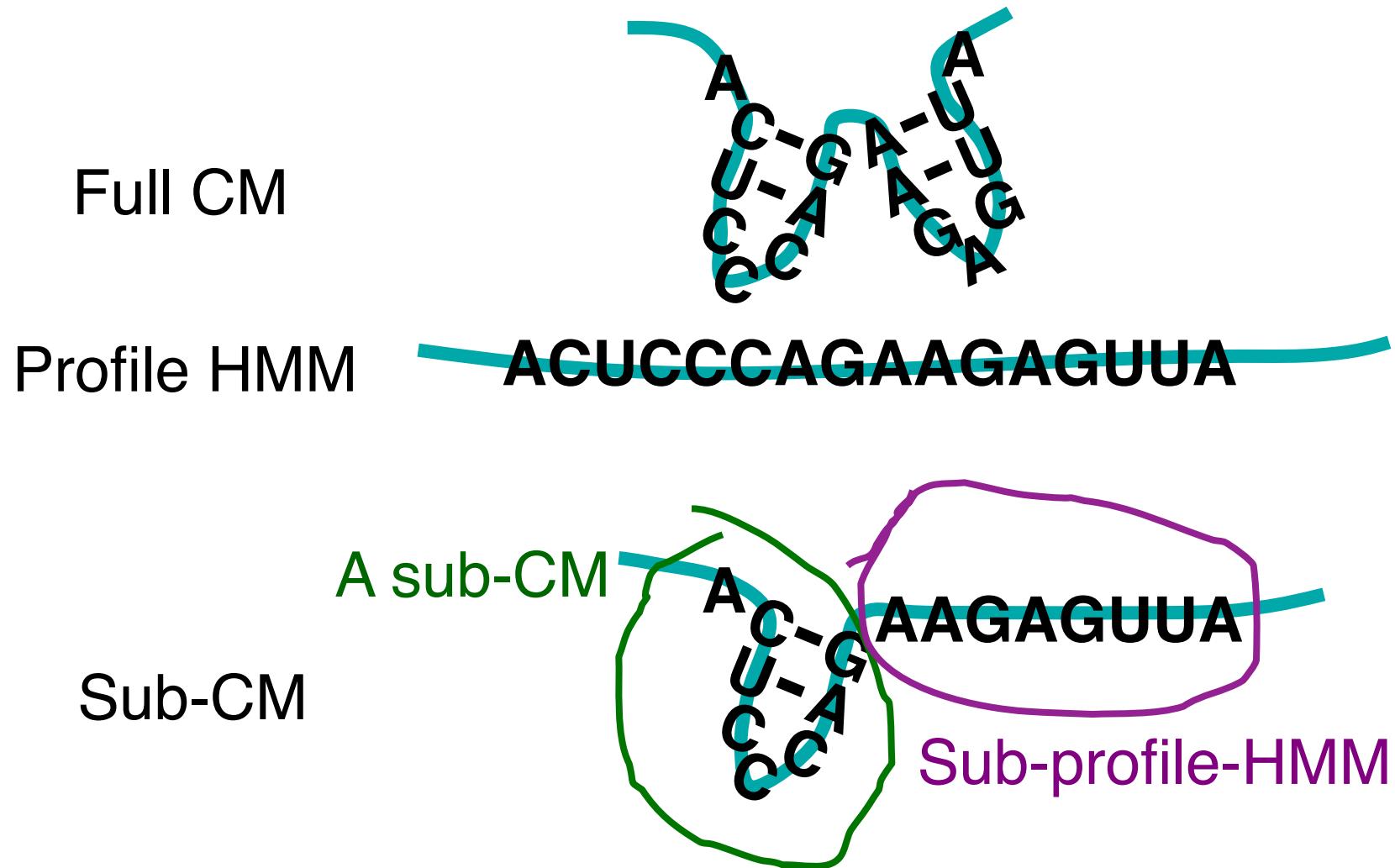
They work well because, tho structure can be critical to function, there is (usually) enough primary sequence conservation to exclude most of DB

But not on all families (and may get worse?)

Can we exploit some structure (quickly)?

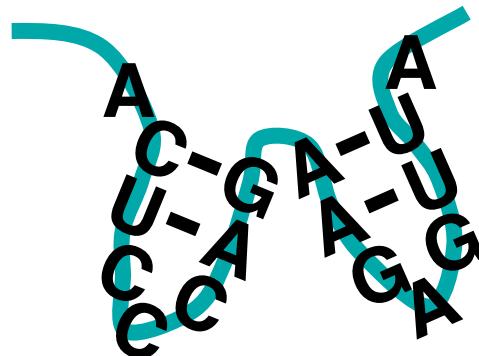
- Idea 1: “sub-CM”
 - Idea 2: extra HMM states remember mate
 - Idea 3: try lots of combinations of “some hairpins”
 - Idea 4: chain together several filters (select via Dijkstra)
- } for some hairpins

Sub-CM filters

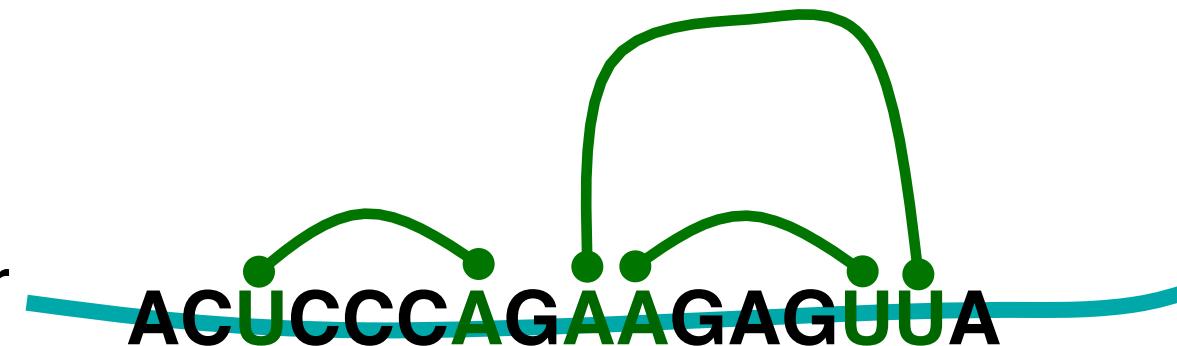


Store-pair filters

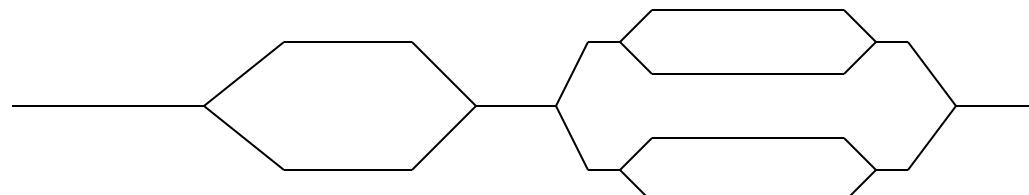
Full CM



Store pair



“Profile” HMM:



ACCGAT
GGACA

Filter Chains

Rigorous filter

Rigorous filter

Rigorous filter

CM



ncRNAs



Why run filters in series?

	Filtering fraction	Run time (sec/Kbase)
Filter 1	0.25	1
Filter 2	0.01	10
CM	N/A	200

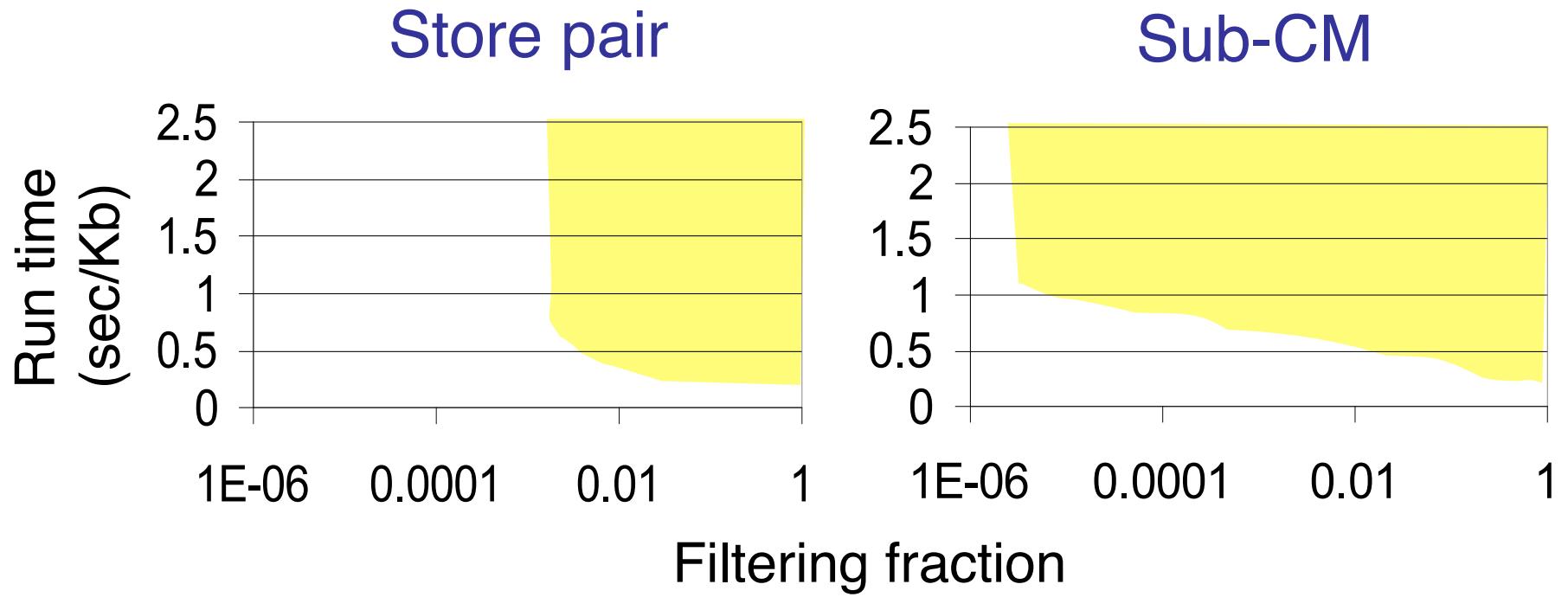
CM alone: 200 s/Kb

Filter 1 → CM: $1 + 0.25 * 200 = 51$ s/Kb

Filter 2 → CM: $10 + 0.01 * 200 = 12$ s/Kb

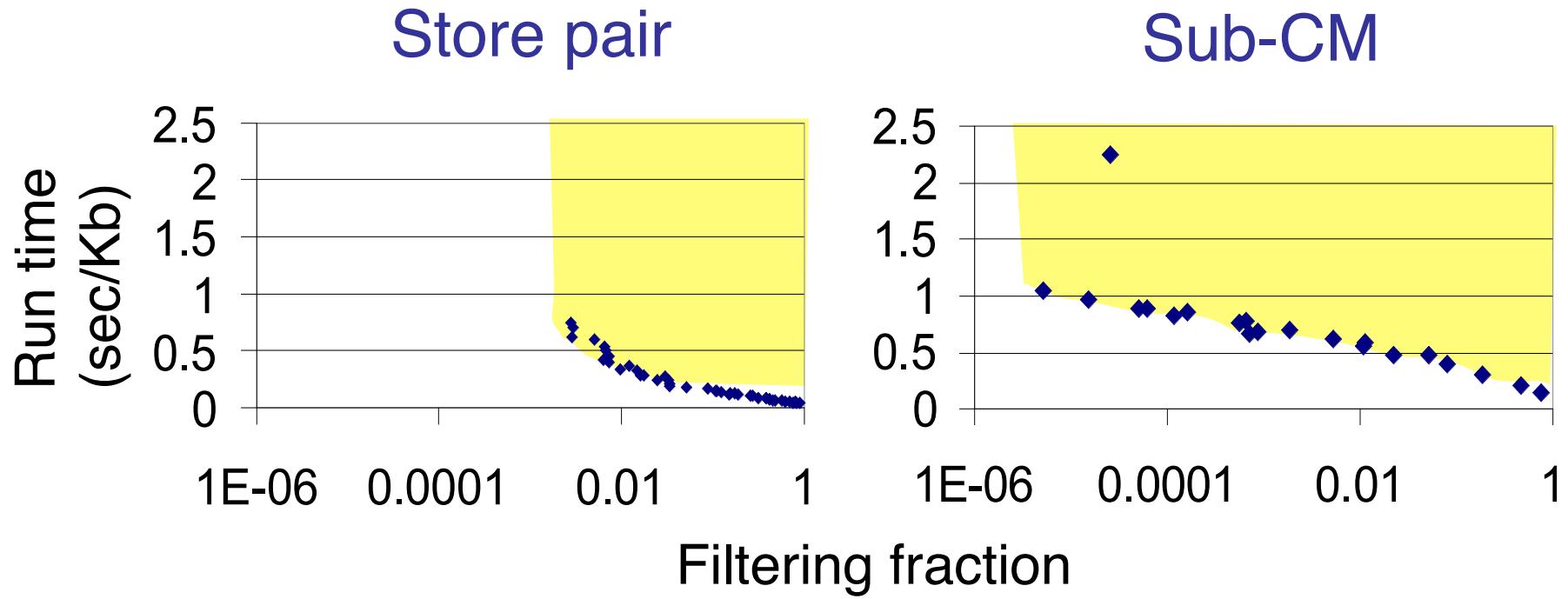
Filter 1 → Filter 2 → CM:

$1 + 0.25 * 10 + 0.01 * 200 = 5.5$ s/Kb

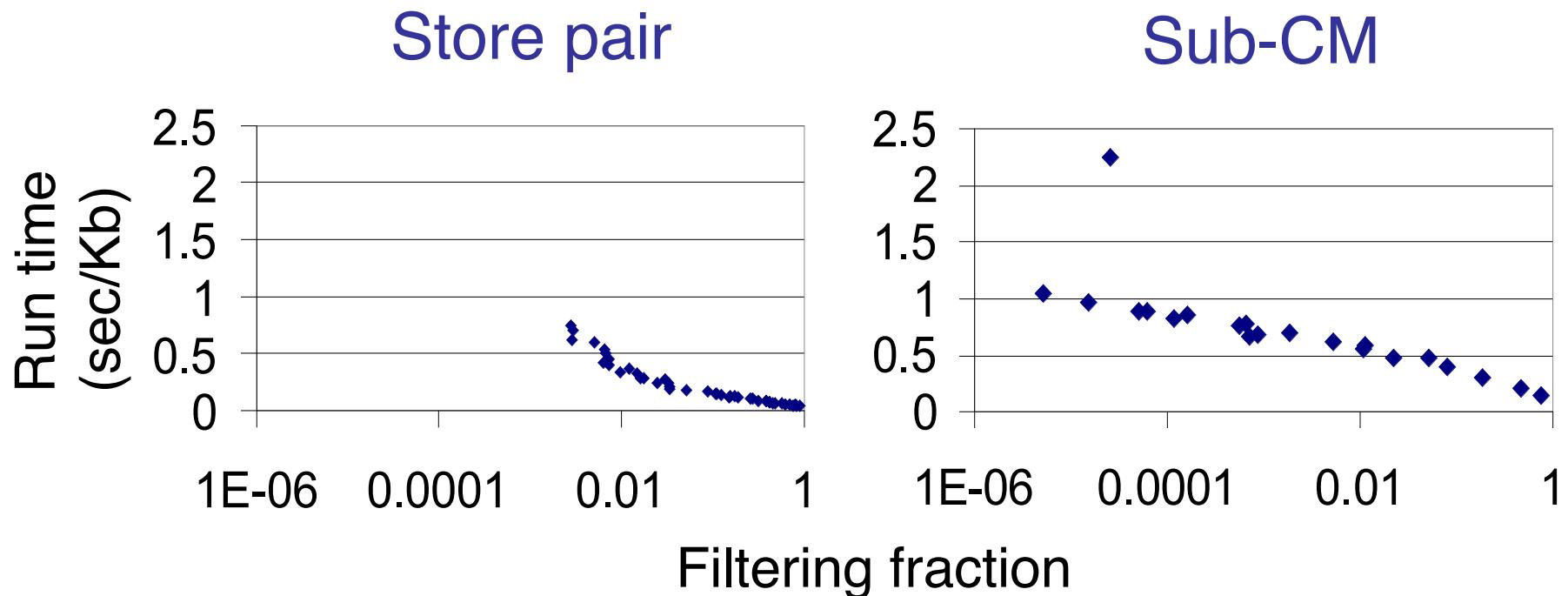


Properties of a filter:

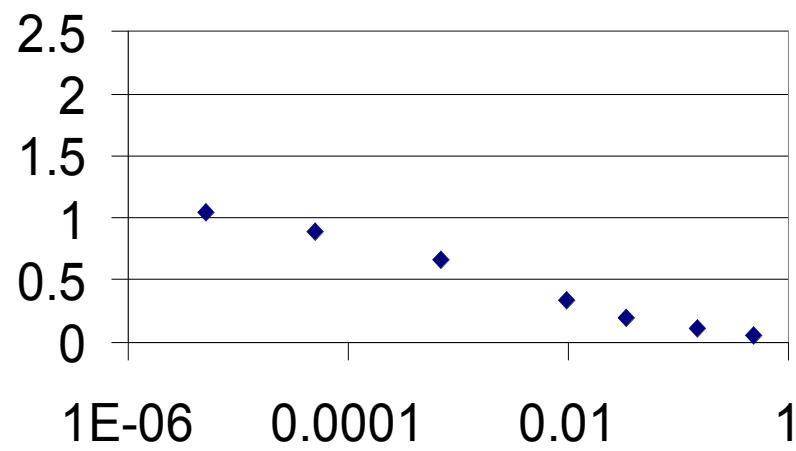
- Filtering fraction
- Run time (sec/Kb)



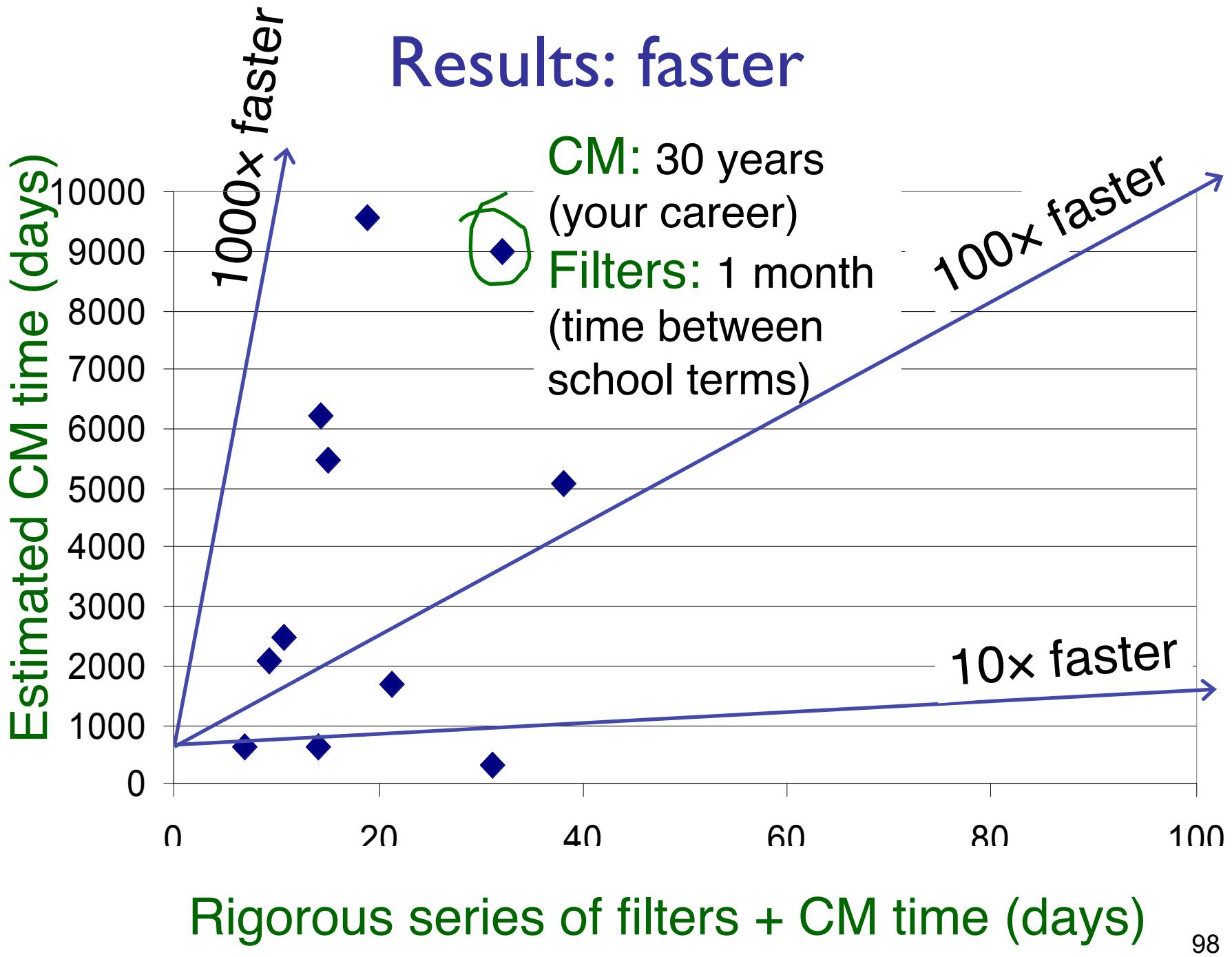
- Simplified performance model (selectivity & speed)
- Independence assumptions for base pairs
- Use dynamic programming to rapidly explore base pair combinations



Selected rigorous filter chain



Results: faster



Results: more sensitive than BLAST

	# with BLAST+CM	# with rigorous filters + CM	# new
Rfam tRNA	58609	63767	5158
Group II intron	5708	6039	331
Iron response element	201	322	121
tmRNA	226	247	21
Lysine riboswitch	60	71	11
And more...			

Is there anything more to do?

Rigorous filters can be too cautious

E.g., 10 times slower than heuristic filters

Yet only 1-3% more sensitive

We want to

Run scans faster with minimal loss of sensitivity

Know empirically what sensitivity we're losing

Heuristic Filters

Rigorous filters optimized for worst case

Possible to trade improved speed for small loss in sensitivity?

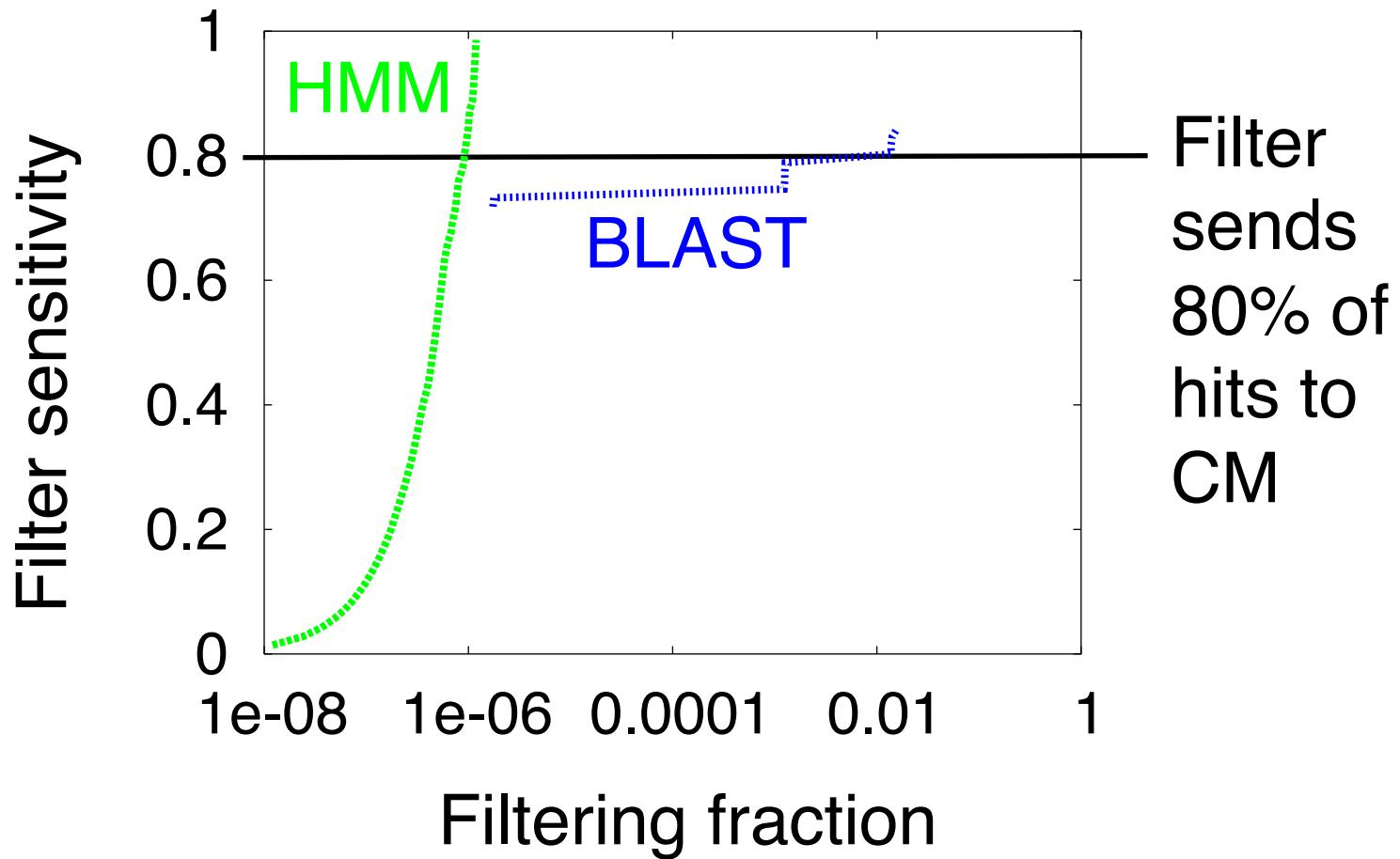
Yes – profile HMMs as before, but optimized for average case

Often 10x faster, modest loss in sensitivity

Heuristic Filters

ROC-like curves

(lysine riboswitch)



Heuristic Filters

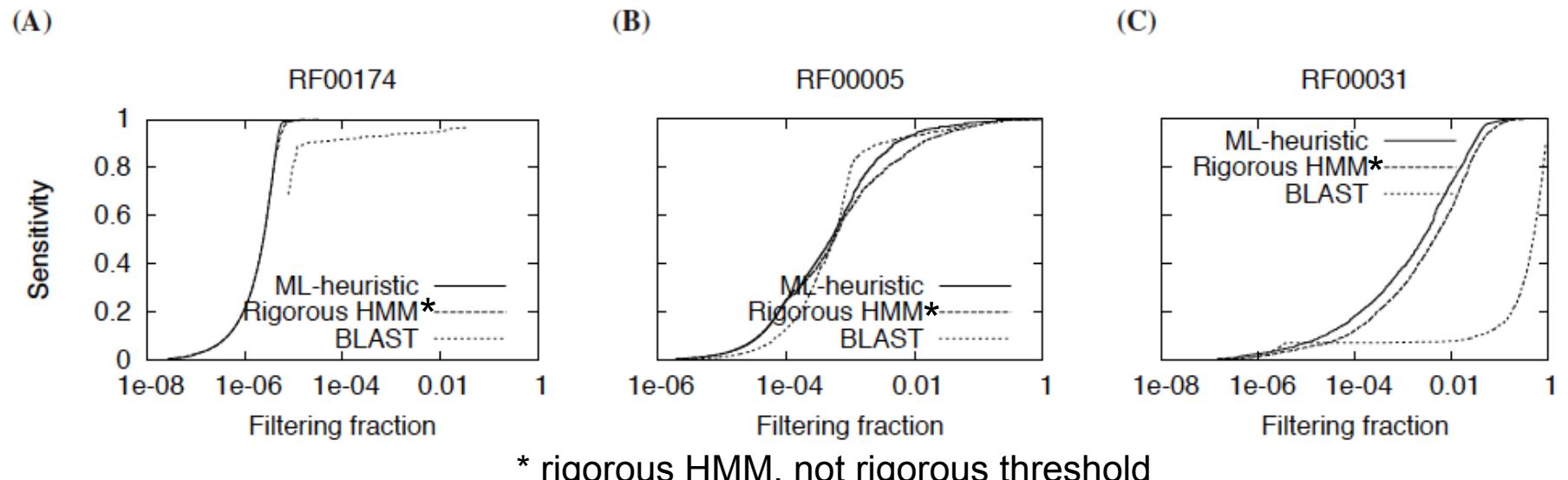
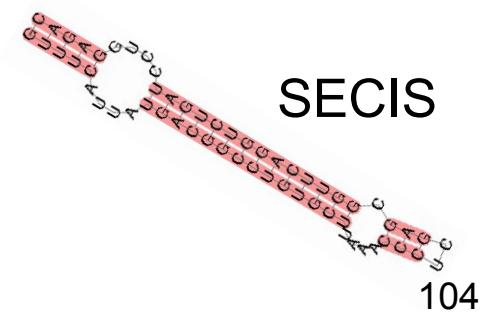
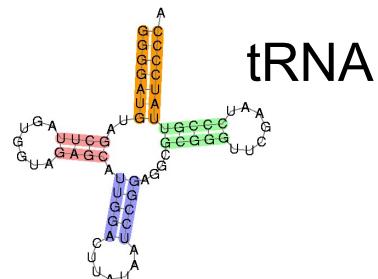
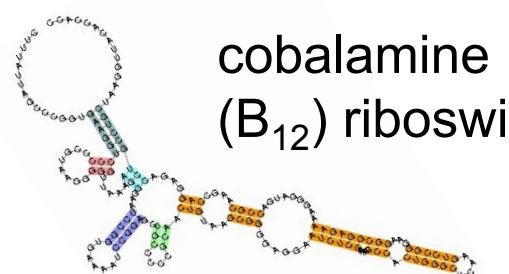


Fig. 1. Selected ROC-like curves. All plot sensitivity against filtering fraction, with filtering fraction in log scale. (A) RF00174 is typical of the other families; the ML-heuristic is slightly better than the rigorous profile HMM, and both often dramatically exceed BLAST. (B) Atypically, in RF00005, BLAST is superior, although only in one region. (C) BLAST performs especially poorly for RF00031. (Recall that rigorous scans were not possible for RF00031, so only $\sim 90\%$ of hits are known; see text.) The supplement includes all ROC-like curves, and the inferior ignore-SS.



Software

Ravenna implements both rigorous and heuristic filters

Infernal (engine behind Rfam) implements heuristic filters and some other (important) accelerations

E,g., dynamic “banding” of dynamic programming matrix based on the insight that large deviations from consensus length must have low scores.

CM Search Summary

Still slower than we might like, but dramatic speedup over raw CM is possible with:

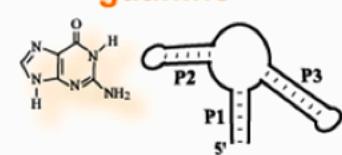
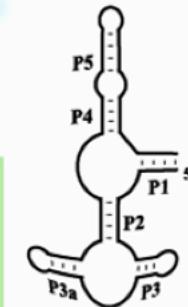
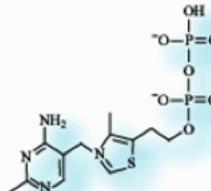
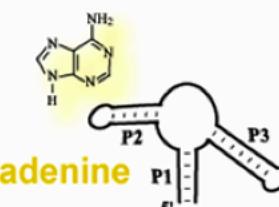
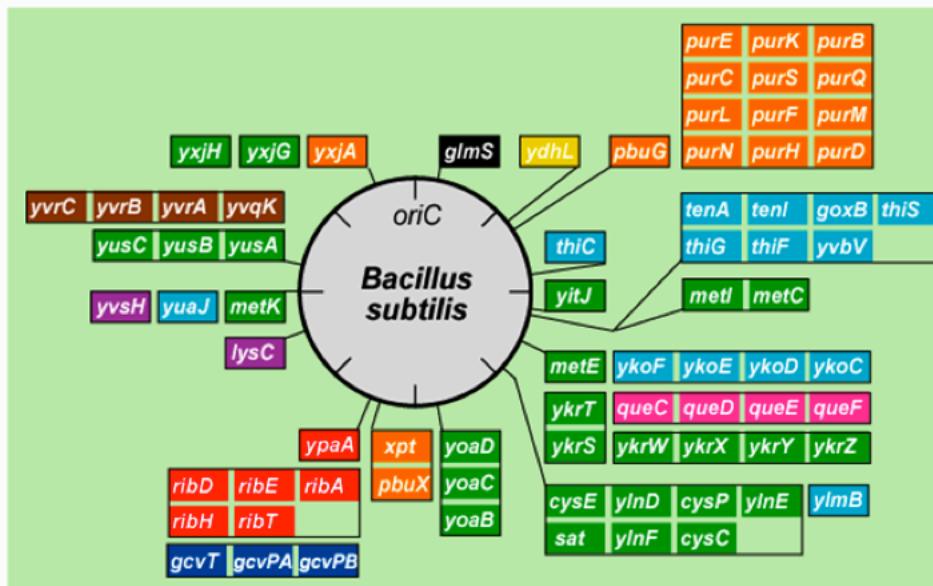
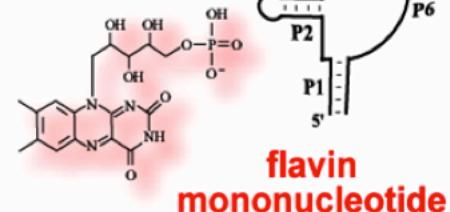
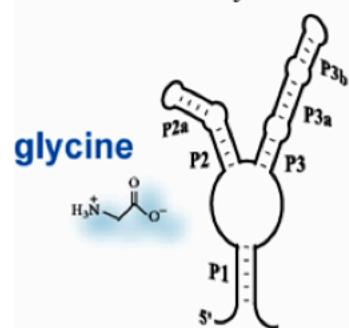
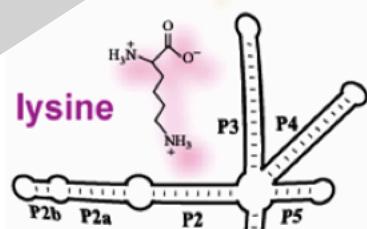
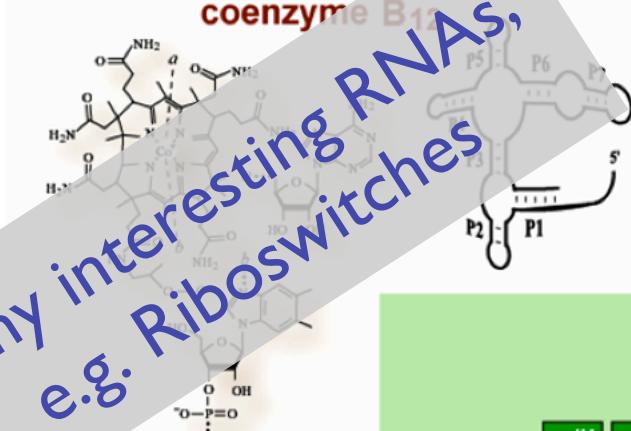
- No loss in sensitivity (provably), or

- Even faster with modest (and estimable) loss in sensitivity

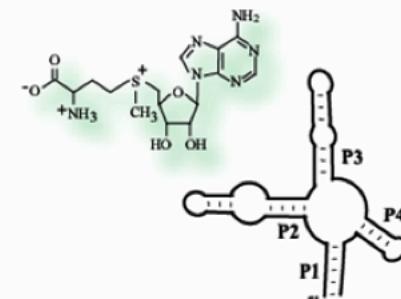
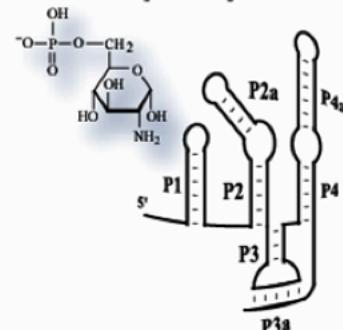
Last Lecture

Part I

Many interesting RNAs,
e.g. Riboswitches



glucosamine-6-phosphate



Nussinov: Structure Prediction

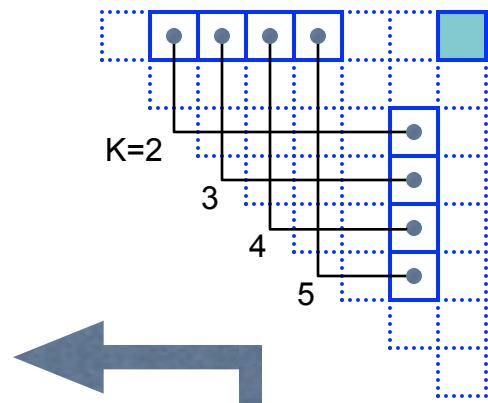
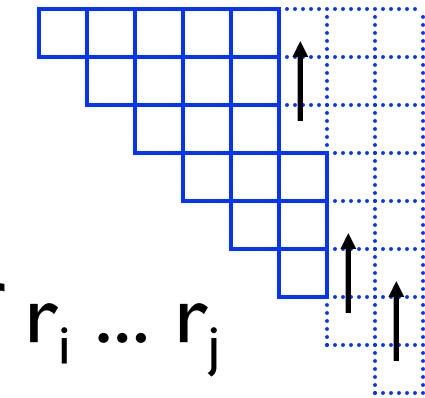
Computation Order

$B(i,j)$ = **# 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:

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



Time: $O(n^3)$

Loop-based energy version is better; recurrences similar, slightly messier

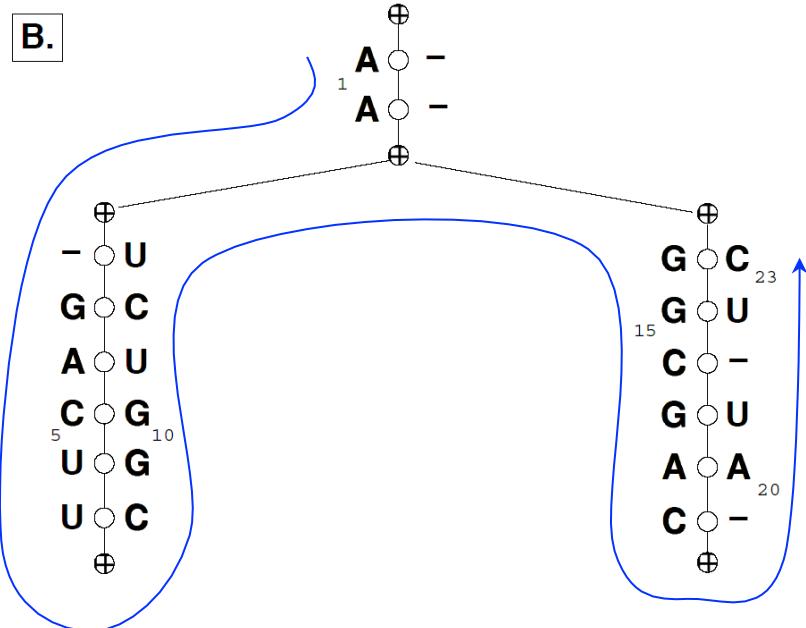
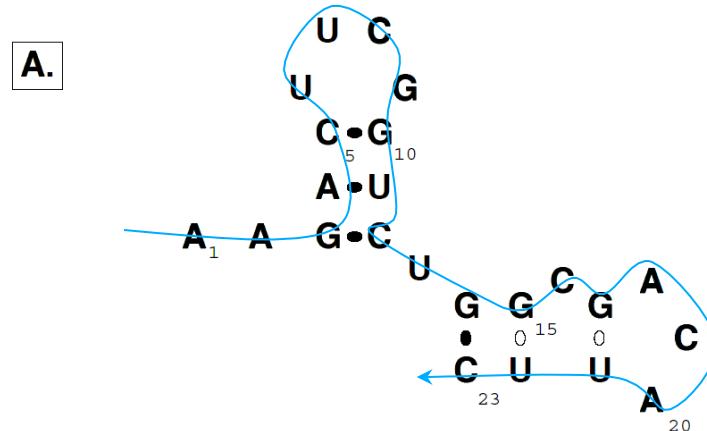
CM's for desc/search Structure

A: Sequence + structure

B: the CM “guide tree”

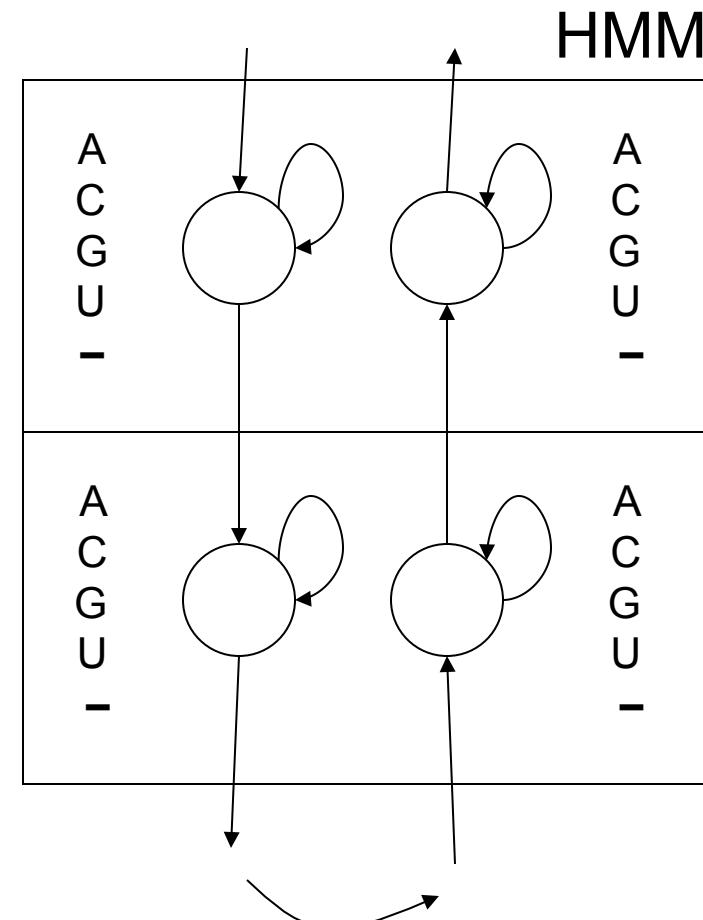
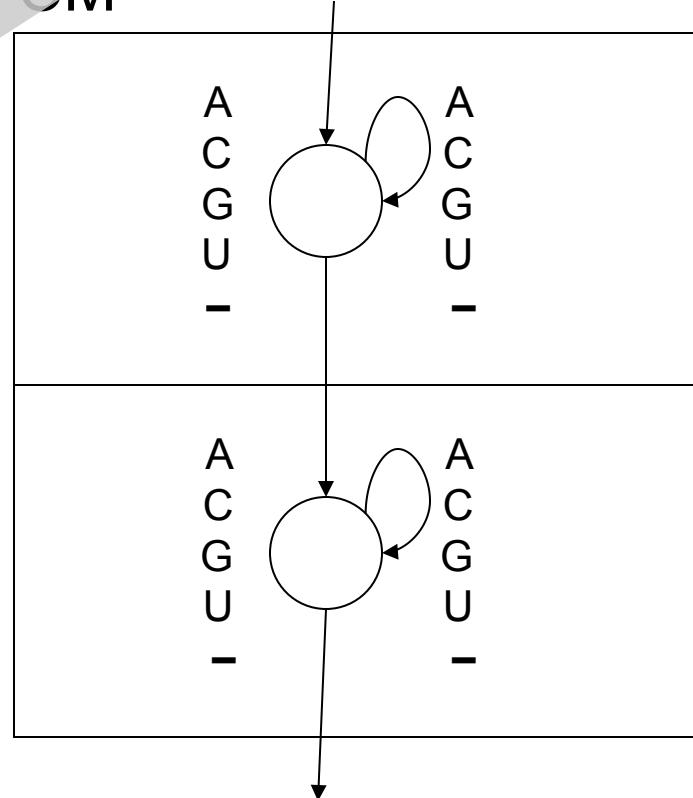
C: probabilities of letters/ pairs & of indels

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



Accelerating search

CM to HMM



25 emissions per state

5 emissions per state, 2x states

Motif Discovery

RNA Motif Discovery

CM's are great, but where do they come from?

Key approach: comparative genomics

Search for motifs with common secondary structure in a set of functionally related sequences.

Challenges

Three related tasks

Locate the motif regions.

Align the motif instances.

Predict the consensus secondary structure.

Motif search space is huge!

Motif location space, alignment space, structure space.

RNA Motif Discovery

Would be great if: given 100 complete genomes from diverse species, we could automatically find all the RNAs.

State of the art: that's hopeless

Hope: can we exploit biological knowledge to narrow the search space?

RNA Motif Discovery

More promising problem: given a 10-20 unaligned sequences of a few kb, most of which contain instances of one RNA motif of 100-200bp -- find it.

Example: 5' UTRs of orthologous glycine cleavage genes from γ -proteobacteria

Example: corresponding introns of orthologous vertebrate genes

Orthologs =
counterparts in
different species

Approaches

Align-First: Align sequences, then look for common structure

Fold-First: Predict structures, then try to align them

Joint: Do both together

Pitfall for sequence-alignment-first approach

Structural conservation \neq Sequence conservation

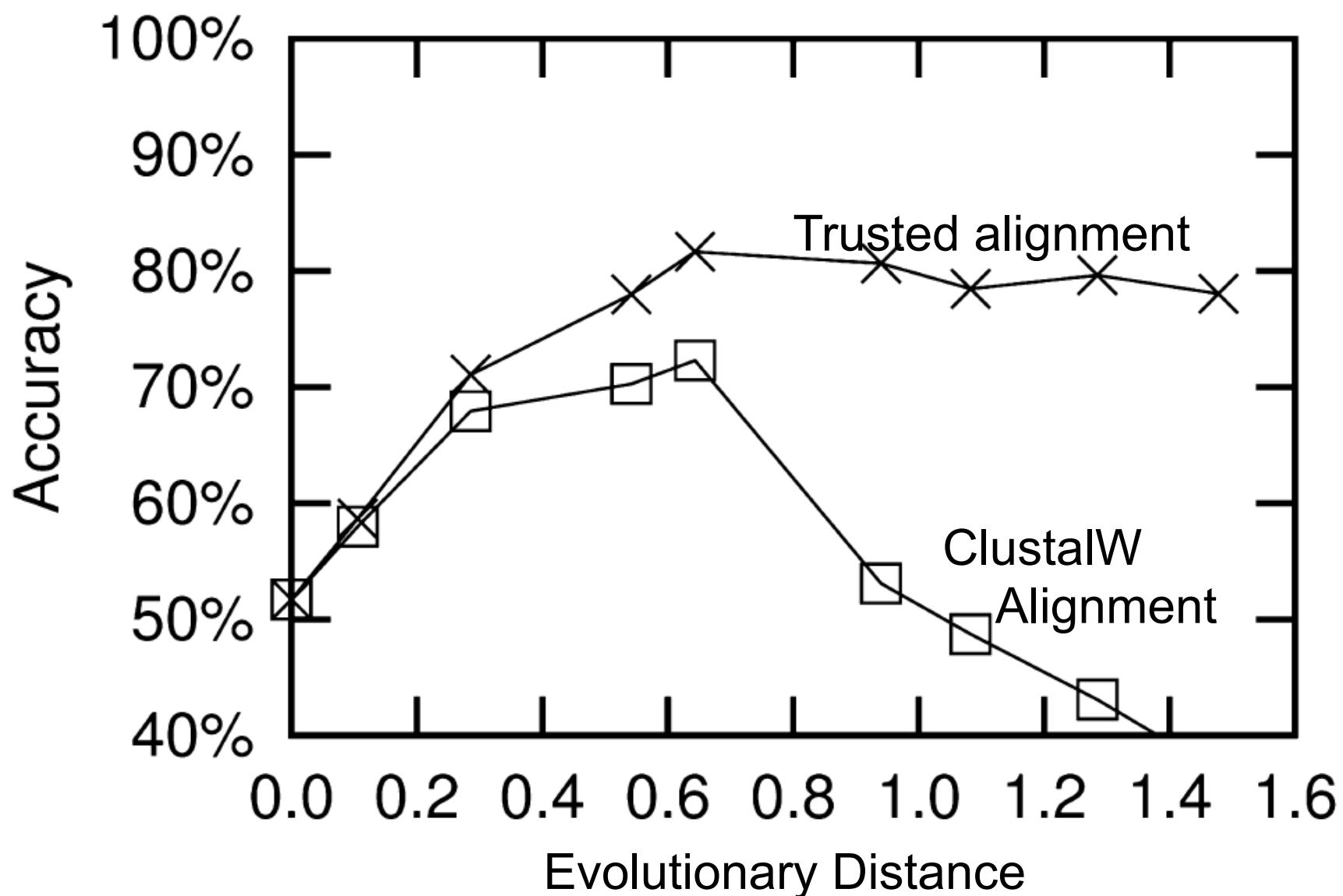
Alignment without structure information is unreliable

CLUSTALW alignment of SECIS elements with flanking regions

```
-----CCCCCCCCAGGCCCTGGTGCAGG--ATGATGACGACCTGGGTG-GAA-A---CCTACCCCTGTGGGCACCC-ATGTCGA-CCCCCCTGGCATT
GGGATCATTCAGCAAGAGCAGCGTG--ACTGACATTA---TGAAGGGCTGTACTGAAGACAGCAA--GCTGTTAGTACAGACC--AGATG---CTTCTTGGCAGGCTCGTTGTAACCTCTGGAAAACCTCAAT
AGGTTTGATTAATGAGGATTACACAGAAAACCTT-TGTAAGGGTTGTGATCTGCTAA- TTGCAAAATTTTATTTTAAAT---ATTCTTACAGAAGAGTTCATTAAAGAATGTTGTGTATAGG
AGTGTGCGGATGATAACTACTGACGAAAGCTCATCGACTCAGTTAGTGGGTGATGTAGTCACATTAGTTGCCCTCCCCCATCTTG---TCTCCCTGGCAAGGAGAATATGCCGACATGATGCTAAGAG
TGGACTGATAGGTA-GCCATGGC- TTCATCTGTC--ATG--TCTGCTTCTTTATATTG-TGTATGATGGTCACAGTGAAAG---TTCCCACAGCTGTGACTTGATTTAA-AAATGTCGGAAGA
TAAACTCGAACCTGGAGCGGGCAATTGCTGATTACGA-TTAACCACGTATCCCTGGGTCGCTGC- TTCTGGCCGTCGCTGGTTCCA-----TTTATCAACTATTAGCTCCAATACATAGCTACAGGTTTTT
AAATTCTCGCTATATGACGATGCCAATCTAAATGT-TCATTGGTTGCCATTGATGAAATCAGTTTGTTGACCTGCAAGAATTGTTGTTACCTTGTCAATTTCATTGAA-ACCACTTCTCAGA
GGGGCGGGAGTACAAGGTGCGTGTGACTGGAGCCA--CCCACTCCGACTCTGCAGGTGGTGTG- CAAATGACGACCGATTGAAATG---GTCTCACGGCCAAAACCTCGTGTCCGACATCAACCCCTTC
TTCTCCAGTGTCTAGTTACATTGATGAGAACAGAA-ACATAAAACTATGACCTAGGGGTTCT- GTTGGATAGCTCTAAATTAAAGAACGGAGAAAGAACAAACAAAGACATATTTCAGTTTTTTCTTAC
CAAACGTGATGGATA-GCCATTGGTATTCTATCTATT--TTAACTCTGTGCTTCTACATTG-TTTATGATGGCCACAGCCTAAA---TACACACGGCTGTGACTTGATTCAAA-GAA-----
TGAGCAACTTGTCT-GATGACTGGGAAAGGAGGAC--CTGCAACCCTGACTTGGTCTCTG- TTAATGACGTCTCCCTCTAA-A---CCC-CATTAAGGACTGGGAGAGGCAGA-GCAAGCCTCAGAG
GATTACTGGCTGACTCTGGGGGCGGTTCTTCCA--TGATGGTGTTCCTCTAAATTGCA- CGGAGAAACACCTGATTTCAGGAAA-ATCCCTCAGATGGCGCTGGTCCATCTCCGATGCCT
AGACCAGGCAAGACAACGTGAGC-GCGATGGCG--TGTACCCCAGGTCAAGGGGTGGTGTG-TCTATGAAGGAGGGGCCGAAG---CCCTTGTGGCGGGCCTCCCTGAGCCCCGTCTGTGGTGCCAG
CACTTCAGAAGGCT-TCTGAATGGAACCCTCTT--GACA-TTGTGTTCTATA-ATATTTG-T-CATGACAGTCACAGCATAAA-G---CGCAGACGGCTGTGACCTGATTAGA-AAATTTTTAGA
```

same-colored boxes *should* be aligned

Pfold (KH03) Test Set D



Knudsen & Hein, Pfold: RNA secondary structure prediction using stochastic context-free grammars, Nucleic Acids Research, 2003, v 31, 3423–3428

Approaches

Align-first: align sequences, then look for common structure

Fold-first: Predict structures, then try to align them

single-seq struct prediction only ~ 60% accurate;
exacerbated by flanking seq; no biologically-validated model for structural alignment

Joint: Do both together

Sankoff – good but slow

Heuristic

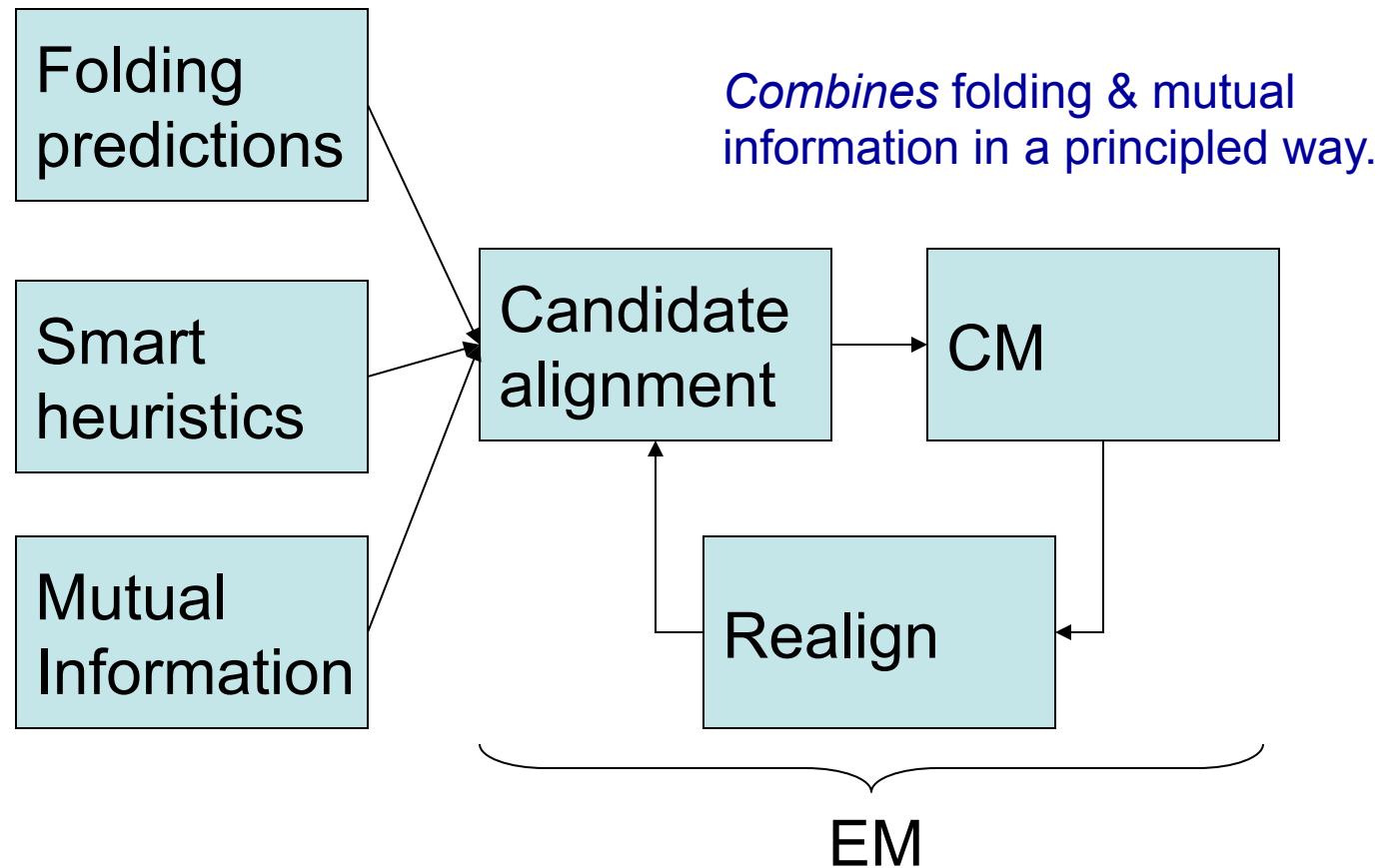
Our Approach: CMfinder

Simultaneous *local* alignment, folding and CM-based motif description using an EM-style learning procedure

Yao, Weinberg & Ruzzo, *Bioinformatics*, 2006

CMFinder

Simultaneous alignment, folding & motif description
Yao, Weinberg & Ruzzo, *Bioinformatics*, 2006



Structure Inference

Part of M-step is to pick a structure that maximizes data likelihood

We combine:

- mutual information

- position-specific priors for paired/unpaired
(based on single sequence thermodynamic folding predictions)

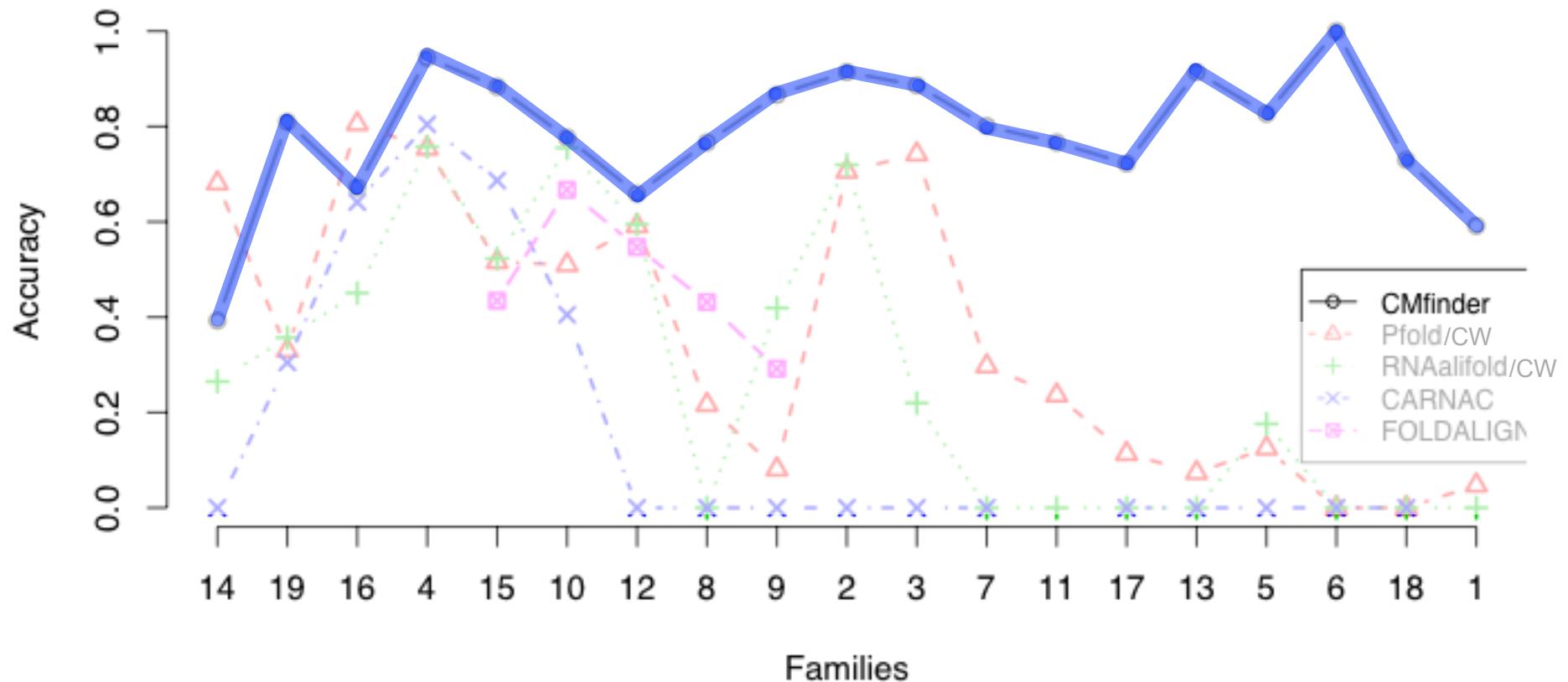
- intuition: for similar seqs, little MI; fall back on single-sequence folding predictions

- data-dependent, so not strictly Bayesian

- Details: see paper

CMfinder Accuracy

(on Rfam families with flanking sequence)



Summary of Rfam test families and results

ID	Family	Rfam ID	#seqs	%id	length	#hp	CMfinder	CW/Pfold	CW/RNAalifold	Carnac	Foldalign	ComRNA
1	Cobalamin	RF00174	71	49	216	4	0.59	0.05	0	X	-	0
2	ctRNA_pGA1	RF00236	17	74	83	2	0.91	0.70	0.72	0	0.86	0
3	Enter_CRE	RF00048	56	81	61	1	0.89	0.74	0.22	0	-	0
4	Enter_OriR	RF00041	35	77	73	2	0.94	0.75	0.76	0.80	0.52	0.52
5	glmS	RF00234	14	58	188	4	0.83	0.12	0.18	0	-	0.13
6	Histone3	RF00032	63	77	26	1	1	0	0	0	-	0
7	Intron_gpII	RF00029	75	55	92	2	0.80	0.30	0	0	-	0
8	IRE	RF00037	30	68	30	1	0.77	0.22	0	0	0.38	0
9	let-7	RF00027	9	69	84	1	0.87	0.08	0.42	0	0.71	0.78
10	lin-4	RF00052	9	69	72	1	0.78	0.51	0.75	0.41	0.65	0.24
11	Lysine	RF00168	48	48	183	4	0.77	0.24	0	X	-	0
12	mir-10	RF00104	11	66	75	1	0.66	0.59	0.60	0	0.48	0.33
13	Purine	RF00167	29	55	103	2	0.91	0.07	0	0	-	0.27
14	RFN	RF00050	47	66	139	4	0.39	0.68	0.26	0	-	0
15	Rhino_CRE	RF00220	12	71	86	1	0.88	0.52	0.52	0.69	0.41	0.61
16	s2m	RF00164	23	80	43	1	0.67	0.80	0.45	0.64	0.63	0.29
17	S_box	RF00162	64	66	112	3	0.72	0.11	0	0	-	0
18	SECIS	RF00031	43	43	68	1	0.73	0	0	0	-	0
19	Tymo_tRNA-like	RF00233	22	72	86	4	0.81	0.33	0.36	0.30	0.80	0.48
Average Accuracy:							0.79	0.36	0.28	0.17	0.60	0.19
Average Specificity:							0.81	0.42	0.57	0.83	0.60	0.65
Average Sensitivity:							0.77	0.36	0.23	0.13	0.61	0.17

Min/Max in col

Bold = best in row

Discovery in Bacteria

OPEN  ACCESS Freely available online

PLOS COMPUTATIONAL BIOLOGY

A Computational Pipeline for High-Throughput Discovery of *cis*-Regulatory Noncoding RNA in Prokaryotes

Zizhen Yao^{1,*}, Jeffrey Barrick^{2✉}, Zasha Weinberg³, Shane Neph^{1,4}, Ronald Breaker^{2,3,5}, Martin Tompa^{1,4}, Walter L. Ruzzo^{1,4}

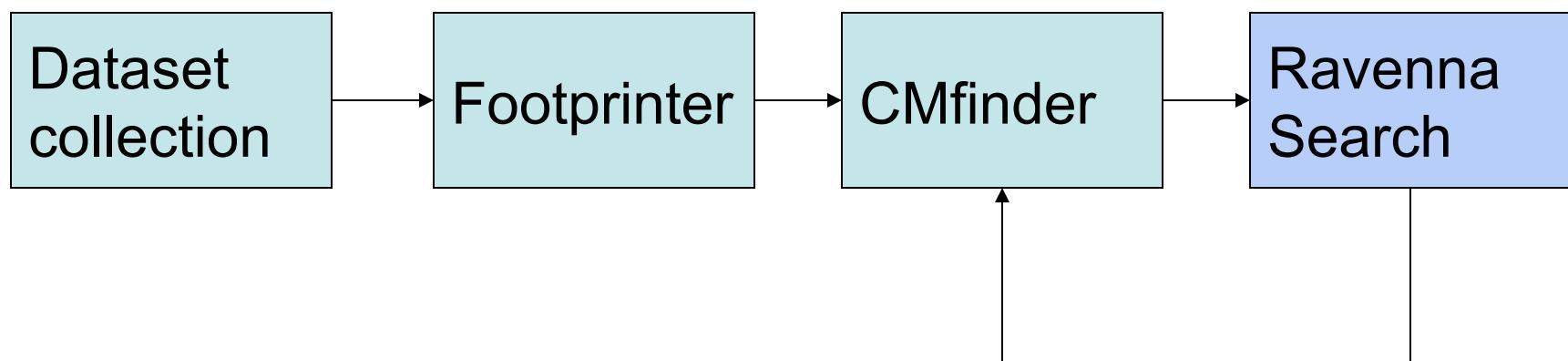
Published online 9 July 2007

Nucleic Acids Research, 2007, Vol. 35, No. 14 4809–4819
doi:10.1093/nar/gkm487

Identification of 22 candidate structured RNAs in bacteria using the CMfinder comparative genomics pipeline

Zasha Weinberg^{1,*}, Jeffrey E. Barrick^{2,3}, Zizhen Yao⁴, Adam Roth², Jane N. Kim¹, Jeremy Gore¹, Joy Xin Wang^{1,2}, Elaine R. Lee¹, Kirsten F. Block¹, Narasimhan Sudarsan¹, Shane Neph⁵, Martin Tompa^{4,5}, Walter L. Ruzzo^{4,5} and Ronald R. Breaker^{1,2,3}

Use the Right Data; Do Genome Scale Search



Right Data: Why/How

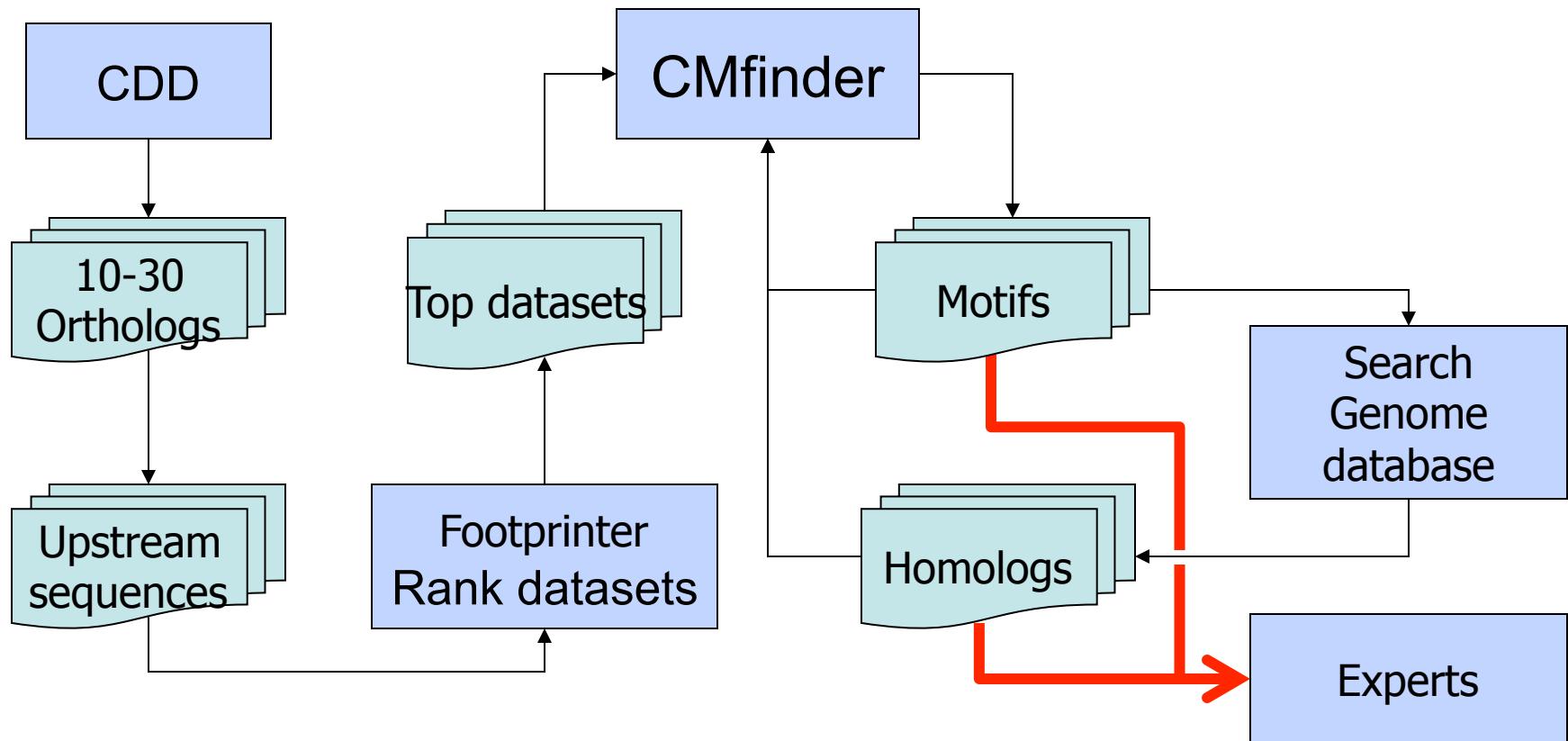
We can recognize, say, 5-10 good examples amidst 20 extraneous ones (but not 5 in 200 or 2000) of length 1k or 10k (but not 100k)

Regulators often near regulatees (protein coding genes), which are usually recognizable cross-species
So, look near similar genes (“homologs”)

Many riboswitches, e.g., are present in ~5 copies per genome

(Not strategy used in vertebrates - 1000x larger genomes)

A pipeline for RNA motif genome scans



Processing Times

Input from ~70 complete Firmicute genomes available in late 2005-early 2006, totaling ~200 megabases

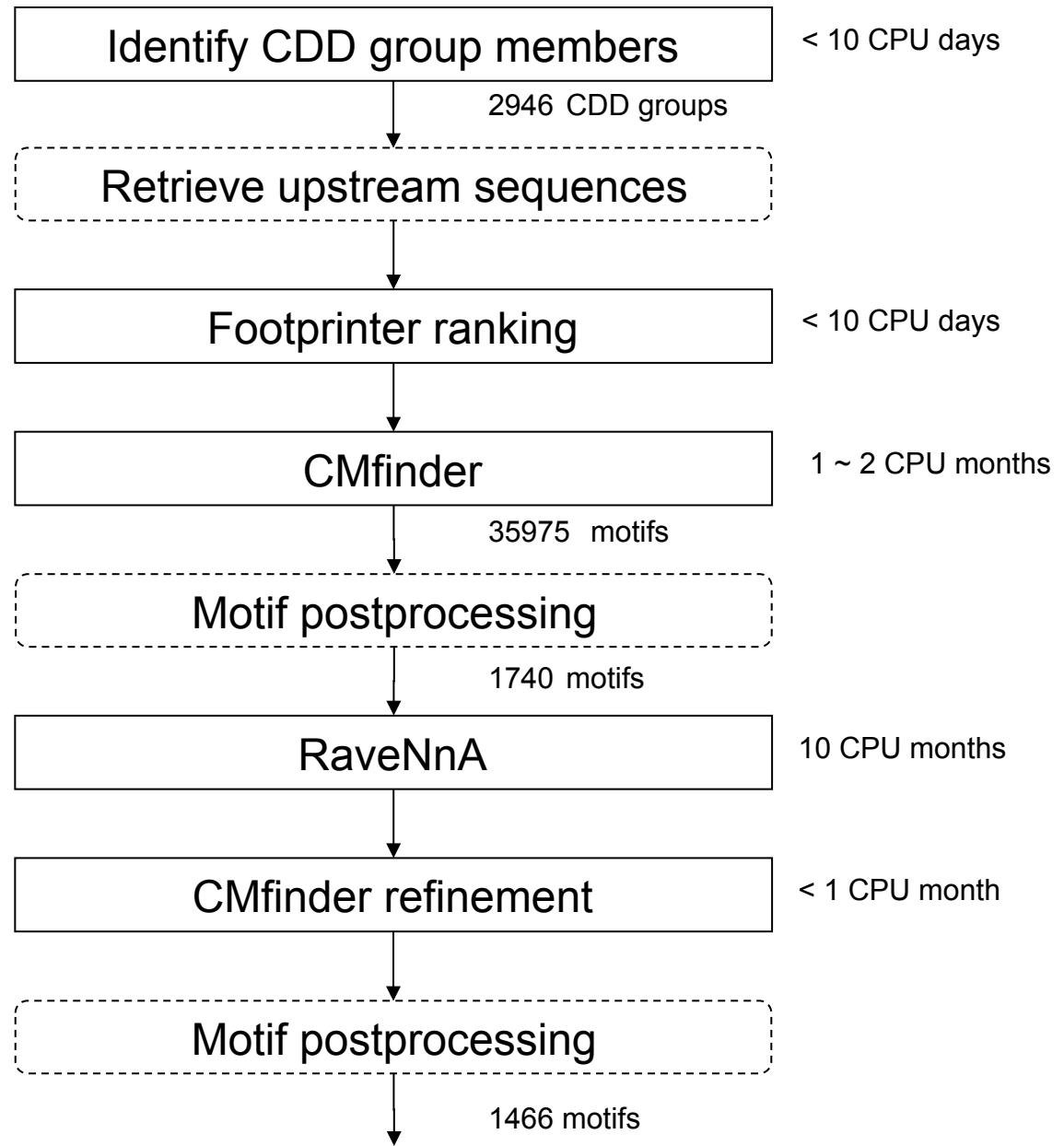


Table I: Motifs that correspond to Rfam families

Rank	Score			#	ID	Gene	CDD Description	Rfam
	RAV	CMF	FP	RAV CMF				
0	43	107	3400	367	11	9904 llvB	Thiamine pyrophosphate-requiring enzymes	RF00230 T-box
1	10	344	3115	96	22	13174 COG3859	Predicted membrane protein	RF00059 THI
2	77	1284	2376	112	6	11125 Meth	Methionine synthase I specific DNA methylase	RF00162 S_box
3	0	5	2327	30	26	9991 COG0116	Predicted N6-adenine-specific DNA methylase	RF00011 RNaseP_bact_b
4	6	66	2228	49	18	4383 DHBP	3,4-dihydroxy-2-butanone 4-phosphate synthase	RF00050 RFN
7	145	952	1429	51	7	10390 GuaA	GMP synthase	RF00167 Purine
8	17	108	1322	29	13	10732 GcvP	Glycine cleavage system protein P	RF00504 Glycine
9	37	749	1235	28	7	24631 DUF149	Uncharacterised BCR, YbaB family COG0718	RF00169 SRP_bact
10	123	1358	1222	36	6	10986 CbiB	Cobalamin biosynthesis protein CobD/CbiB	RF00174 Cobalamin
20	137	1133	899	32	7	9895 LysA	Diaminopimelate decarboxylase	RF00168 Lysine
21	36	141	896	22	10	10727 TerC	Membrane protein TerC	RF00080 yybP-ykoY
39	202	684	664	25	5	11945 MgtE	Mg/Co/Ni transporter MgtE	RF00380 ykoK
40	26	74	645	19	18	10323 glmS	Glucosamine 6-phosphate synthetase	RF00234 glmS
53	208	192	561	21	5	10892 OpuBB	ABC-type proline/glycine betaine transport systems	RF00005 tRNA ¹
122	99	239	413	10	7	11784 EmrE	Membrane transporters of cations and cationic drug	RF00442 ykkC-yxkD
255	392	281	268	8	6	10272 COG0398	Uncharacterized conserved protein	RF00023 tmRNA

Table 1: Motifs that correspond to Rfam families. “Rank”: the three columns show ranks for refined motif clusters after genome scans (“RAV”), CMfinder motifs before genome scans (“CMF”), and FootPrinter results (“FP”). We used the same ranking scheme for RAV and CMF. “Score”

Rfam		Membership			Overlap			Structure		
		#	Sn	Sp	nt	Sn	Sp	bp	Sn	Sp
RF00174	Cobalamin	183	0.74 ¹	0.97	152	0.75	0.85	20	0.60	0.77
RF00504	Glycine	92	0.56 ¹	0.96	94	0.94	0.68	17	0.84	0.82
RF00234	glmS	34	0.92	1.00	100	0.54	1.00	27	0.96	0.97
RF00168	Lysine	80	0.82	0.98	111	0.61	0.68	26	0.76	0.87
RF00167	Purine	86	0.86	0.93	83	0.83	0.55	17	0.90	0.95
RF00050	RFN	133	0.98	0.99	139	0.96	1.00	12	0.66	0.65
RF00011	RNaseP_bact_b	144	0.99	0.99	194	0.53	1.00	38	0.72	0.78
RF00162	S_box	208	0.95	0.97	110	1.00	0.69	23	0.91	0.78
RF00169	SRP_bact	177	0.92	0.95	99	1.00	0.65	25	0.89	0.81
RF00230	T-box	453	0.96	0.61	187	0.77	1.00	5	0.32	0.38
RF00059	THI	326	0.89	1.00	99	0.91	0.69	13	0.56	0.74
RF00442	ykkC-yxkD	19	0.90	0.53	99	0.94	0.81	18	0.94	0.68
RF00380	ykoK	49	0.92	1.00	125	0.75	1.00	27	0.80	0.95
RF00080	yybP-ykoY	41	0.32	0.89	100	0.78	0.90	18	0.63	0.66
mean		145	0.84	0.91	121	0.81	0.82	21	0.75	0.77
median		113	0.91	0.97	105	0.81	0.83	19	0.78	0.78

Tbl 2: Prediction accuracy compared to prokaryotic subset of Rfam full alignments.

Membership: # of seqs in overlap between our predictions and Rfam's, the sensitivity (Sn) and specificity (Sp) of our membership predictions. Overlap: the avg len of overlap between our predictions and Rfam's (nt), the fractional lengths of the overlapped region in Rfam's predictions (Sn) and in ours (Sp). Structure: the avg # of correctly predicted canonical base pairs (in overlapped regions) in the secondary structure (bp), and sensitivity and specificity of our predictions. ¹After 2nd RaveNnA scan, membership Sn of Glycine, Cobalamin increased to 76% and 98% resp., Glycine Sp unchanged, but Cobalamin Sp dropped to 84%.

Table 3: High ranking motifs not found in Rfam

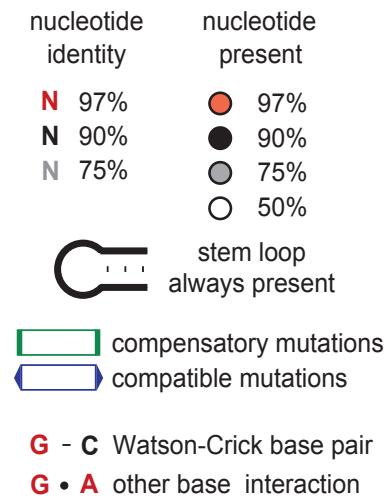
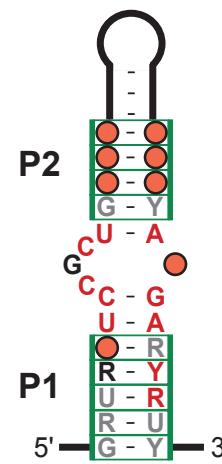
Rank	#	CDD	Gene: Description	Annotation
6	69	28178	DHOase IIa: Dihydroorotate	PyrR attenuator [22]
15	33	10097	RpL: Ribosomal protein L7/L1	L10 r-protein leader; see Supp
19	36	10234	RpsF: Ribosomal protein S6	S6 r-protein leader
22	32	10897	COG1179: Dinucleotide-utilizing enzymes	6S RNA [25]
27	27	9926	RpsJ: Ribosomal protein S10	S10 r-protein leader; see Supp
29	11	15150	Resolvase: N terminal domain	
31	31	10164	InfC: Translation initiation factor 3	IF-3 r-protein leader; see Supp
41	26	10393	RpsD: Ribosomal protein S4 and related proteins	S4 r-protein leader; see Supp [30]
44	30	10332	GroL: Chaperonin GroEL	HrcA DNA binding site [46]
46	33	25629	Ribosomal L21p: Ribosomal prokaryotic L21 protein	L21 r-protein leader; see Supp
50	11	5638	Cad: Cadmium resistance transporter	[47]
51	19	9965	RpIB: Ribosomal protein L2	S10 r-protein leader
55	7	26270	RNA pol Rpb2 1: RNA polymerase beta subunit	
69	9	13148	COG3830: ACT domain-containing protein	
72	28	4174	Ribosomal S2: Ribosomal protein S2	S2 r-protein leader
74	9	9924	RpsG: Ribosomal protein S7	S12 r-protein leader
86	6	12328	COG2984: ABC-type uncharacterized transport system	
88	19	24072	CtsR: Firmicutes transcriptional repressor of class III	CtsR DNA binding site [48]
100	21	23019	Formyl trans N: Formyl transferase	
103	8	9916	PurE: Phosphoribosylcarboxyaminoimidazole	
117	5	13411	COG4129: Predicted membrane protein	
120	10	10075	RpI(O): Ribosomal protein L15	L15 r-protein leader
121	9	10132	RpmJ: Ribosomal protein L36	IF-1 r-protein leader
129	4	23962	Cna B: Cna protein B-type domain	
130	9	25424	Ribosomal S12: Ribosomal protein S12	S12 r-protein leader
131	9	16769	Ribosomal L4: Ribosomal protein L4/L1 family	L3 r-protein leader
136	7	10610	COG0742: N6-adenine-specific methylase	yhbH putative RNA motif [4]
140	12	8892	Pencillinase R: Penicillinase repressor	Blai, Mecl DNA binding site [49]
157	25	24415	Ribosomal S9: Ribosomal protein S9/S16	L13 r-protein leader; Fig 3
160	27	1790	Ribosomal L19: Ribosomal protein L19	L19 r-protein leader; Fig 2
164	6	9932	GapA: Glyceraldehyde-3-phosphate dehydrogenase/erythrose	
174	8	13849	COG4708: Predicted membrane protein	
176	7	10199	COG0325: Predicted enzyme with a TIM-barrel fold	
182	9	10207	RpmF: Ribosomal protein L32	L32 r-protein leader
187	11	27850	LDH: L-lactate dehydrogenases	
190	11	10094	CspR: Predicted rRNA methylase	
194	9	10353	FusA: Translation elongation factors	EF-G r-protein leader

Example: Ribosomal Autoregulation:
Excess L19 represses L19 (RF00556; 555-559 similar)

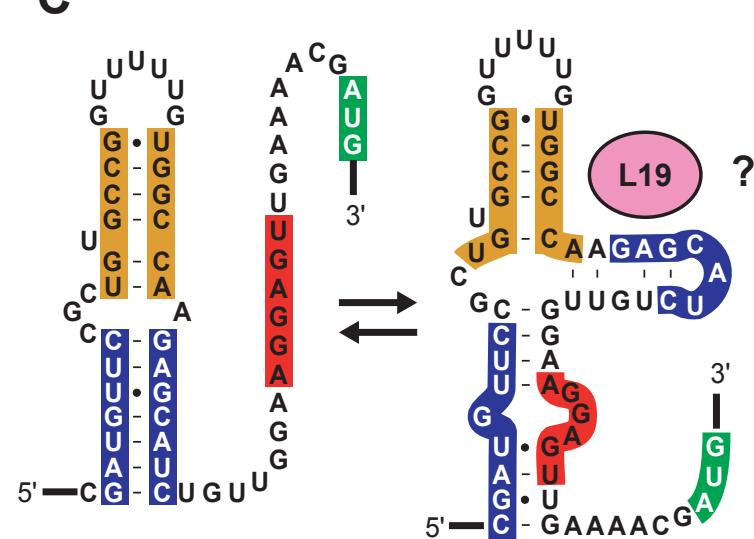
A L19 (*rplS*) mRNA leader

	-35	-10	TSS →	P1	P2	RBS	Start			
<i>Bsu</i>	TTGCAT.	17.	TAAGAT.	40. AAAAC	GAUGUUC	CGCUGUGCCG.. GUUUUUG...	GGC. CAAAGAGCAUC	UG. 05. AGGAGU. 08. AUG		
<i>Bha</i>	TTGTTTC.	17.	TCTTCT.	17. AUUAC	GAUGUUC	CGCUG. CAG... GGGUAGAAG...	CUGUCAUGAGCAUC	UG. 06. AGGAGG. 11. AUG		
<i>Oih</i>	TTGAAC.	17.	TATATT.	31. UAAAC	GAUGUUC	CGCUG. UC... CCAUACUU...	GUUCAUGAGCAUU	AG. 06. AGGAGU. 07. AUG		
<i>Bce</i>	TTGCTA.	18.	TATGCT.	36. UUAAC	GAUGUUC	CGCUG. UAA. UUUUAUAAGACU...	UUA. UAAGAGCAUC	UG. 05. AGGAGA. 09. AUG		
<i>Gka</i>	TTGCCT.	17.	TATCAT.	38. AAAAC	GAUGUUC	CGCUG. CAAUGA. AGAGA...	UCAUUGGCAUGAACAU	UG. 04. AGGAGU. 08. AUG		
<i>Bcl</i>	TTGTGC.	17.	TATGAT.	45. AUUAC	GAUAUUC	CGCUG. CUG... CAGUGU...	UGG. CAUGAAUGUC	UG. 06. AGGAGG. 10. AUG		
<i>Bac</i>	ATGACA.	17.	GATACT.	35. AUUAC	GAUGUUC	CGCUG. CA. AUAAAAGAAAGUCUG...	UG. CAAGAGCAUC	UG. 05. AGGAGU. 08. AUG		
<i>Lmo</i>	TTTACA.	17.	TAACCT.	28. AUUAC	GAUAUUC	CGCUG. CAU... UAUUAAU...	AUG. AAUGAAUGUU	UG. 05. AGGAGA. 07. AUG		
<i>Sau</i>	TTGAAA.	17.	TAACAT.	23. AUCAC	UAUCAUCC	CGCUG. CU... AUUAUAUUGUCG...	AGGCAAGAACAU	AG. 04. AGAGGA. 09. AUG		
<i>Cpe</i>	TTAAAG.	18.	TAACAT.	08. GUACC	GGCGGUCC	CUCUGUCACA...	UGUGUUAAGAACGUCA	AA. 17. AGGAGG. 08. AUG		
<i>Chy</i>	TTGCAT.	17.	TATAAT.	09. UACCAA	ACGUUC	CGCUG. GA... CAGGGC...	UC. CAUGAACGUGCC	03. AGGAGG. 09. AUG		
<i>Swo</i>	TTGAGA.	17.	TAAAAT.	16. AAAAA	GGUGGUCC	CGCUG. CAUU...	AAUG. UAUGAACACC	UU. 05. AGGAGG. 07. AUG		
<i>Ame</i>	TTGCGG.	17.	TATAAT.	10. UUACG	GGCGGUCC	CUCUA. UAC...	GU. UAAGAACGUCA	UA. 07. AGGAGG. 07. AUG		
<i>Dre</i>	TTGCC.	17.	TATAAT.	16. UUACG	GACGGUCC	CGCUG. CCU...	CUGGAA...	AGG. UAAGAACGUCA. 04. AGGAAG. 12. GUG		
<i>Spn</i>	TTTACT.	17.	TAAACT.	28. AUAC	GUUAUCC	CGCUG. AGGA...	AGAU...	UCCU. CAAGAUUGACAA. 04. AGGAGA. 05. AUG		
<i>Smu</i>	TTTACA.	17.	TACAAT.	26. AAACG	GCUAUAC	CGCUG. AG...	ACAGAGCA...	CU. UAUGAUUAAGUA. 04. AGGAGA. 07. AUG		
<i>Lpl</i>	TTGCGT.	18.	TATTCT.	21. UUAC	GAUGUUC	CGCUG. AC...	CAGGUU...	GU. CACGAAUGUC	GG. 04. AGGAAG. 09. AUG	
<i>Efa</i>	TTTACA.	17.	TAAACT.	28. AUUAC	AAUAUUC	CGCUG. UGG. CA...	GAAG...	UGACCA. UAAGAUAU	UG. 06. AGGAGA. 08. AUG	
<i>Ljo</i>	TTTACA.	17.	TAAACT.	25. UUAUG	GGUAUUC	CGCUG. GCAC...	AAG...	GUGUUGAU	GAAUGCC	GU. 03. AGGAGA. 07. AUG
<i>Sth</i>	TAGACA.	17.	TAAGAT.	29. UUACG	GGCUAAUC	CGCUG. AGA. CACAGAGGU...	UGCUCU...	UAAGAUUA	GUAA. 03. AGGAGU. 08. AUG	
<i>Lac</i>	TTAAAA.	17.	TTACTT.	39. UUAUG	GGGUAUUC	CGCUG. ACG...	CUGGU...	CGUUGAU	GAAUGCC	GA. 03. AGGAGA. 10. AUG
<i>Spy</i>	TTTACA.	17.	TAGAAAT.	29. UUACG	GGCUAAUC	CGCUG. AG...	ACAAGUA...	CU. UAAGAUUA	GUAA. 03. AGGAGA. 06. AUG	
<i>Lsa</i>	TTTTAA.	17.	TAAAAT.	26. ACAAC	GAUAUUC	CGCUG. GCG...	CAAGA...	CGUUAU	GAAUAUC	UG. 06. AGGAGA. 07. AUG
<i>Lsl</i>	TTTACT.	17.	TATTCT.	24. AUUAC	GAUAUUC	CGCUG. C...	AACUG...	GACAU	GAAUGUC	GG. 04. AGGAAA. 07. AUG
<i>Fnu</i>	TTGACA.	17.	TTAAAT.	12. AAUUC	GAUAUUC	CGCUG. UAA...	AAAA...	UUA. AAU	GAAUAUC	UU. 04. AGGAAG. 02. AUG

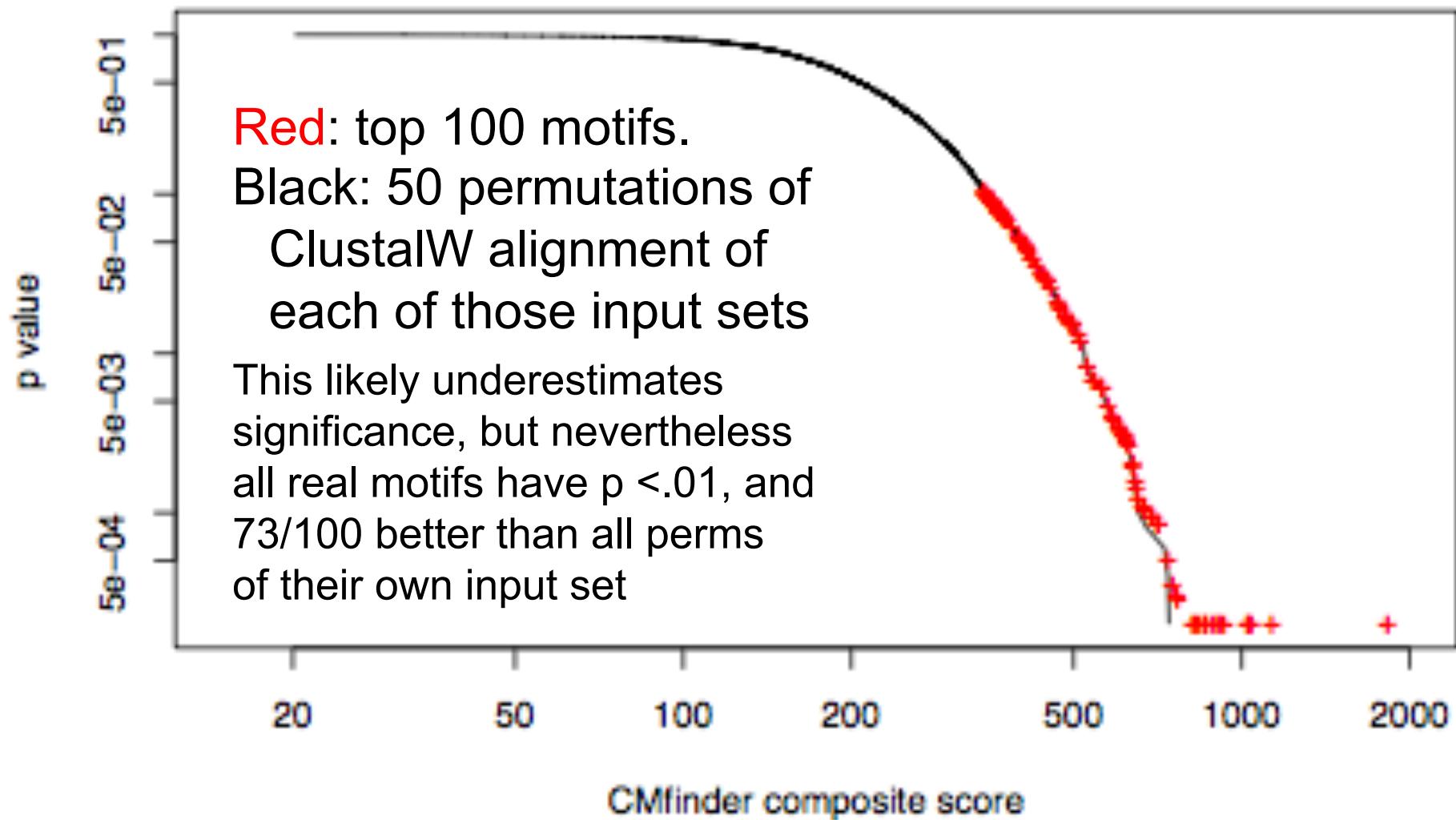
B



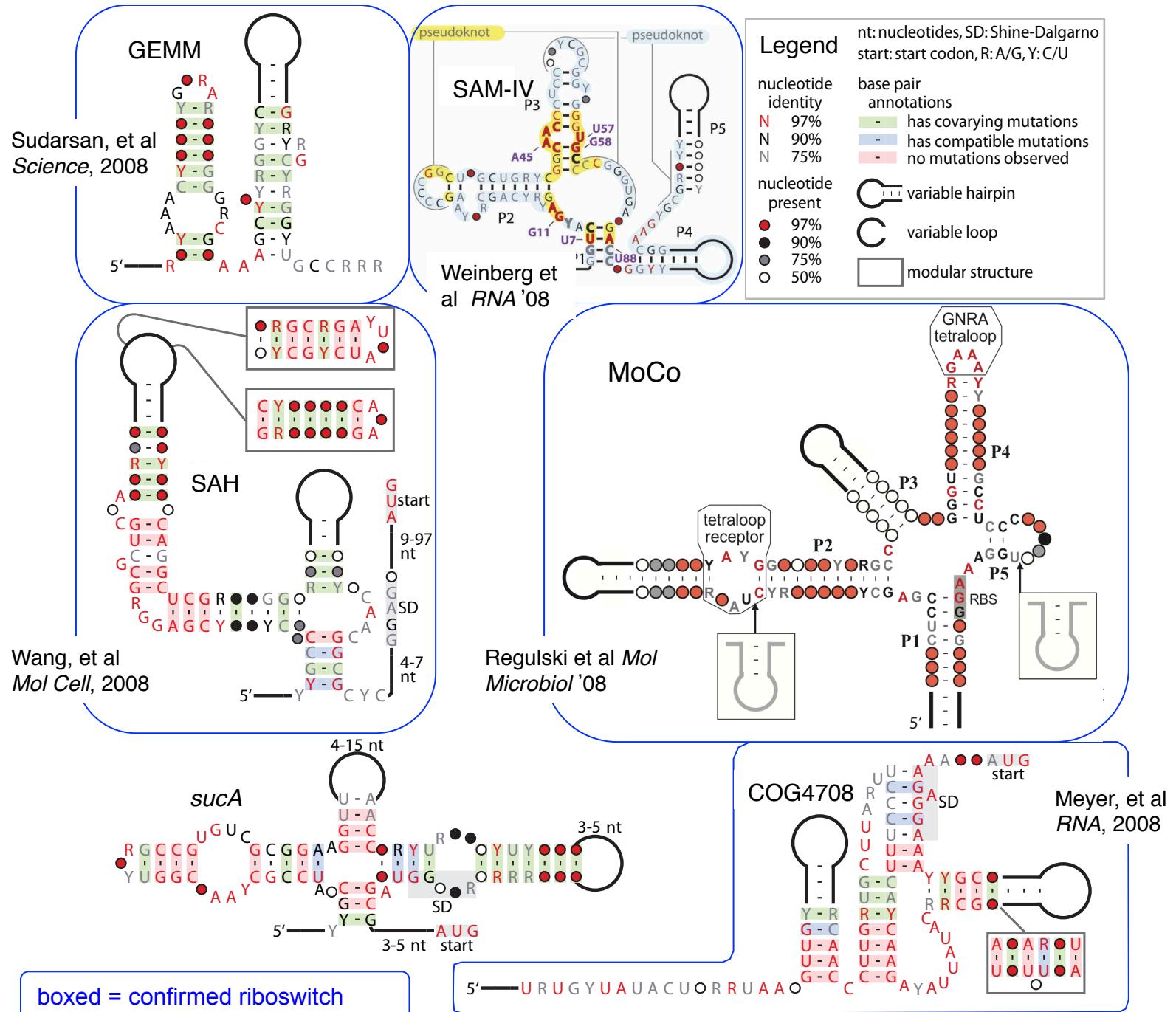
C



Estimating Motif Significance



Examples: 6 (of 22) Representative motifs



Vertebrate ncRNAs

Some Results

Human Predictions

Evofold

S Pedersen, G Bejerano, A Siepel, K Rosenbloom, K Lindblad-Toh, ES Lander, J Kent, W Miller, D Haussler, "Identification and classification of conserved RNA secondary structures in the human genome."

[PLoS Comput. Biol., 2, #4 \(2006\) e33.](#)

48,479 candidates (~70% FDR?)

FOLDALIGN

E Torarinsson, M Sawaya, JH Havgaard, M Fredholm, J Gorodkin, "Thousands of corresponding human and mouse genomic regions unalignable in primary sequence contain common RNA structure."

[Genome Res., 16, #7 \(2006\) 885-9.](#)

1800 candidates from 36970 (of 100,000) pairs

RNAz

S Washietl, IL Hofacker, M Lukasser, A Hutenhofer, PF Stadler, "Mapping of conserved RNA secondary structures predicts thousands of functional noncoding RNAs in the human genome."

[Nat. Biotechnol., 23, #11 \(2005\) 1383-90.](#)

30,000 structured RNA elements

~1000 conserved across *all* vertebrates.

~1/3 in introns of known genes, ~1/6 in UTRs

~1/2 located far from any known gene

CMfinder

Torarinsson, Yao, Wiklund, Bramsen, Hansen, Kjems, Tommerup, Ruzzo and Gorodkin. Comparative genomics beyond sequence based alignments: RNA structures in the ENCODE regions.

[Genome Research, Feb 2008, 18\(2\):242-251](#) PMID: 18096747

6500 candidates in ENCODE alone (better FDR, but still high)

Some details below

CMfinder Search in Vertebrates

Extract ENCODE* Multiz alignments

Remove exons, most conserved elements.

56017 blocks, 8.7M bps.

Apply CMfinder to both strands.

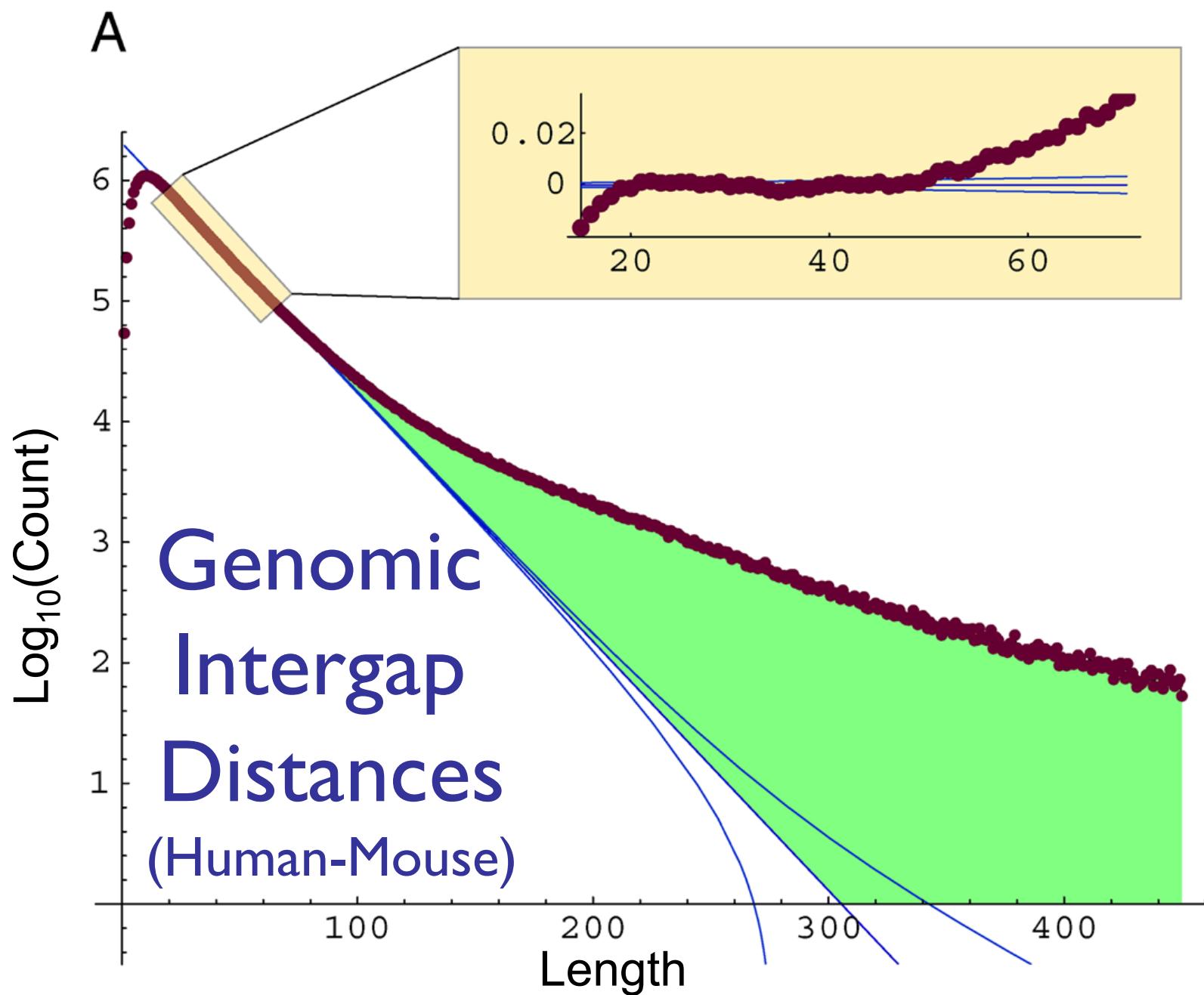
10,106 predictions, 6,587 clusters.

High false positive rate, but still suggests 1000's of RNAs.

Trust 17-way alignment for orthology, not for detailed alignment

(We've applied CMfinder to whole human genome:
many 100's of CPU years. Analysis in progress.)

* ENCODE: deeply annotated 1% of human genome



Genome-Wide Identification of Human Functional DNA Using a Neutral Indel Model
Gerton Lunter, Chris P. Ponting, Jotun Hein, PLoS Comput Biol 2006, 2(1): e5.

Overlap w/ Indel Purified Segments

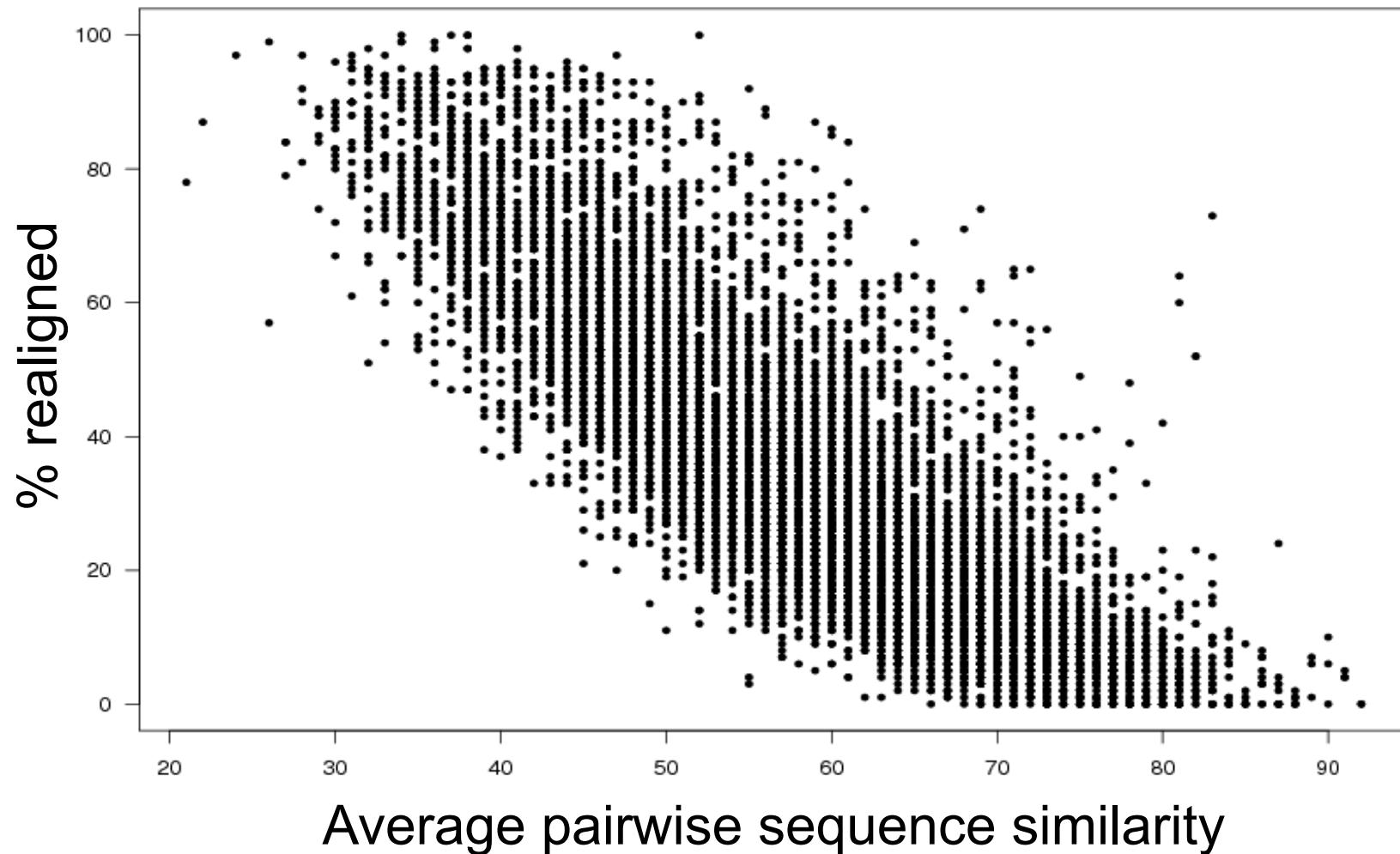
IPS presumed to signal purifying selection

Majority (64%) of candidates have >45% G+C

Strong P-value for their overlap w/ IPS

G+C	data	P	N	Expected	Observed	P-value	%
0-35	igs	0.062	380	23	24.5	0.430	5.8%
35-40	igs	0.082	742	61	70.5	0.103	11.3%
40-45	igs	0.082	1216	99	129.5	0.00079	18.5%
45-50	igs	0.079	1377	109	162.5	5.16E-08	20.9%
50-100	igs	0.070	2866	200	358.5	2.70E-31	43.5%
all	igs	0.075	6581	491	747.5	1.54E-33	100.0%

Realignment



Alignment Matters

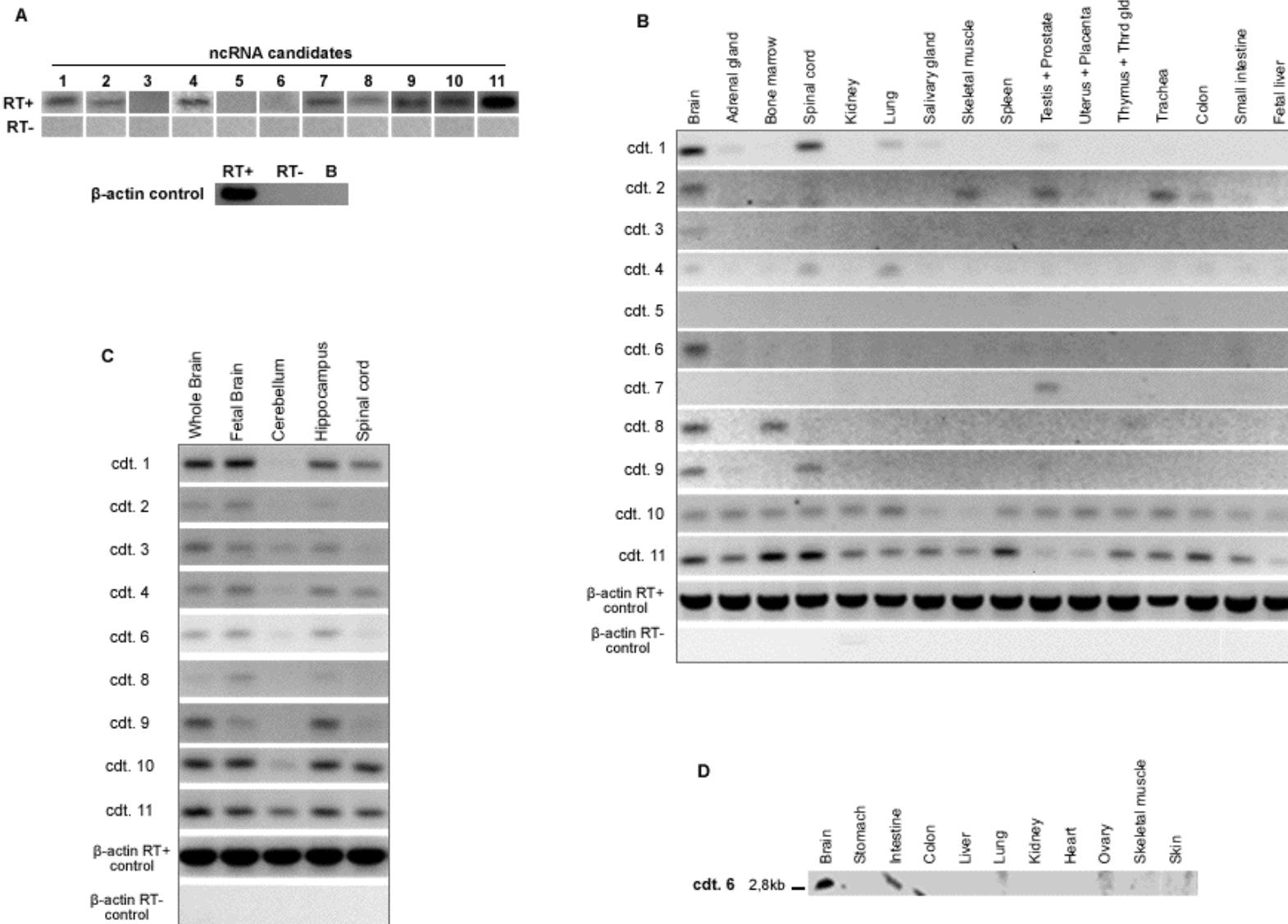
The original MULTIZ alignment without flanking regions. **RNAz Score: 0.132 (no RNA)**

Human	GGTCACTTCAAAGAGGGCTT-GTGGGGCTGTGAAACCA	AGAGGT-----CTT	AACAGTATGACCAAAA	ACTGAAGT	
Chimp	GGACATTTCAATGCAGGGCTC-ATGGGGCTGTGAAGCCA	AGAGCT-----ATT	AACACTATGACCAAGG	ACTGAAAT	
Cow	GGTCATTTCAAAGAGGGCTT-ATGAGACCA--AAACCG	GGAGCT-----CTT	AATGCTGTGACCAAA	AGATTGAAGT	
Dog	GGTCATTTCAAAGAGGGCTTGTGGAAC	TA--AAACCA	AGGGCT-----CTT	AACTCTGTGACCAA	ATATTAGAGT
Rabbit	GATCATTCAAAGAGGGTTT-GTGGTGCTGTGAAGTCA	AGAACT-----CTT	AACTGTATGCCA	AAAGATTAAAGT	
Rhesus	GGTCACTTCAAAGAGGGCTT-GTGGGGCTGTGAAACCA	AGAGGTAGGTCTT	AACAGTATAACC	AAAGACTGAAGT	
Str	(((((.....((((((....))))....))))....)))	

The local CMfinder re-alignment of the MULTIZ block. **RNAz Score: 0.709 (RNA)**

Human	GGTCACTTCAAAGAGGGCTT-GTGGGGCTGTGAAA-CCA	-----AGAGGTCTT	AACAGTATGACCAAAA	ACTGAAC	
Chimp	GGACATTTCAATGCAGGGCTC-ATGGGGCTGT-GAAGCCA	-----AGAGCTATT	AACACTATGACCAAGG	ACTGAAAT	
Cow	GGTCATTTCAAAGAGGGCTT-ATGAGACCA--AAA-CCG	-----GGAGCTCTT	AATGCTGTGACCAAA	AGATTGAAC	
Dog	GGTCATTTCAAAGAGGGCTTGTGGAAC	TA--AAA-CCA	-----AGGGCTCTT	AACTCTGTGACCAA	ATATTAGAC
Rabbit	GATCATTCAAAGAGGGTTT-GTGGTGCTGT-GAAGTCA	-----AGAACTCTT	AACTGTATGCCA	AAAGATTAAAC	
Rhesus	GGTCACTTCAAAGAGGGCTT-GTGGGGCTGTGAAA-CCA	AGAGG-TAGGTCTT	AACAGTATAACC	AAAGACTGAAC	
Str	(((((.....((((((....))))....))))....)))	

10 of 11 top (differentially) expressed



Open Problems - Better CM's

Optional- and variable-length stems

Riboswitches & other regulatory RNAs often switch between conformations; better search & alignment exploiting both alternatives?

“Augmented” CM handling pseudoknots probably too slow for scan, but plausibly could be used for alignment

Better use of prior knowledge? (GNRA tetraloops, single-stranded A's...)

Open Problems - Better algorithms & scoring

incorporating phylogeny in model construction & scoring

e.g. “mutual information” ignores it

improve scoring by “shuffling”

other ideas for scan filtering

comparing & clustering RNA structures

search/alignment/inference with splicing

Open Problems - Applications & Biology

clustering intergenic sequences, esp
prokaryotic

systematic look at eukaryotic UTRs

how to cluster? how to score?

“swiss-cheese phylogenies”

evidence for selection (no dN/dS)

ncRNA Summary

ncRNA is a “hot” topic

For family homology modeling: CMs

Training & search like HMM (but slower)

Dramatic acceleration possible

Automated model construction possible

New computational methods yield new discoveries

Many open problems