

Genome 559

Wi 2009

RNA

Function, Search, Discovery

The Message

Cells make lots of ~~RNA~~ *noncoding* RNA

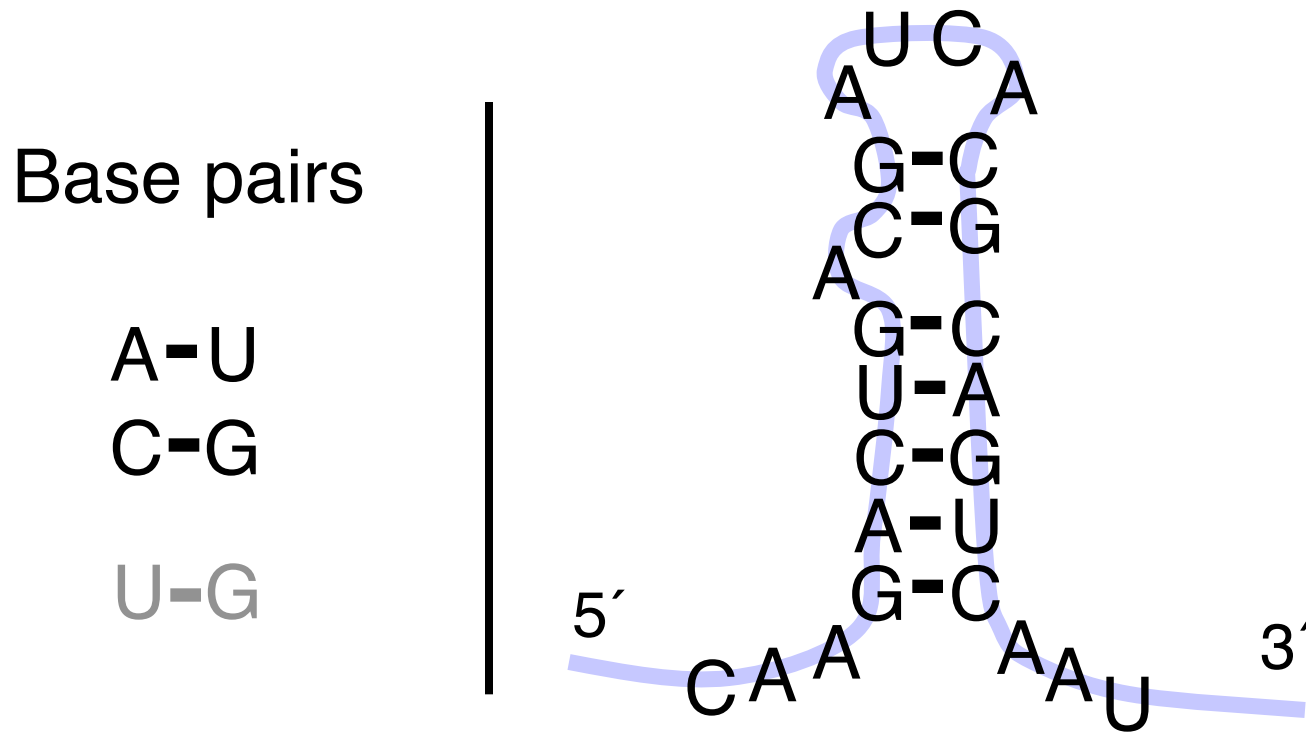
Functionally important, functionally diverse

Structurally complex

New tools required

alignment, discovery, search, scoring, etc.

RNA Secondary Structure: RNA makes helices too



Usually *single* stranded

Central Dogma of Molecular Biology

by

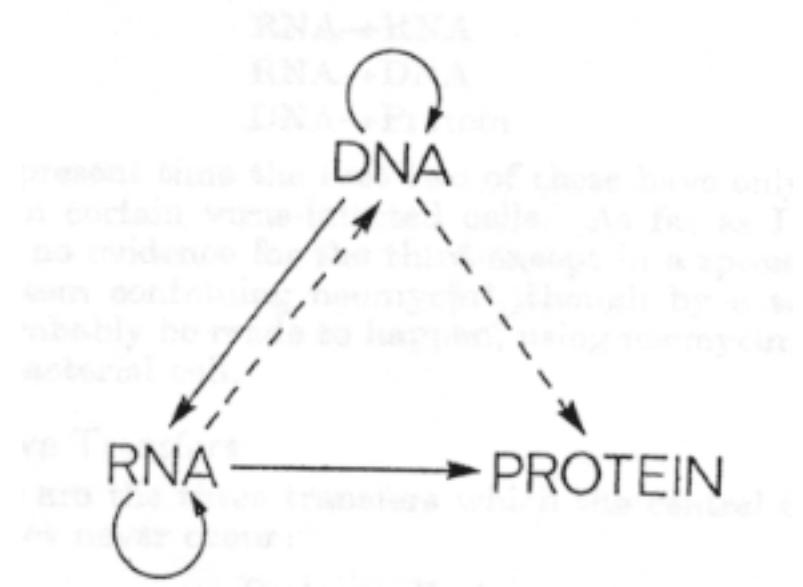
FRANCIS CRICK

MRC Laboratory
Hills Road,
Cambridge CB2 2QH

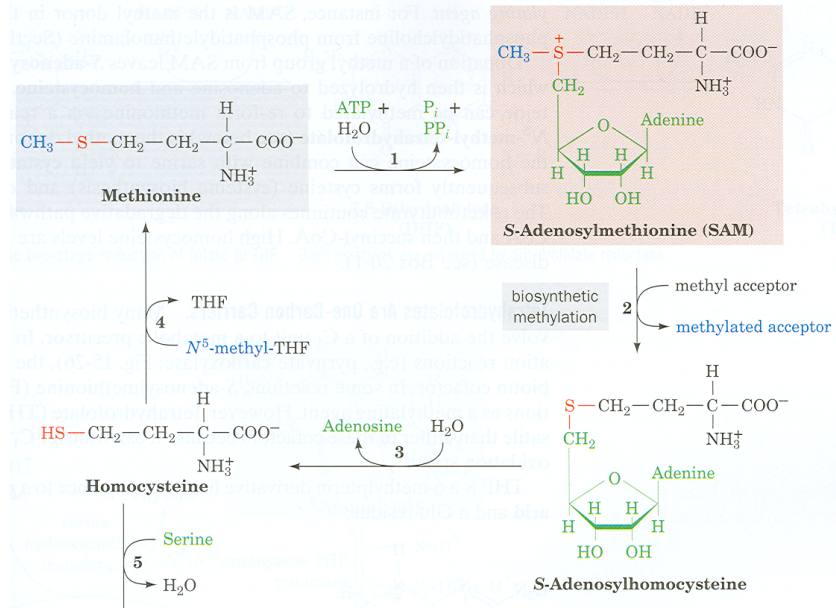
The central dogma of molecular biology deals with the detailed residue-by-residue transfer of sequential information. It states that such information cannot be transferred from protein to either protein or nucleic acid.

“The central dogma, enunciated by Crick in 1958 and the keystone of molecular biology ever since, is likely to prove a considerable over-simplification.”

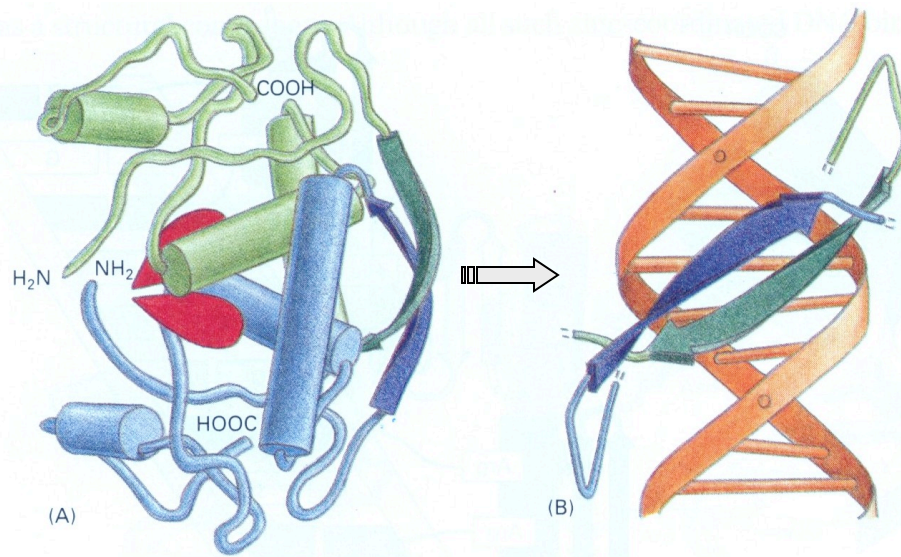
Fig. 2. The arrows show the situation as it seemed in 1958. Solid arrows represent probable transfers, dotted arrows possible transfers. The absent arrows (compare Fig. 1) represent the impossible transfers postulated by the central dogma. They are the three possible arrows starting from protein.



Proteins catalyze & regulate biochemistry

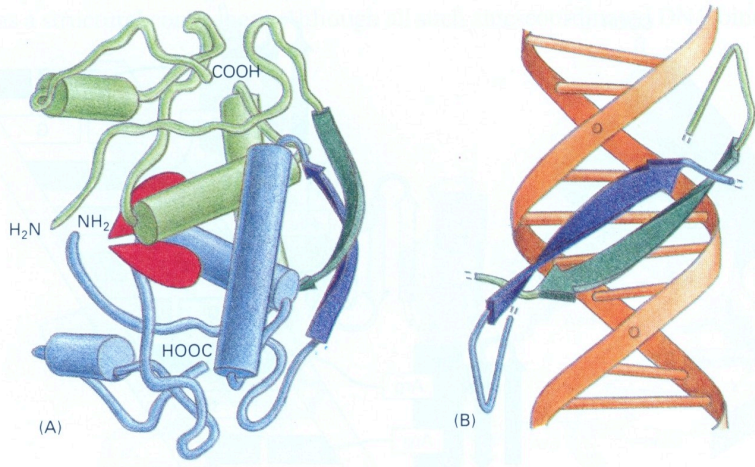


SAM



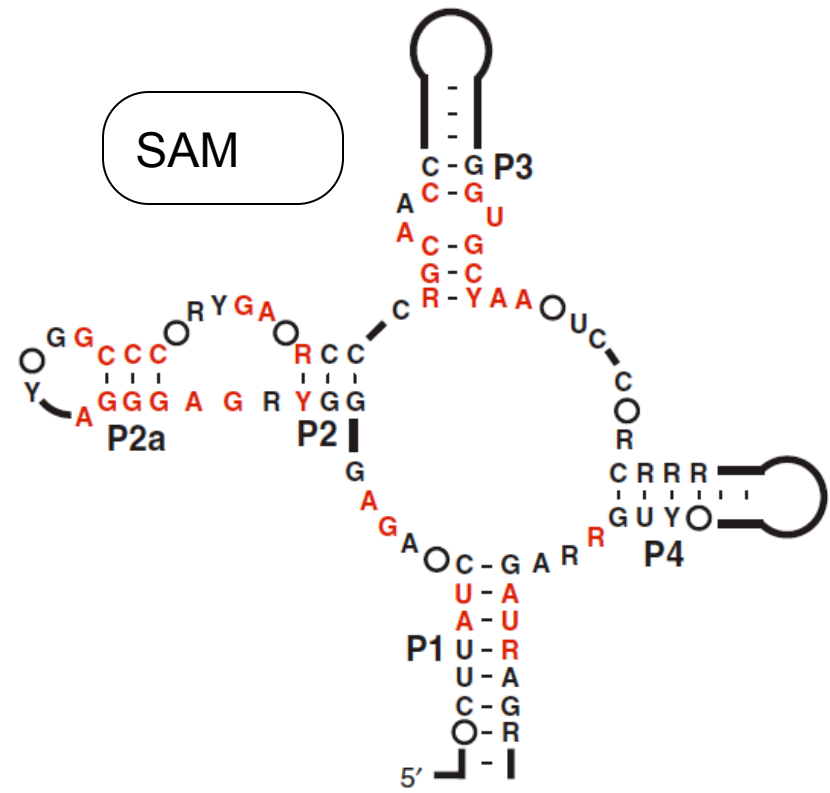
The Met Repressor

Alberts, et al, 3e.



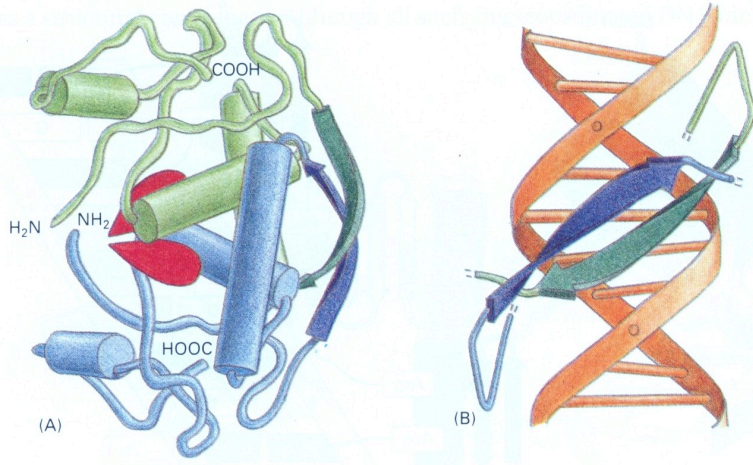
← The protein way

Riboswitch alternative



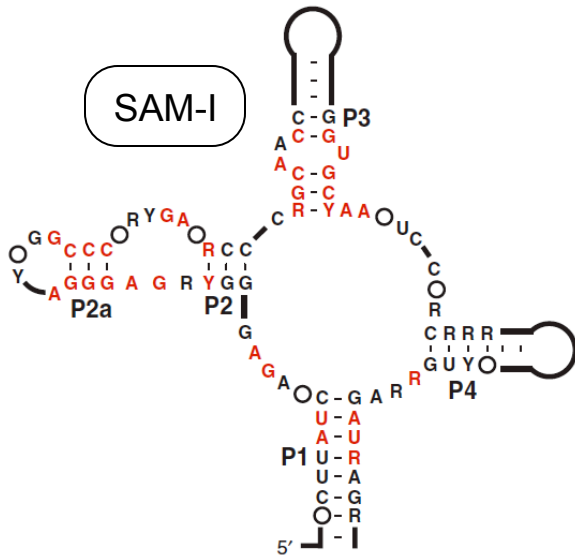
Grundy & Henkin, Mol. Microbiol 1998
Epshtein, et al., PNAS 2003
Winkler et al., Nat. Struct. Biol. 2003

Alberts, et al, 3e.

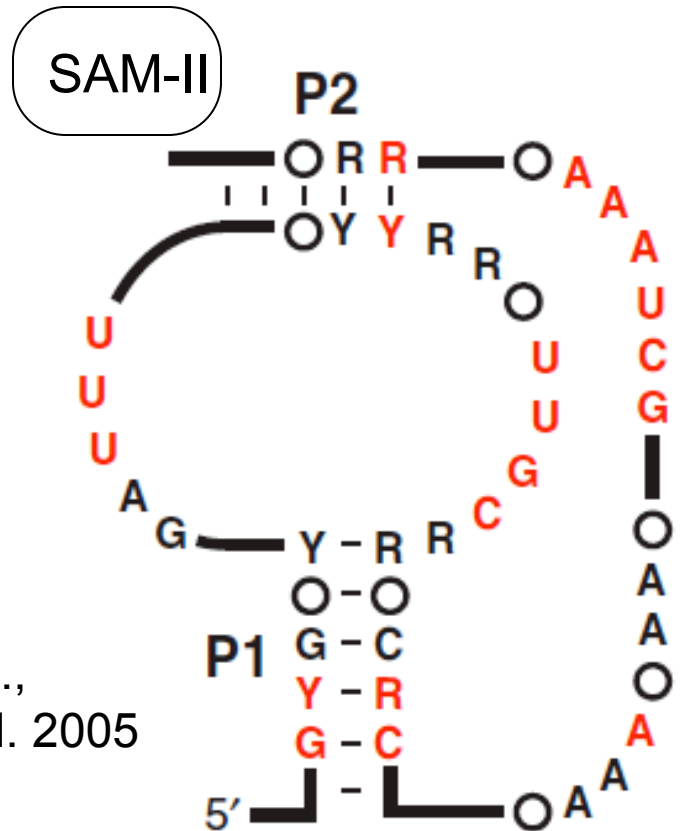


← The protein way

Riboswitch alternatives

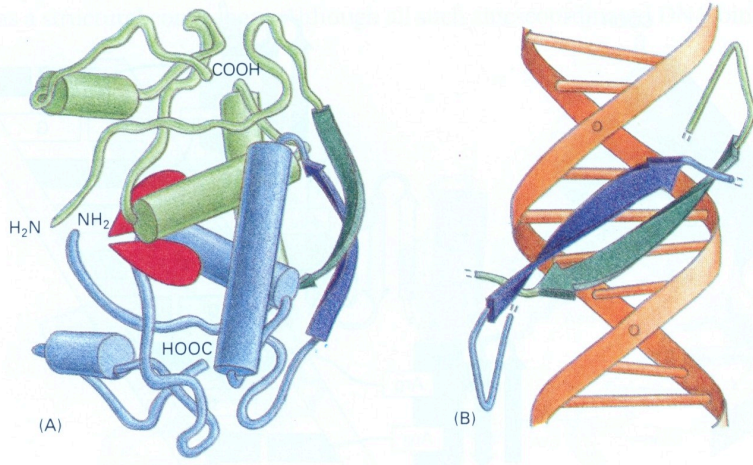


Grundy, Epshtein, Winkler et al., 1998, 2003



Corbino et al.,
Genome Biol. 2005

Alberts, et al, 3e.

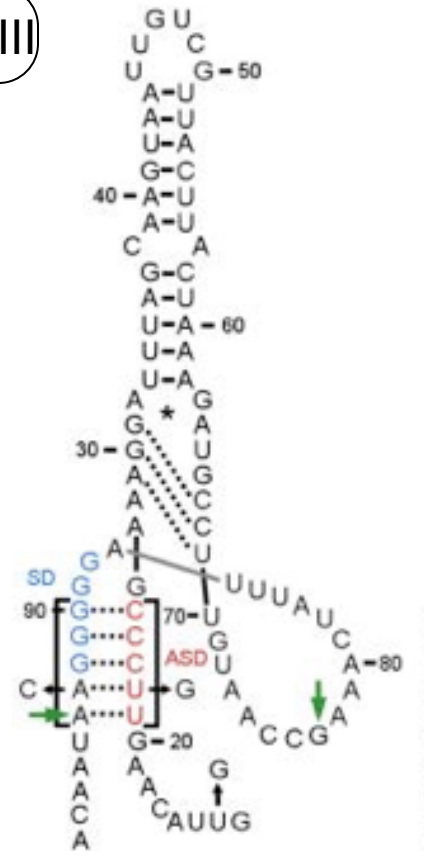


← The protein way

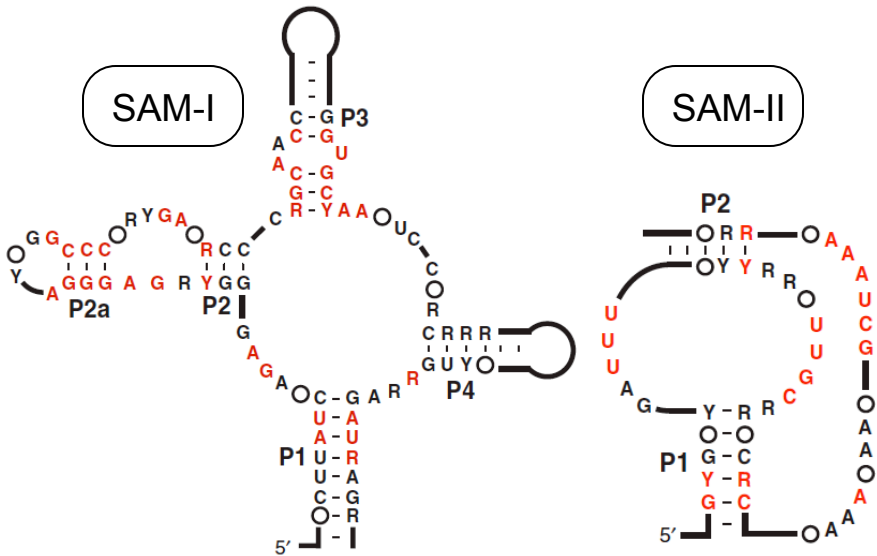
Riboswitch alternatives



SAM-III



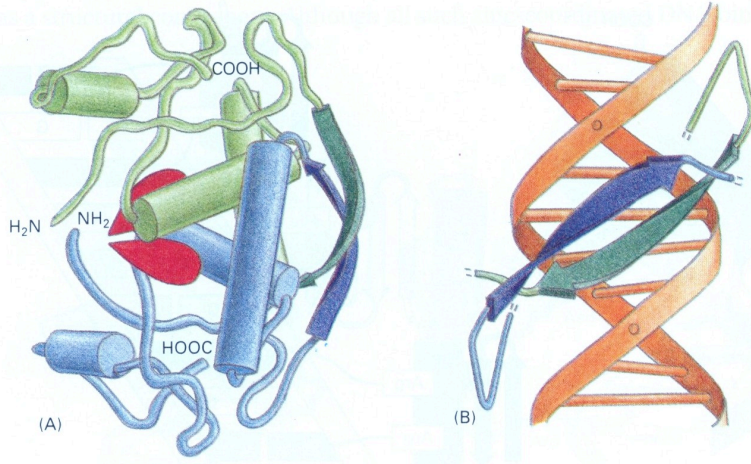
Fuchs et al.,
NSMB 2006



Grundy, Epshtein, Winkler
et al., 1998, 2003

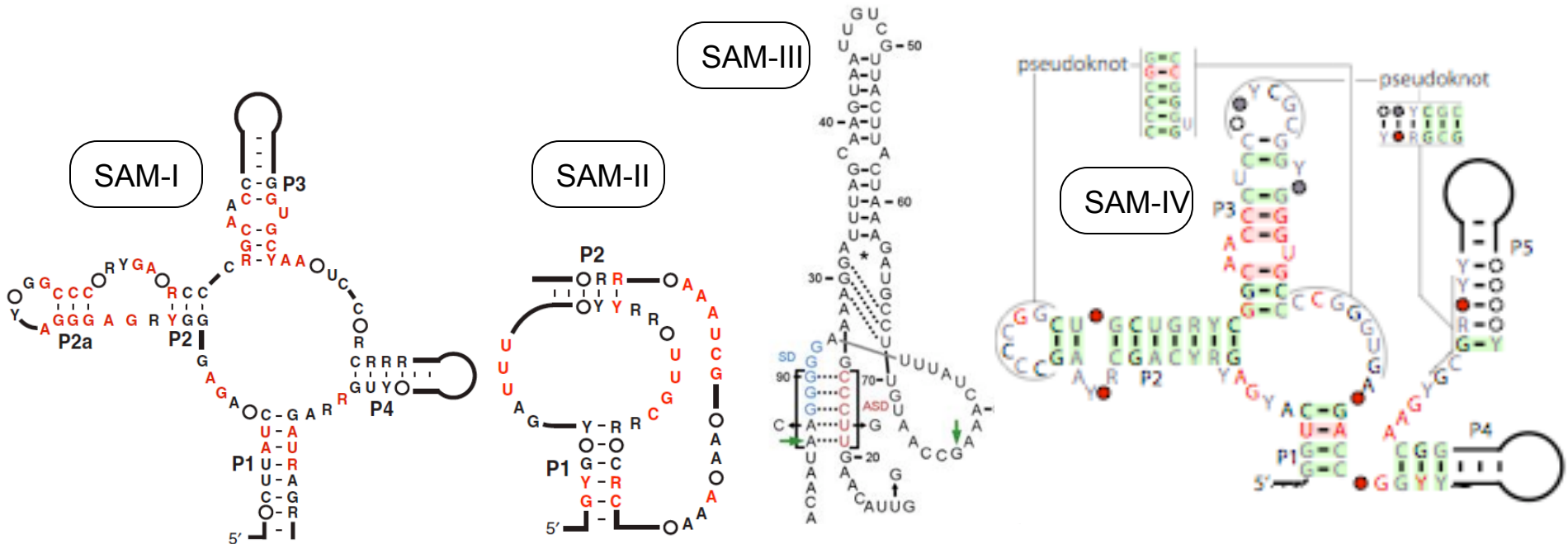
Corbino et al.,
Genome Biol. 2005

Alberts, et al, 3e.



The protein way

Riboswitch alternatives

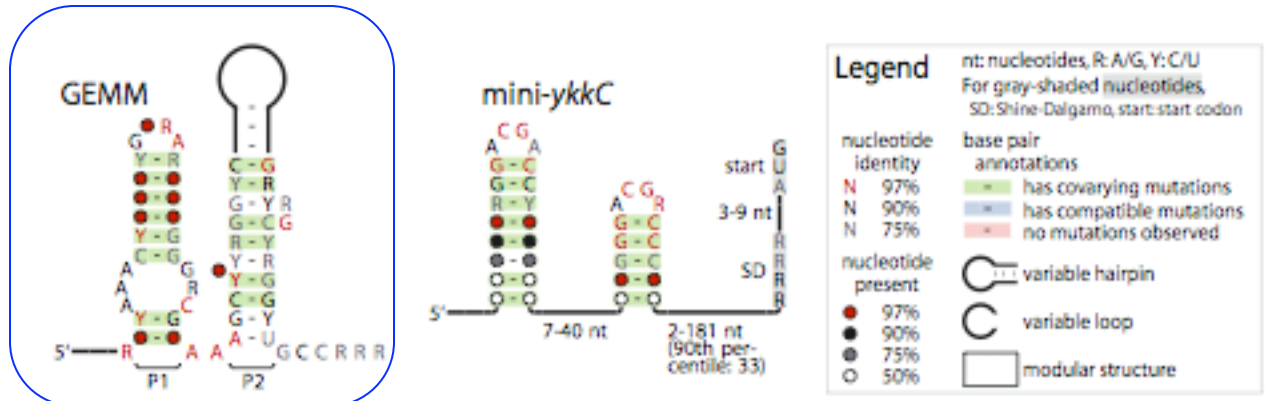


Grundy, Epshtein, Winkler et al., 1998, 2003

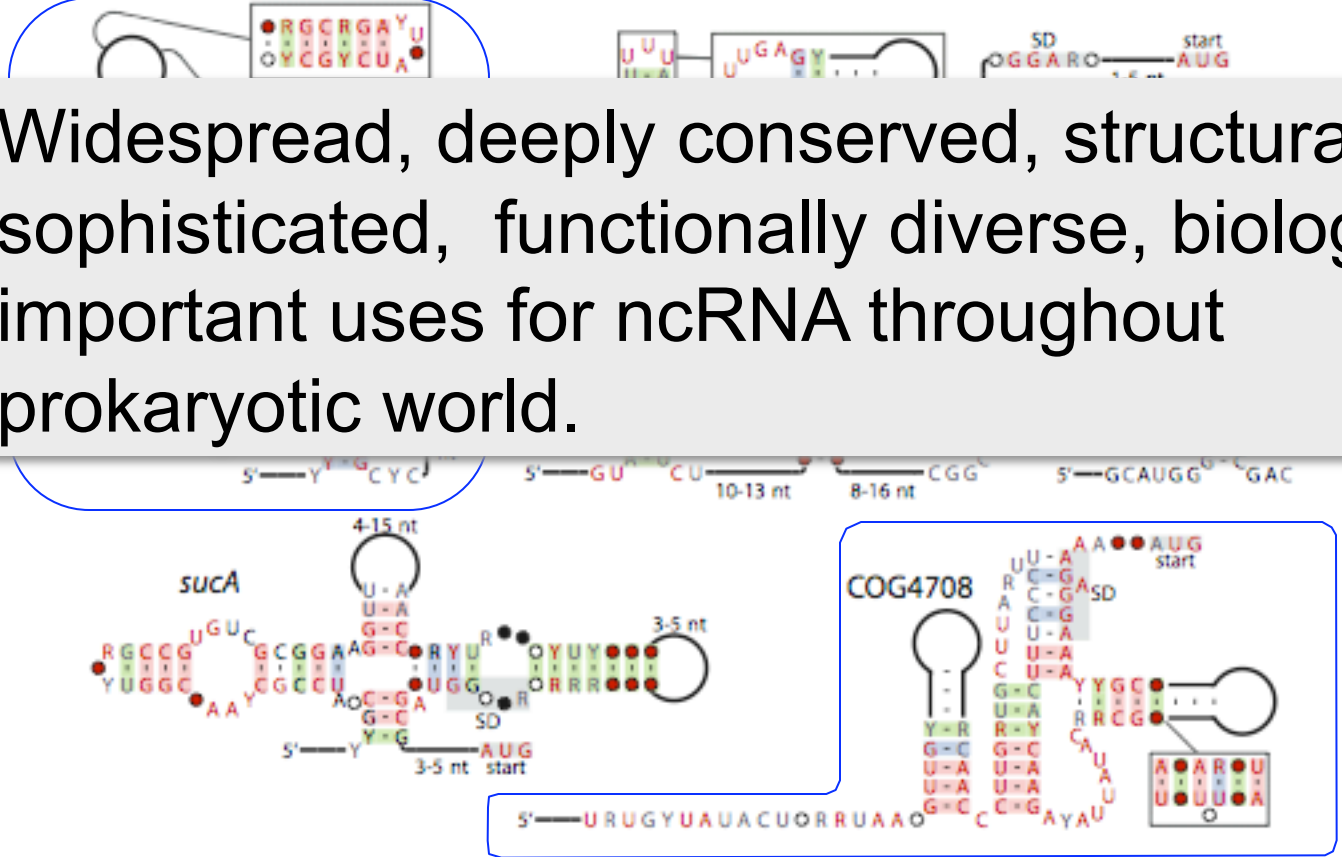
Corbino et al., Genome Biol. 2005

Fuchs et al., NSMB 2006

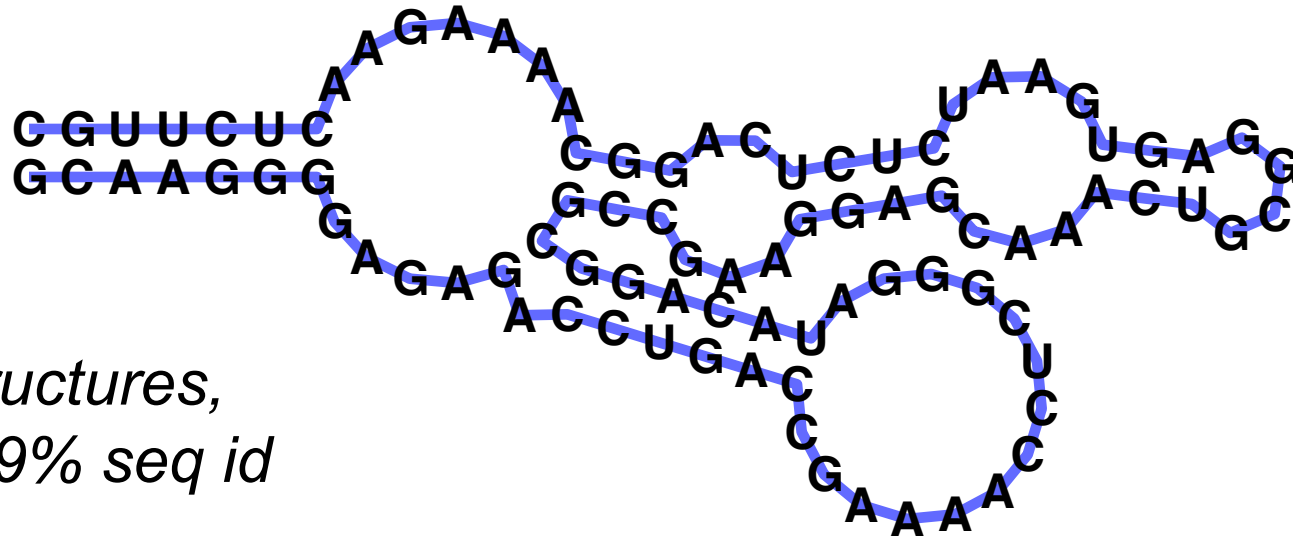
Weinberg et al., RNA 2008 ⁹



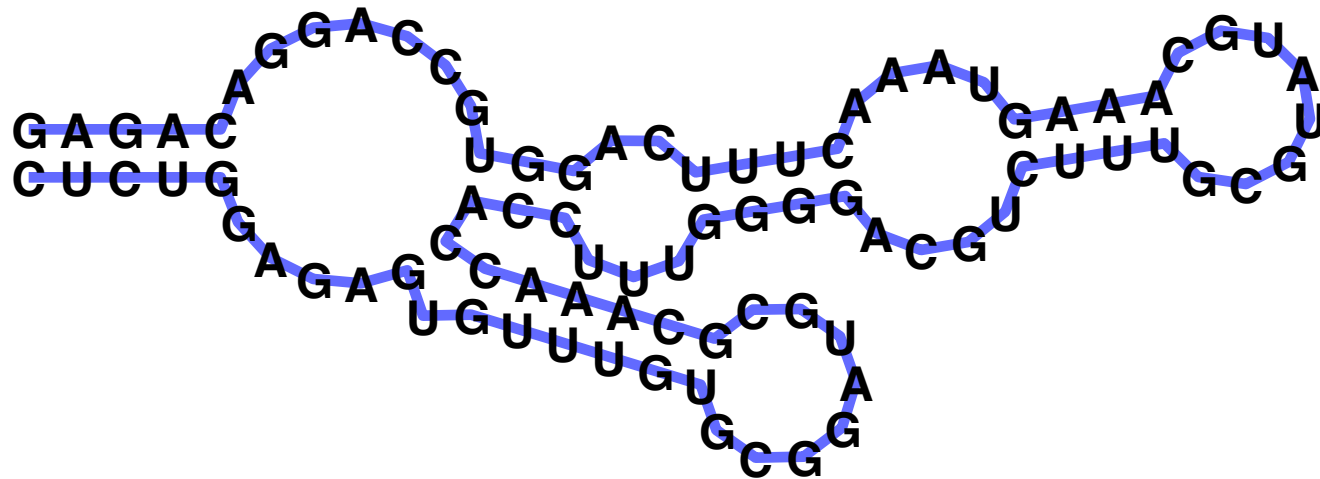
Widespread, deeply conserved, structurally sophisticated, functionally diverse, biologically important uses for ncRNA throughout prokaryotic world.



Why is RNA hard to deal with?



*Similar structures,
but only 29% seq id*



A: *Structure* often more important than *sequence*¹¹

Motif Description & Inference

RNA Motif Models

“Covariance Models” (Eddy & Durbin 1994)

aka profile stochastic context-free grammars

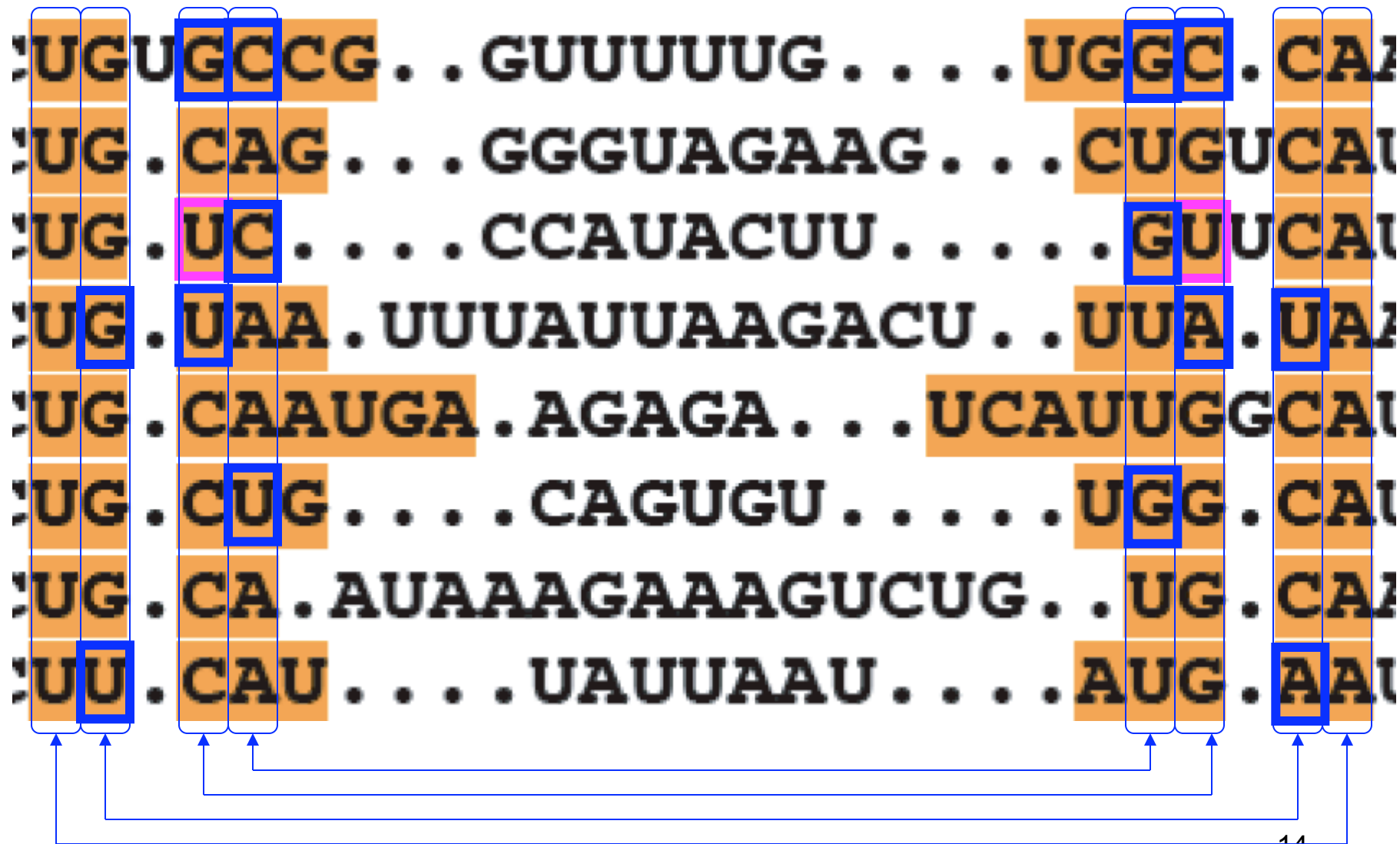
aka hidden Markov models on steroids

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

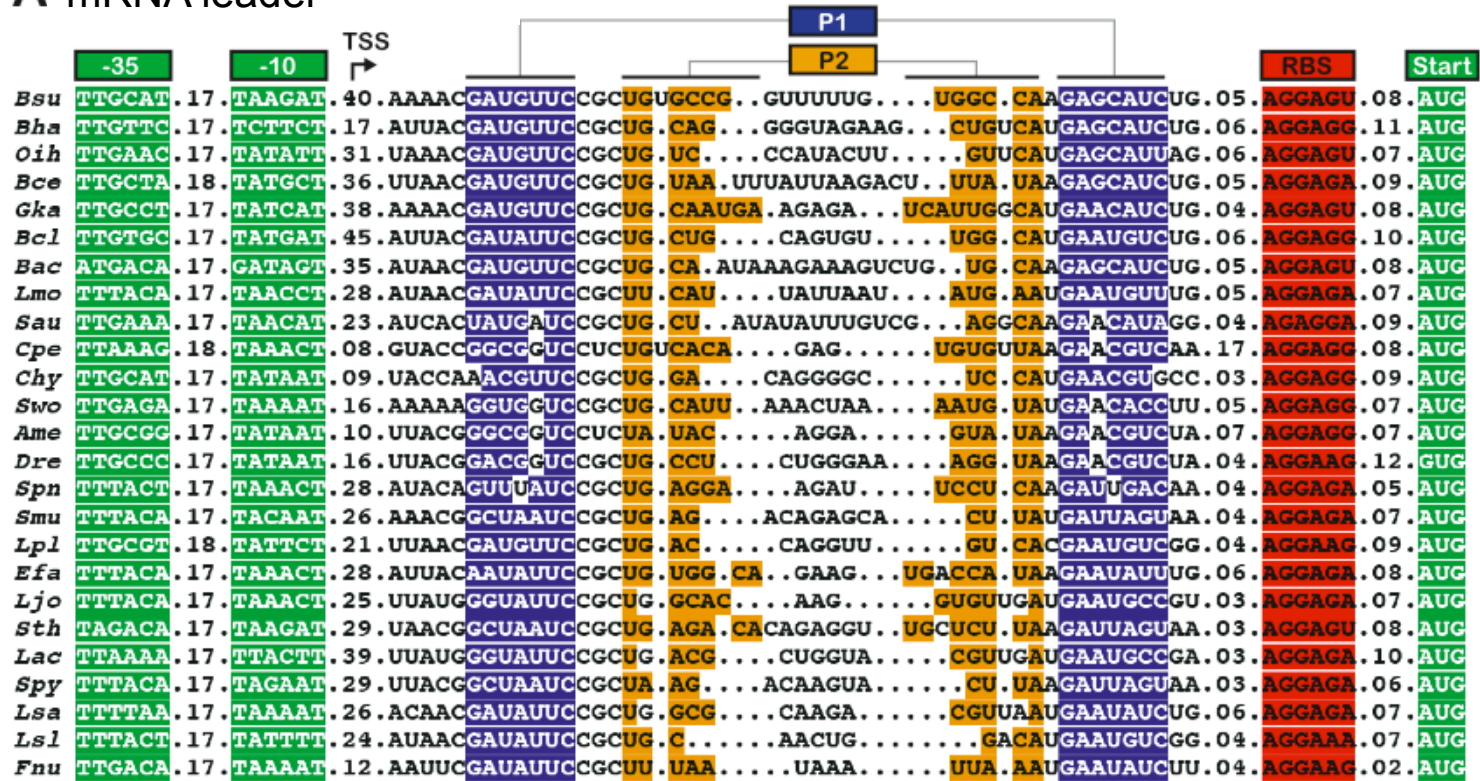
Pro: accurate

Con: model building hard, search sloooow

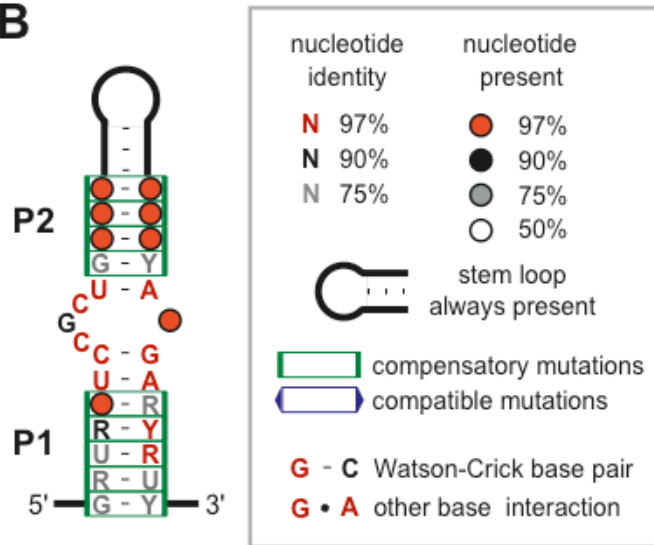
P2



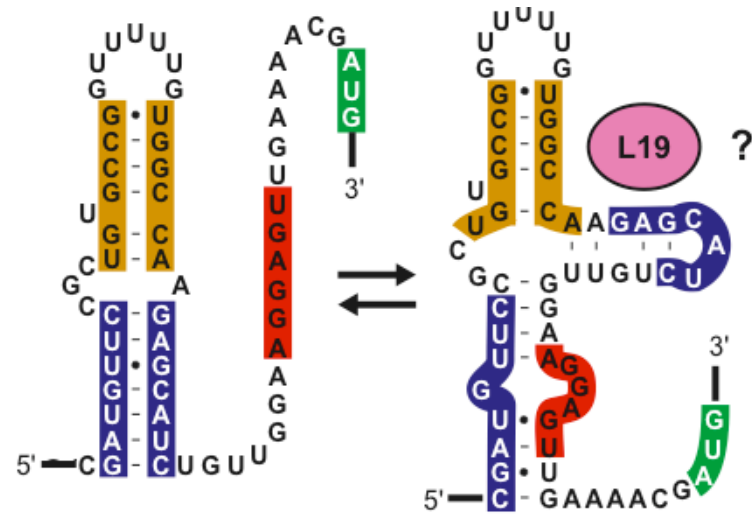
A mRNA leader



B



C mRNA leader switch?



Mutual Information

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

Max when *no* seq conservation but perfect pairing;

Expected score gain from modeling i & j as paired.

Given columns, finding optimal pairing *without pseudoknots* can be done by dynamic programming

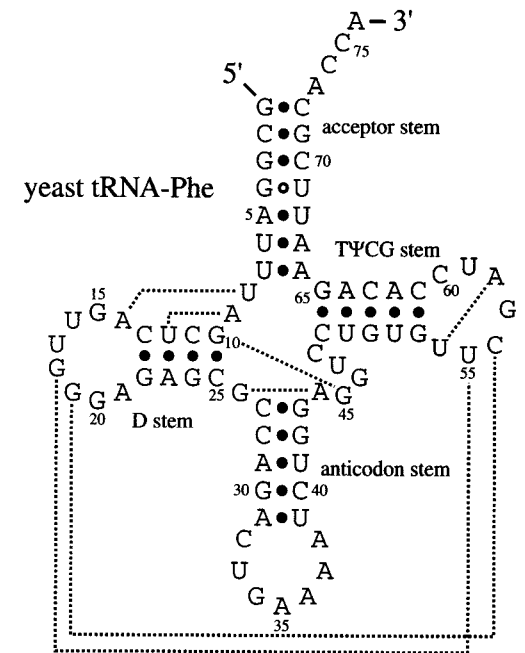
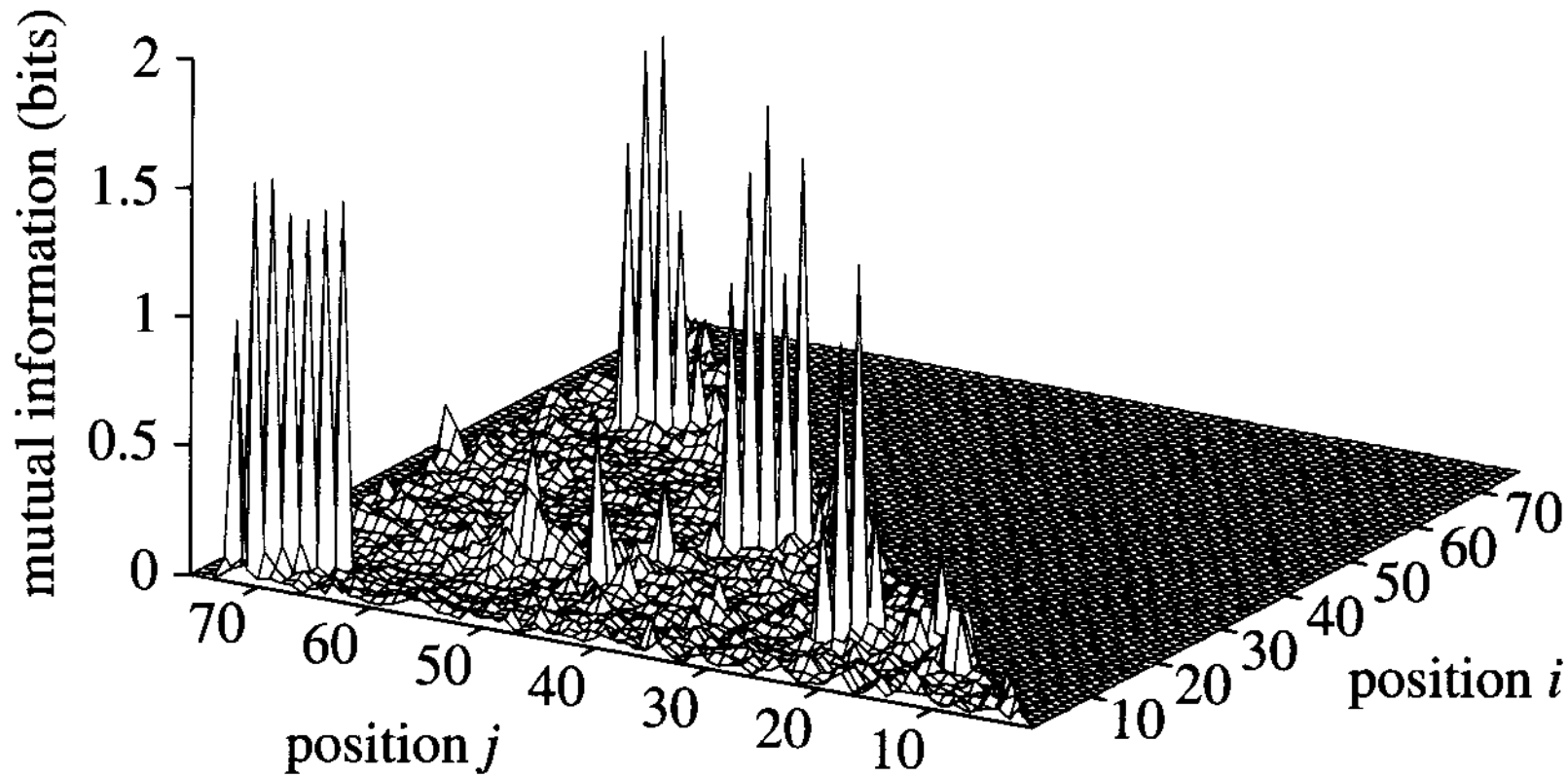


Figure 10.6 A mutual information plot of a tRNA alignment (top) shows four strong diagonals of covarying positions, corresponding to the four stems of the tRNA cloverleaf structure (bottom; the secondary structure of yeast phenylalanine tRNA is shown). Dashed lines indicate some of the additional tertiary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

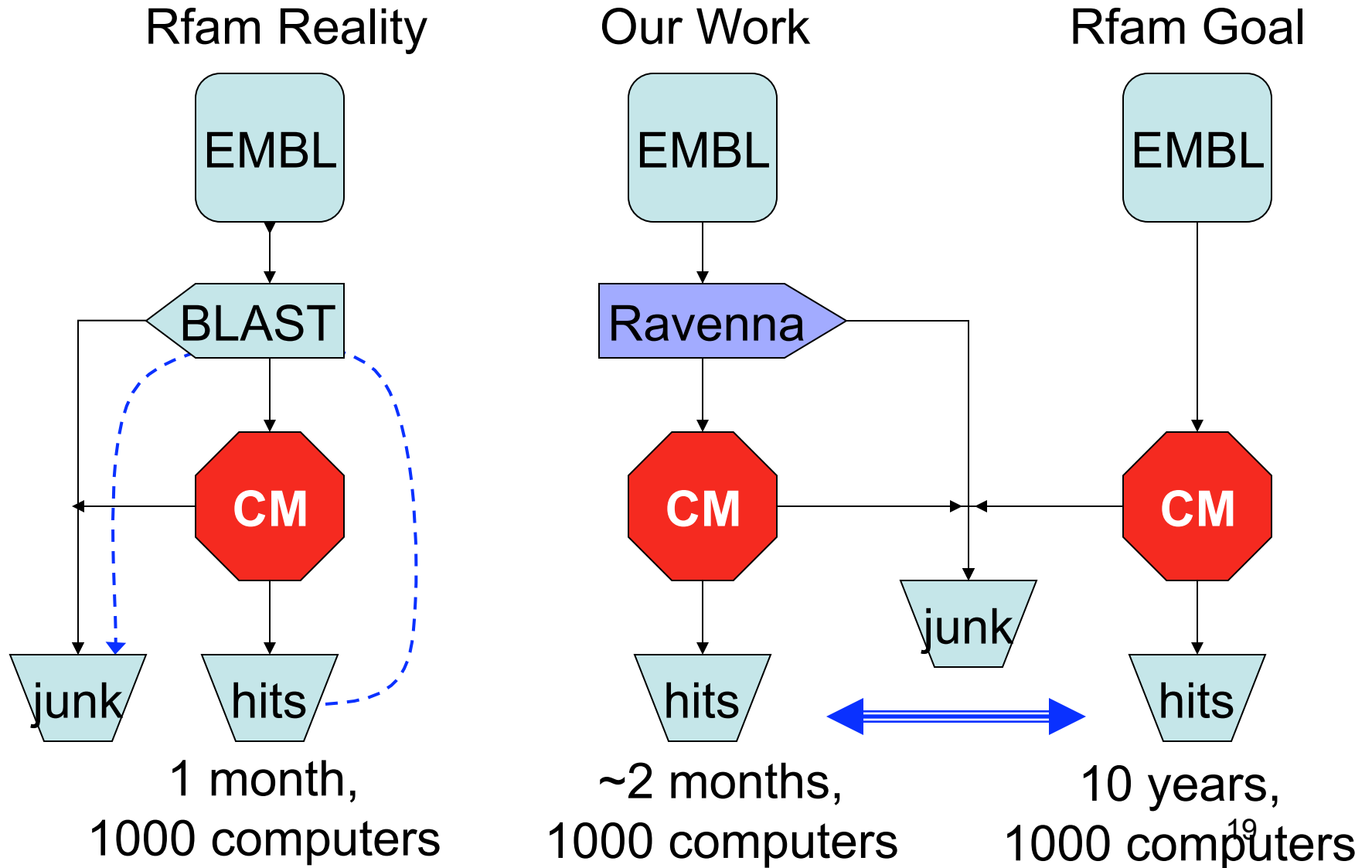
Fast Motif Search

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

Weinberg & Ruzzo

Recomb '04, ISMB '04, Bioinformatics '06

CM's are good, but slow



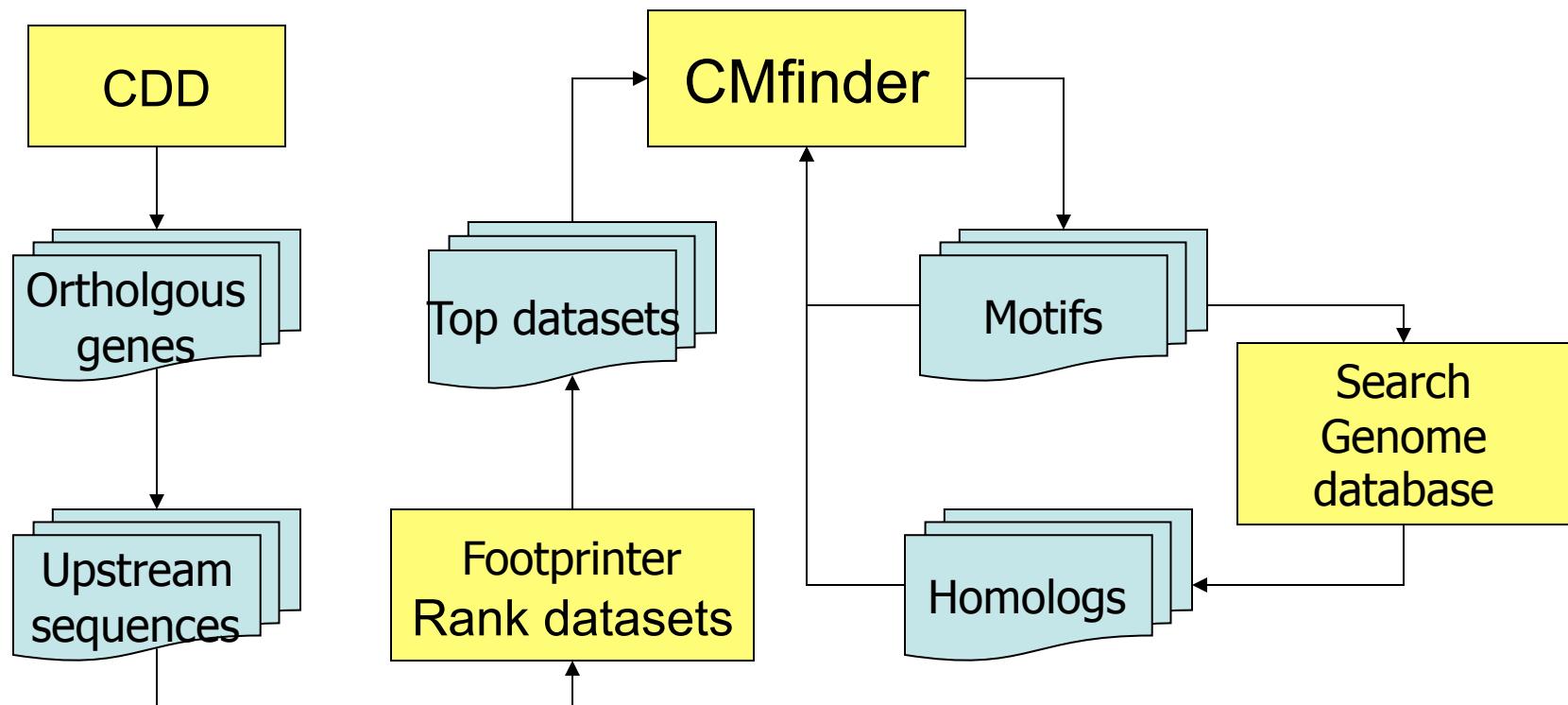
Results: New ncRNA's?

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

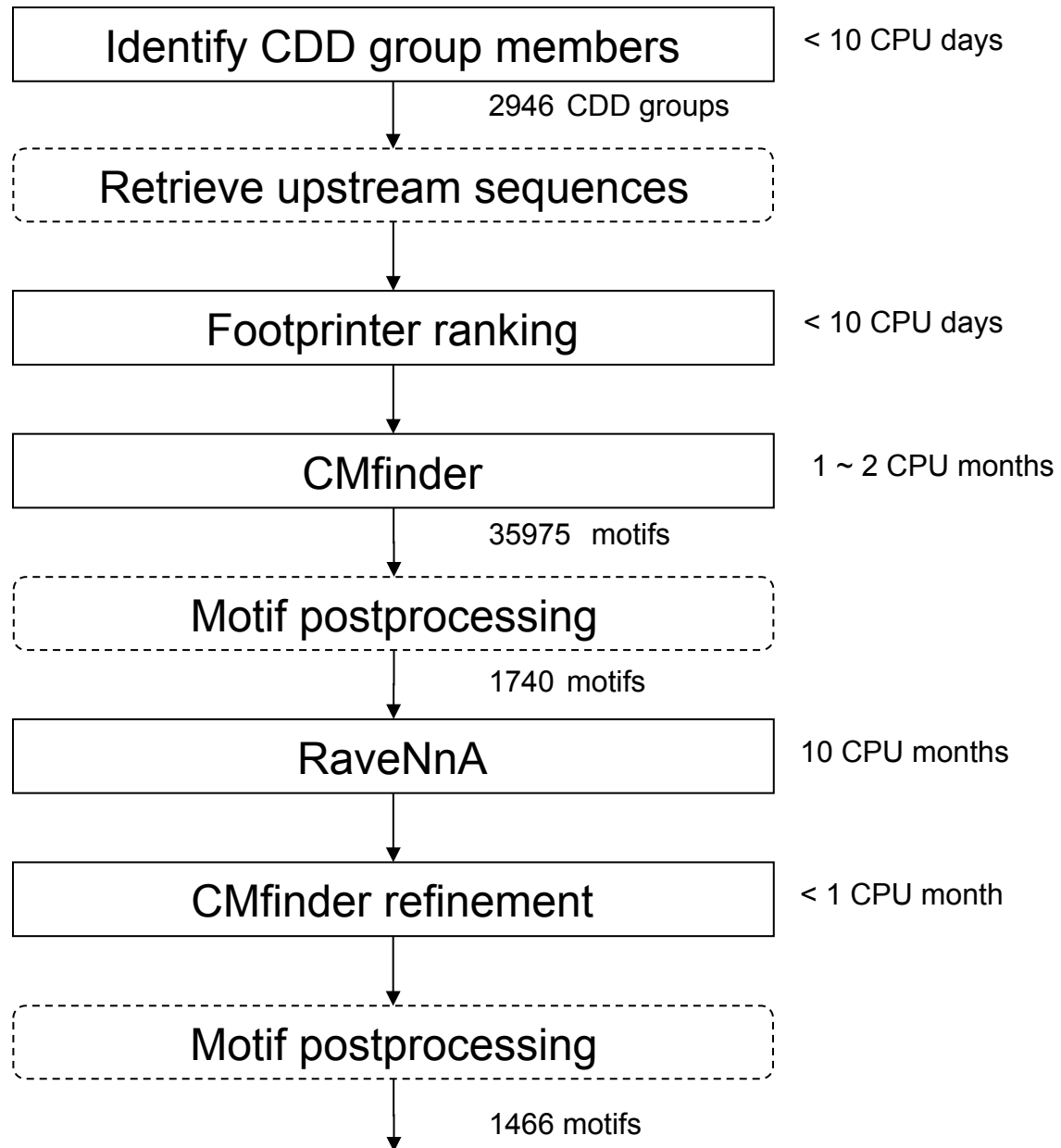
Motif Discovery In Prokaryotes

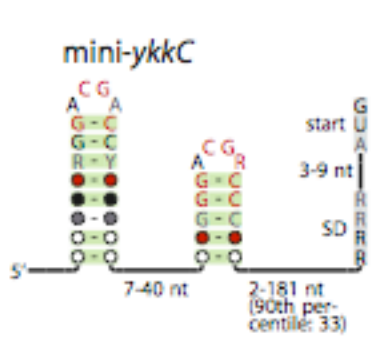
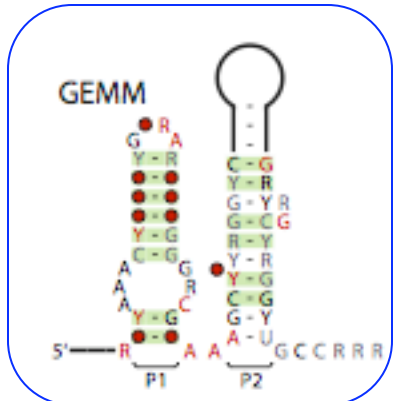
(Vertebrates too, but no time today...
see, e.g., Torarinsson, et al.
Genome Research, Jan 2008)

A pipeline for RNA motif genome scans



Yao, Barrick, Weinberg, Neph, Breaker, Tompa and Ruzzo. A Computational Pipeline for High Throughput Discovery of cis-Regulatory Noncoding RNA in Prokaryotes. *PLoS Computational Biology*. 3(7): e126, July 6, 2007. ²²





Legend

nt: nucleotides, R: A/G, Y: C/U
 For gray-shaded nucleotides,
 SD: Shine-Dalgarno, start: start codon

nucleotide identity

N	97%
N	90%
N	75%

nucleotide present

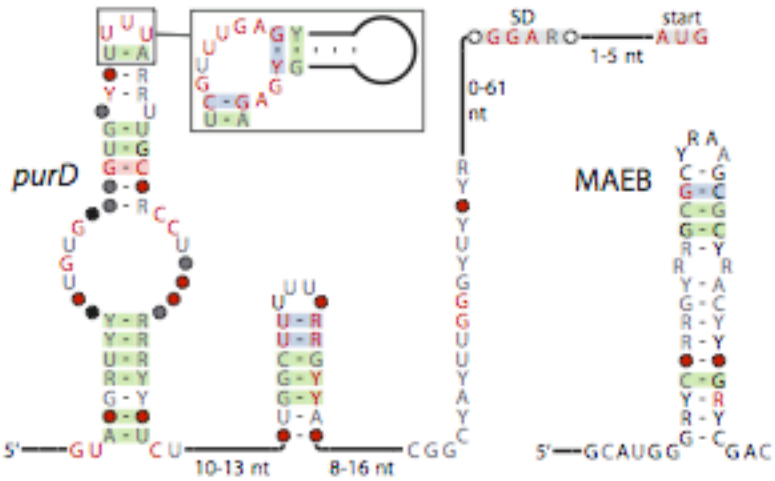
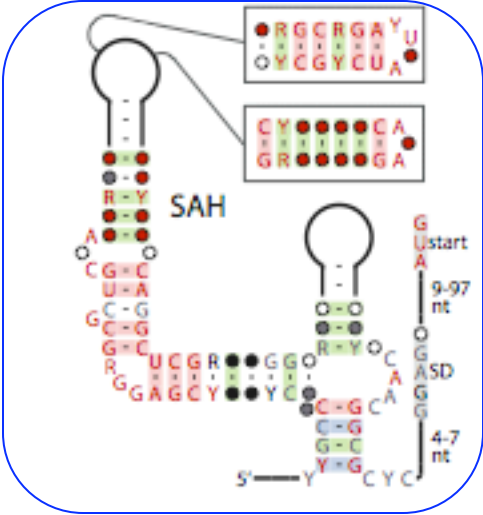
●	97%
●	90%
●	75%
○	50%

base pair annotations

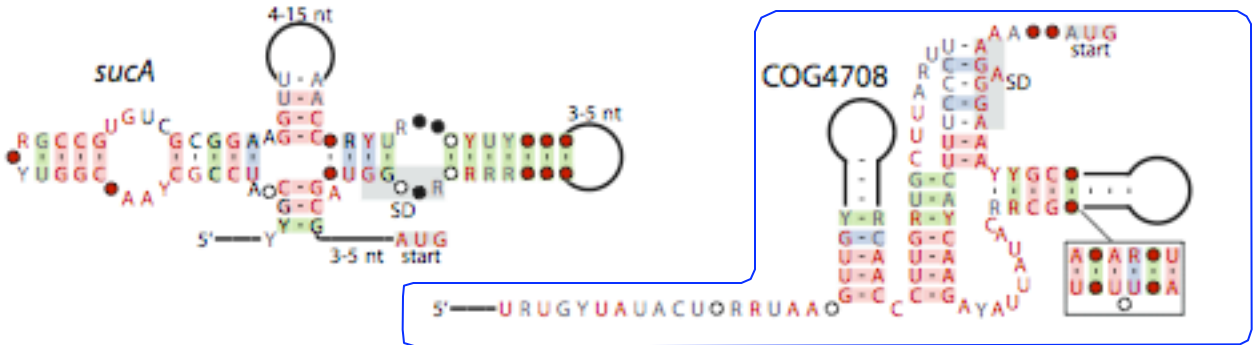
—	has covarying mutations
—	has compatible mutations
—	no mutations observed

nucleotide annotations

 	variable hairpin
	variable loop
 	modular structure



boxed = confirmed riboswitch (+2 more)



Summary

ncRNA - apparently widespread, much interest

Covariance Models - powerful but expensive

RaveNnA filtering - search ~100x faster with no/little loss

CMfinder - CM-based motif discovery in unaligned sequences

Pipelines integrating comp and bio for ncRNA discovery

Many vertebrate ncRNAs? *structural*, not seq conservation;
functional significance unclear

BIG CPU demands...

Still need for further methods development & application

Course Wrap Up

Modern biology is suddenly very data-rich

Mathematical & computational tools needed

We showed: sequence modeling, alignment & search, phylogeny, linkage/association mapping, some data bases

Python is a good tool for doing much of this

There's lots more!

Check out, e.g., GENOME 540/I, CSE 527...

We hope you enjoyed it.

Thanks!