

RNA Search and Motif Discovery

Lecture 9
CSEP 590A
Summer 2006

Outline

Whirlwind tour of ncRNA search & discovery
Covariance Model Review
Algorithms for Training
“Mutual Information”
Algorithms for searching
Rigorous & heuristic filtering
Motif discovery
Wrap up
Course Evals

The Human Parts List, circa 2001

```
1 gagccggccc cgggggacgg gccgggggat agcgggaccc cggcgcggcg gtgcgcttca
61 gggcgcagcg gggcccgag accgagcccc gggcgggca agaggcggcg ggagccggtg
121 gggcgcagcg gggcccgag accgagcccc gggcgggca agaggcggcg ggagccggtg
181 gggcgcagcg gggcccgag accgagcccc gggcgggca agaggcggcg ggagccggtg
241 aacagagccc actcggcccc agagagagdc cgttggagga caccgacgcg ttaaaggacc
301 cggcgcagcg gggcccgag accgagcccc gggcgggca agaggcggcg ggagccggtg
361 cggcgcagcg gggcccgag accgagcccc gggcgggca agaggcggcg ggagccggtg
421 gggcgcagcg gggcccgag accgagcccc gggcgggca agaggcggcg ggagccggtg
481 acacacac ac aatattcgct gtagaatg aggtagctgc agtgacgato actgtctatg
541 cggcgcagcg gggcccgag accgagcccc gggcgggca agaggcggcg ggagccggtg
601 gggcgcagcg gggcccgag accgagcccc gggcgggca agaggcggcg ggagccggtg
661 gggcgcagcg gggcccgag accgagcccc gggcgggca agaggcggcg ggagccggtg
721 gggcgcagcg gggcccgag accgagcccc gggcgggca agaggcggcg ggagccggtg
781 ctggggccac cctgtgaga tgtgtcctgc ccagcctcac cctgcgcgc gtggcttoat
841 tccaaatata cgcacgggag cttgtcaaga tgtggatgaa tgccaggcca tcccgggctc
901 ctgtcagggg ggaattgca ttaatactgt tgggtctttt gaggccaast gccctgctgg
961 acacaaactt aatgaagtgt cacaaaaatg tgaagatatt gatgaatgca gaccattcc
1021 ...
```

3 billion nucleotides, containing:

- 25,000 protein-coding genes (only ~1% of the DNA)
- Messenger RNAs made from each
- Plus a double-handful of other RNA genes

Noncoding RNAs



Dramatic discoveries in last 5 years

100s of new families

Many roles: Regulation, transport, stability, catalysis, ...

1% of DNA codes for protein, but 30% of it is copied into RNA, i.e. ncRNA >> mRNA

“RNA sequence analysis using covariance models”

Eddy & Durbin
Nucleic Acids Research, 1994
vol 22 #11, 2079-2088
(see also, Ch 10 of Durbin *et al.*)

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

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”

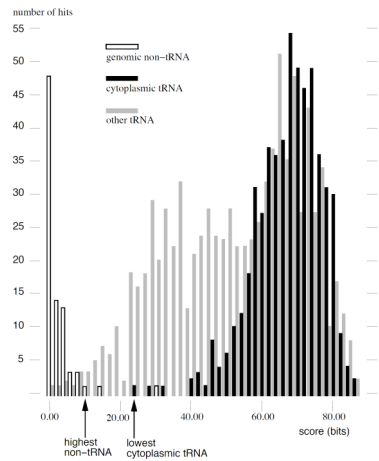
Scoring:

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

Viterbi approximation - find single best path

(Bonus: alignment & structure prediction)

Example: searching for tRNAs



Alignment Quality

Trusted:

```

DF6280 GCGGAUUUAGCUCAGU GGG AGAGC0CCGAGACUGAAG AUCUGGAG GUCCUGGUUCGUAUCCACAGAAUUCGACCA
DF6280 GCGGAUUUAGCUCAGU GGG AGAGC0CCGAGACUGAAGAAUAUCUCUGUCAAUUUUCUGGAG GUCCUGGUUCGUAUCCACAGAAUUCGCA
DF6280 UCUCGGAUAGUUDAAU GGUAGAAU0GGCGUUG GUCUUGGAG A UCUGGUAUUCUCCUCCUGGAGCCCA
DX1661 CUCUGGUGGAGCAGCCUGU AGCUCUUCUGGUCUCAGA ACCCGAAG GUCCUGGUUCAAAUCCGCCCCCGCACCA
DS6280 GGCACUUGGGCCGAG GGUUAAAGGC0AAAGAUA AGAUUUU GGCUUUUGCCCG CCGAGGUUCGAGUCCGAGUUGGGCCCA
    
```

U100:

```

DF6280 GCGGAUUUAGCUCAG UGGGAGAGCC0CCGAGU GA AG AUCUGGA GUCCUGGUUCGUAUCCACAGAAUUCGCA
DF6280 GCGGAUUUAGCUCAG UGGGAGAGCC0CCGAGUgaaagaaauu0CGUCAAUUUUCUGGA GUCCUGGUUCGUAUCCACAGAAUUCGCA
DF6280 UCUCGGAUAGUUDAAU GGUAGAAU0GGCGUUG GU CU CUGGCA GAU UCUGGUAUUCUCCUCCUGGAGCCCA
DX1661 CUCUGGUGGAGCAGCCUGUAGCUCUUCUGGUCU CA UA ACCCGAA GUCCUGGUUCAAAUCCGCCCCCGCACCA
DS6280 GGCACUUGGGCCGAG UGGUAAAGGC0AAAGAUA AG AA AUUUUUGGGUUUGCCG CCGAGGUUCGAGUCCGAGUUGGGCCCA
    
```

ClustalV:

```

DF6280 GCGGAUUUAGCUCAGUUGGGAGAGCC0CCGAGCUGAAGA UCUGGAGUUCGUAUCCACAGAAUUCGACCA
DF6280 GCGGAUUUAGCUCAGUUGGGAGAGCC0CCGAGCUGAAGAAUAUCUCUGUCAAUUUUCUGGAG GUCCUGGUUCGUAUCCACAGAAUUCGCA
DX1661 UCUCGGAUAGUUDAAU G GUAGAAUUGGGCGUUG CUUS UCUGGAG AGAUUGG GUUUCAAUUCUCCUCCUGGAGCCCA
DX1661 CUCUGGUGGAGCAGCC CUGGAGCUCUUCUGG CUGA UAACCGGA AUUUCUGGUAUUCUCCUCCUGGAGCCCA
DS6280 GGCACUUGGGCCGAGUGGUUAAAGGC0AAAGAUA AGAAAUUUUGGGC UUDGCCG CCGAGGUUCGAGUCCGAGUUGGGCCCA
    
```

Comparison to TRNASCAN

Fichant & Burks - best heuristic then

- 97.5% true positive
- 0.37 false positives per MB

CM A1415 (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

Profile HMM Structure

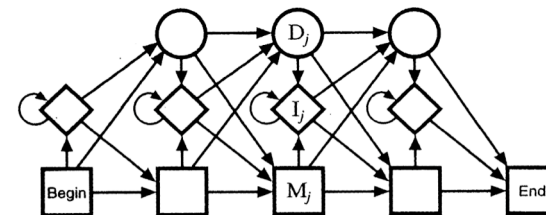


Figure 5.2 The transition structure of a profile HMM.

M_j: Match states (20 emission probabilities)
I_j: Insert states (Background emission probabilities)
D_j: Delete states (silent - no emission)

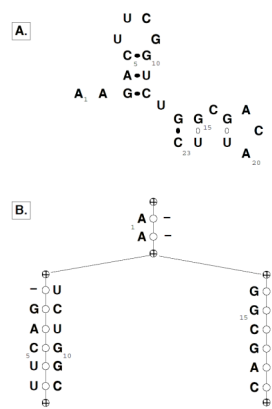
CM Structure

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)

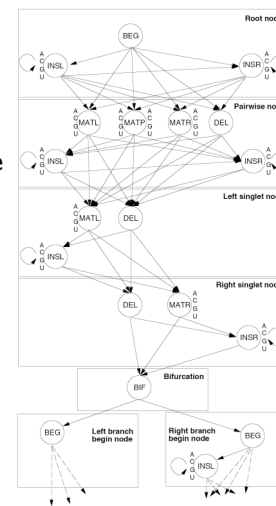


Overall CM Architecture

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

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

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

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

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

MI-Based Structure-Learning

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

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

“Just like Nussinov/Zucker folding”

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

Pseudoknots
disallowed allowed $(\sum_{i=1}^n \max_{j \neq i} M_{i,j})/2$

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

Table 1: Statistics of the training and test sets of 100 tRNA sequences each. The average identity in an alignment is the average pairwise identity of all aligned symbol pairs, with gap/symbol alignments counted as mismatches. Primary sequence information content is calculated according to [48]. Calculating pairwise mutual information content is an NP-complete problem of finding an optimum partition of columns into pairs. A lower bound is calculated by using the model construction procedure to find an optimal partition subject to a non-pseudoknotting restriction. An upper bound is calculated as sum of the single best pairwise covariation for each position, divided by two; this includes all pairwise tertiary interactions but overcounts because it does not guarantee a disjoint set of pairs. For the meaning of multiple alignment accuracy of ClustalV, see the text.

Model	training set	iterations	score (bits)	alignment accuracy
A1415	all sequences (aligned)	3	58.7	95%
A100	SIM100 (aligned)	3	57.3	94%
A65	SIM65 (aligned)	3	46.7	93%
U100	SIM100 (degapped)	23	56.7	90%
U65	SIM65 (degapped)	29	47.2	91%

Table 2: Training and multiple alignment results from models trained from the trusted alignments (A models) and models trained from no prior knowledge of tRNA (U models).

Rfam – an RNA family DB

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

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

Rel 7.0, 3/05: 503 families, >300k instances

Rfam

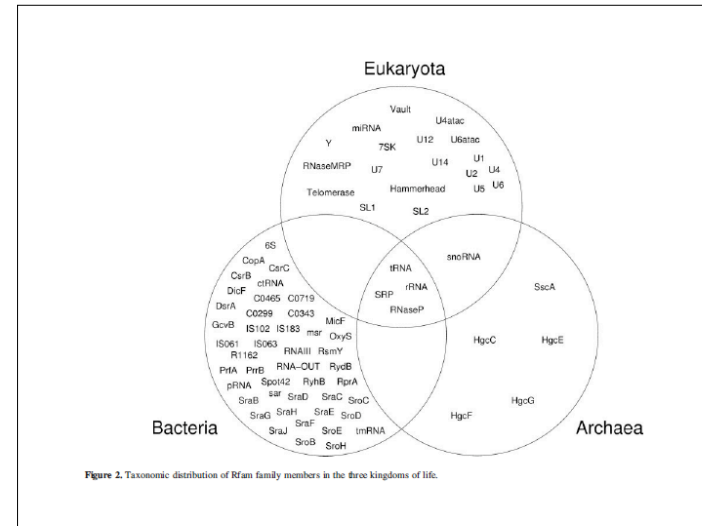
04400 G 400JJ JGU
 33303 300 JJ0RA GA
 33303 300 JJ0RA GA

Input (hand-curated):
 MSA "seed alignment"
 SS_cons
 Score Thresh T
 Window Len W

Output:
 CM
 scan results & "full alignment"

IRE (partial seed alignment):

Hom. sap.	GUUCCUGCUUCAACAGUGUUUGGAUGGAAC
Hom. sap.	UUUCUUC . UUCAACAGUGUUUGGAUGGAAC
Hom. sap.	UUUCCUGUUUCAACAGUGCUUGGA . GGAAC
Hom. sap.	UUUAUC . . AGUGACAGAGUUCACU . AUA
Hom. sap.	UCUCUUGCUUCAACAGUGUUUGGAUGGAAC
Hom. sap.	AUUAUC . . GGAACAGUGUUUCC . AUAU
Hom. sap.	UCUUC . . UUCAACAGUGUUUGGACGGAAG
Hom. sap.	UGUAUC . . GGAGACAGUGAUUCC . AUAUG
Hom. sap.	AUUAUC . . GGAAGCAGUGCCUCC . AUAU
Cav. por.	UCUCCUGCUUCAACAGUGCUUGGACGGAAC
Mus. mus.	UAUAUC . . GGAGACAGUGAUUCC . AUAUG
Mus. mus.	UUUCCUGCUUCAACAGUGCUUGGACGGAAC
Mus. mus.	GUACUUGCUUCAACAGUGUUUGGACGGAAC
Rat. nor.	UAUAUC . . GGAGACAGUGAUUCC . AUAUG
Rat. nor.	UAUCUUGCUUCAACAGUGUUUGGACGGAAC
SS_cons	<<<<<. . .<<<<. . .>>>>. >>>>. >>>>



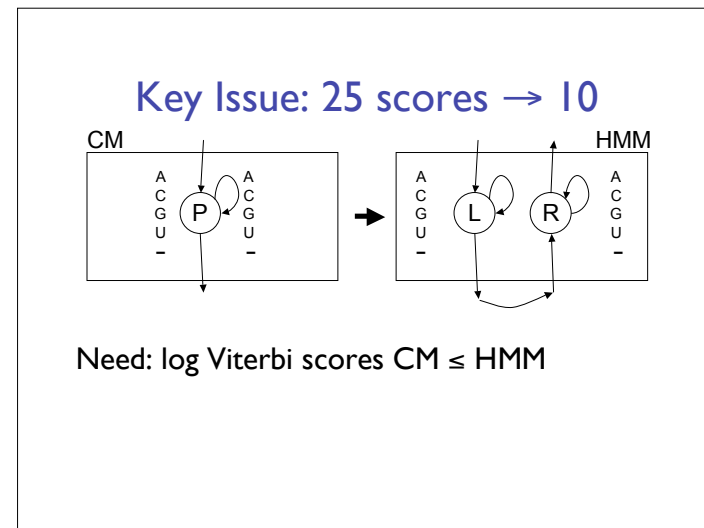
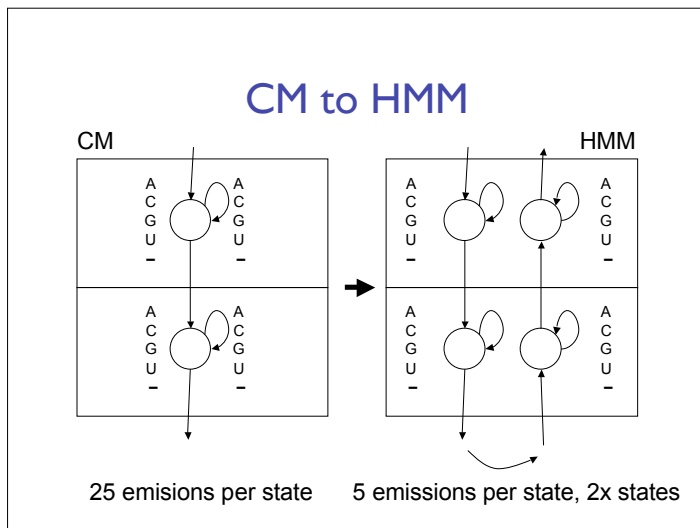
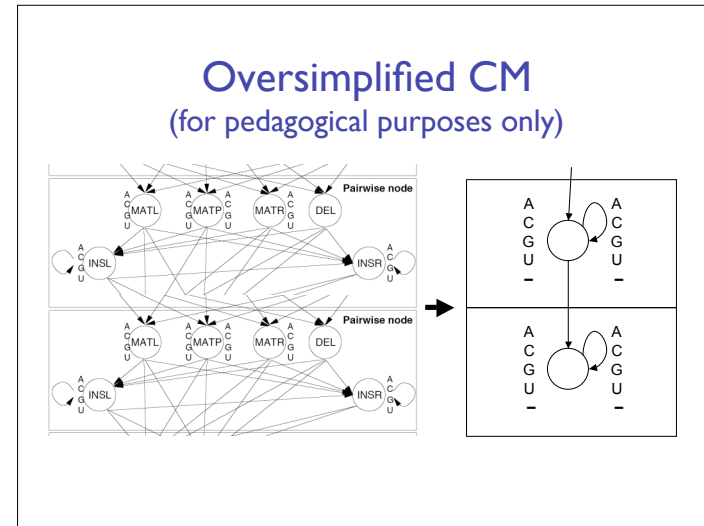
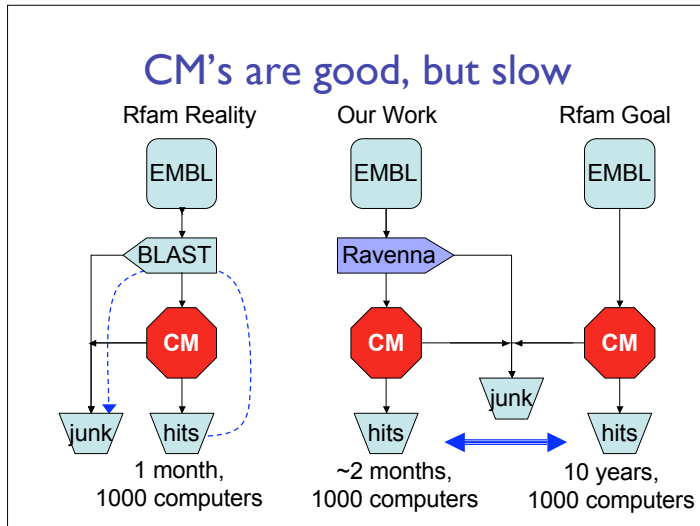
**Faster Genome Annotation
 of Non-coding RNAs
 Without Loss of Accuracy**

Zasha Weinberg
 & W.L. Ruzzo

Recomb '04, ISMB '04, Bioinfo '06

Covariance Model

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



Viterbi/Forward Scoring

Path π defines transitions/emissions

Score(π) = product of “probabilities” on π

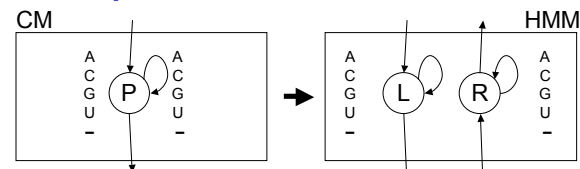
NB: ok if “probs” aren’t, e.g. $\sum \neq 1$
(e.g. in CM, emissions are odds ratios vs 0th-order background)

For any nucleotide sequence x :

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

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

Key Issue: 25 scores \rightarrow 10



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

$$\begin{array}{ll}
 P_{AA} \leq L_A + R_A & P_{CA} \leq L_C + R_A \quad \dots \\
 P_{AC} \leq L_A + R_C & P_{CC} \leq L_C + R_C \quad \dots \\
 P_{AG} \leq L_A + R_G & P_{CG} \leq L_C + R_G \quad \dots \\
 P_{AU} \leq L_A + R_U & P_{CU} \leq L_C + R_U \quad \dots \\
 P_{A-} \leq L_A + R_- & P_{C-} \leq L_C + R_- \quad \dots
 \end{array}$$

NB: HMM not a prob. model!

Rigorous Filtering

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

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

$$P_{UA} = 1 \leq L_U + R_A$$

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

Option 1:

$$L_U = R_A = R_G = 2$$

Option 2:

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

Assuming $\text{ACGU} \approx 25\%$

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

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

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

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

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

} ~100x speedup

Results: New ncRNA's?

Name	# found BLAST + CM	# found rigorous filter + CM	# new
<i>Pyrococcus</i> snoRNA	57	180	123
Iron response element	201	322	121
Histone 3' element	1004	1106	102
Purine riboswitch	69	123	54
Retron msr	11	59	48
Hammerhead I	167	193	26
Hammerhead III	251	264	13
U4 snRNA	283	290	7
S-box	128	131	3
U6 snRNA	1462	1464	2
U5 snRNA	199	200	1
U7 snRNA	312	313	1

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

Filter Chains

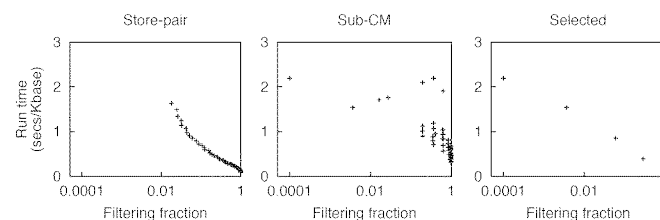


Fig. 2. Filter creation and selection. Filters for Rfam tRNA (RF00005) generated by the store-pair and sub-CM techniques and those selected for actual filtering are plotted by filtering fraction and run time. The CM runs at 3.5 secs/kbase. The four selected filters are run one after another, from highest to lowest fraction.

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

“ML heuristic”: train HMM from the infinite alignment generated by the CM

Often 10x faster, modest loss in sensitivity

Heuristic Filters

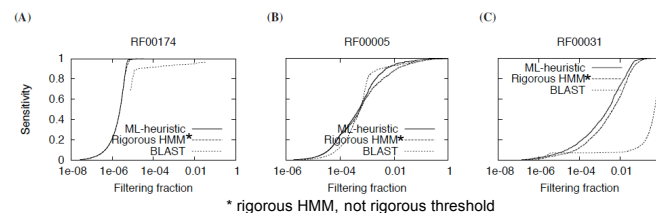
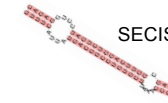
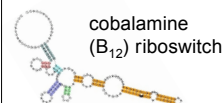


Fig. 1. Selected ROC-like curves. All plot sensitivity against filtering fraction, with filtering fraction in log scale. (A) RF00174 is typical of the other families; the ML-heuristic is slightly better than the rigorous profile HMM, and both often dramatically exceed BLAST. (B) Atypically, in RF00005, 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 ~90% of hits are known; see text.) The supplement includes all ROC-like curves, and the inferior ignore-SS.



Cmfinder--A Covariance Model Based RNA Motif Finding Algorithm

[Bioinformatics, 2006, 22\(4\): 445-452](#)

Zizhen Yao
Zasha Weinberg
Walter L. Ruzzo
University of Washington, Seattle

Searching for noncoding RNAs

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

An 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.

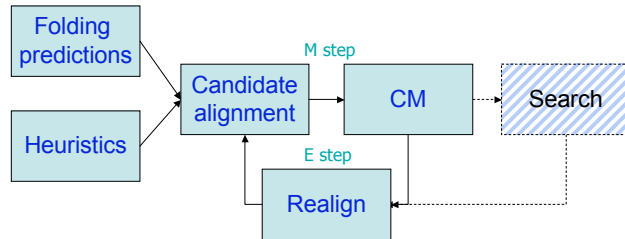
Align the motif instances.

Predict the consensus secondary structure.

Motif search space is huge!

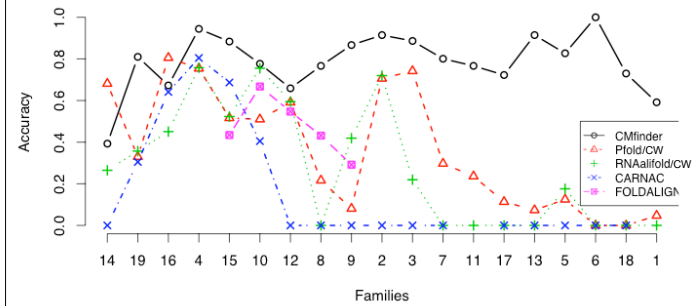
Motif location space, alignment space, structure space.

CMfinder Outline

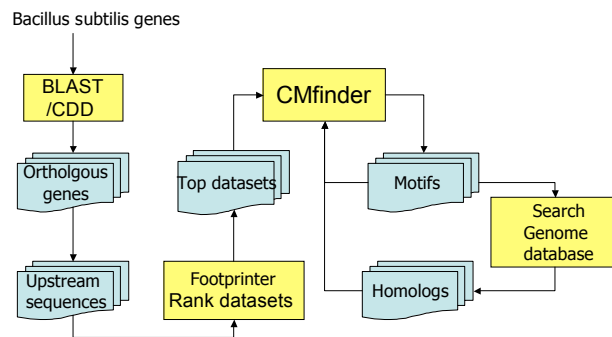


M-step uses M.I. + folding energy for structure prediction

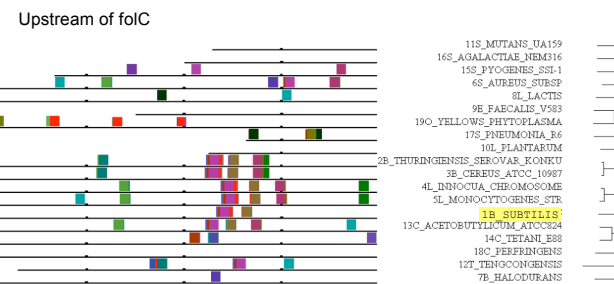
CMfinder Accuracy (on Rfam families with flanking sequence)



A pipeline for RNA motif genome scans



Footprinter finds patterns of conservation



A blind test

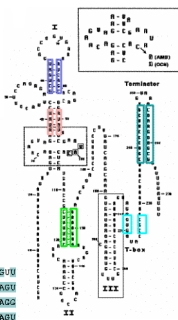
1ST genome scan: 234 sequences
 2ND genome scan: 447 sequences
The motif turned out to be T box
 Match to Rfam T box family: 299 OF 342
 False Positives: 89/148 are probable (upstream of annotated RNA-synthetase genes)

```

AUAUUC...GUUACGU...UCCAGAGACGUCAGUUGCCCGGUGAAA...AUCCGACGAGCGGAUAUAU
CGAAAU...GUCCUUCUUAUAUAAGAUUCGUAUGUUGUGGAA...AUCCGAAAG...AAACAUFUUU
AGAAAU...AGAAACCG...AUCUUGCGAUUGAGGAAU...GGUGUGAG...GCGAAAGGUUUU
CGAAAU...GUCCUUCUUAUAUAAGAUUCGUCAGUUGGUGUGGAAA...AUCCGAAAG...AAACAUFUUU

CGAA...UACACUICAUAGAACCCUUUUGGAAACRABCCGGGCGCGUUCAGUA...GUGGAGG
UGAA...UCCAUUCUGGAAU...GSRUAUUGGAAUUCUUAUUGGAAU...AGUAAGCAUUGC
ASAAAUUC...ACUCUUGAGU...UUCAUUAGGAAA...CA.....AGUAAGCAUUGG
UGAA...UCCAUUCUGGAAU...GSRUAUUGGAAUUCUUAUUGGAAU...AGUAAGCAUUGC

ACGGAC...CUGAUCUGUUAUCAGGCAAAAGUACCGCGCAUAUAUC...AUUCGUCUUCUUGUUAAGCGAAGGGGCGU
...CGGUS...AAGAGCGGUUAU...UUAUAGCGGCGG...GUUAGUCGCGGCGGUUAUUAUAGGAGCGGAGU
...CGGUUCAGC...UCCGUUAUCGUAUCUUAAGCGGCGCA...GUCUUCUUCUUCUUCUUAU...UGGUUAGAAAGC
...CGGUS...AAGAGCGGUUAU...UUAUAGCGGCGG...GUUAGUCGCGGCGGUUAUUAUAGGAGCGGAGU
    
```



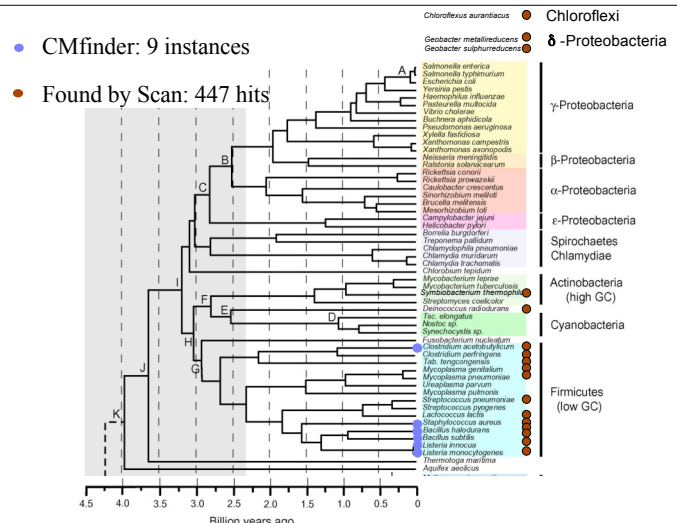
tyrs T box structure

Some Preliminary Actino Results

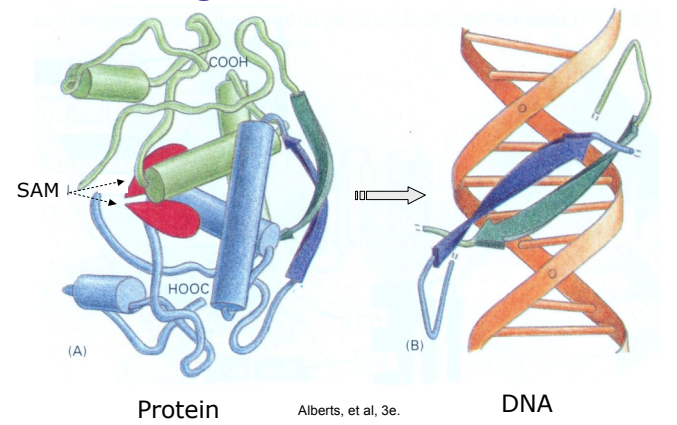
8 of 10 Rfam families found

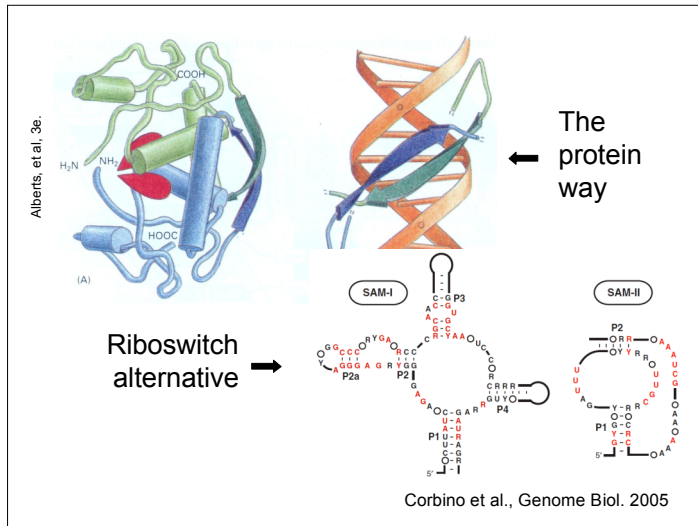
Rfam Family	Type (metabolite)	Rank
THI	riboswitch (thiamine)	4
ydaO-yuaA	riboswitch (unknown)	19
Cobalamin	riboswitch (cobalamin)	21
SRP_bact	gene	28
RFN	riboswitch (FMN)	39
yybP-ykoY	riboswitch (unknown)	48
gcvT	riboswitch (glycine)	53
S_box	riboswitch (SAM)	401
tmRNA	gene	Not found
RNaseP	gene	Not found

not cis-regulatory (got one anyway) (pointing to SRP_bact and S_box)



Gene Regulation: The MET Repressor





More Prelim Actino Results

Many others (not in Rfam) are likely real of top 50:

known (Rfam, 23S)	10
probable (Tbox, CIRCE, LexA, parP, pyrR)	7
probable (ribosomal genes)	9
potentially interesting	12
unknown or poor	12

One bench-verified, 2 more in progress

Preliminary results of genome scan

Top 115 datasets (some are redundant)
 13 T box, 22 riboswitches, 30 ribosomal genes
 RNase P, tRNA, CIRCE elements and other DNA binding sites

Gene	#motif	hits	RFAM	#seed	#full	#TP	specificity	sensitivity
metK	13	150	S_box	71	151	145	0.967	0.960
ribB	9	106	RFN	48	114	97	0.915	0.851
folC	9	447	T_box	67	342	299	0.669	0.874
xpt	14	106	Purine	37	100	97	0.915	0.970
glmS	16	33	glmS	14	37	33	1.000	0.892
thiA	16	305	THI	237	366	305	1.000	0.833
ykoY	10	34	yybP-ykoY	74	127	33	0.971	0.260

Summary

ncRNA - apparently widespread, much interest

Covariance Models - powerful but expensive tool for ncRNA motif representation, search, discovery

Rigorous/Heuristic filtering - typically 100x speedup in search with no/little loss in accuracy

CMfinder - CM-based motif discovery in unaligned sequences

Course Wrap Up

What is DNA? RNA?
How many Amino Acids are there?
Did human beings, as we know them, develop from earlier species of animals?
What are stem cells?
What did Viterbi invent?
What is dynamic programming?
What is a likelihood ratio test?
What is the EM algorithm?
How would you find the maximum of $f(x) = ax^3 + bx^2 + cx + d$ in the interval $-10 < x < 25$?

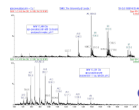
“High-Throughput BioTech”

Sensors

DNA sequencing
Microarrays/Gene expression
Mass Spectrometry/Proteomics
Protein/protein & DNA/protein interaction

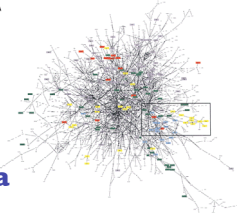
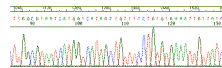
Controls

Cloning
Gene knock out/knock in
RNAi



Floods of data

“Grand Challenge” problems



CS Points of Contact

Scientific visualization

Gene expression patterns

Databases

Integration of disparate, overlapping data sources
Distributed genome annotation in face of shifting underlying coordinates

AI/NLP/Text Mining

Information extraction from journal texts with inconsistent nomenclature, indirect interactions, incomplete/inaccurate models,...

Machine learning

System level synthesis of cell behavior from low-level heterogeneous data (DNA sequence, gene expression, protein interaction, mass spec,

Algorithms

...

Frontiers & Opportunities

New data:

Proteomics, SNP, arrays CGH, comparative sequence information, methylation, chromatin structure, ncRNA, interactome

New methods:

graphical models? rigorous filtering?

Data integration

many, complex, noisy sources

Frontiers & Opportunities

Open Problems:

splicing, alternative splicing
multiple sequence alignment (genome scale, w/ RNA etc.)
protein & RNA structure
interaction modeling
network models
RNA trafficking
ncRNA discovery
...

Exciting Times

Lots to do
Various skills needed
I hope I've given you a taste of it

Thanks!