

# **Modeling Structural RNA Families with Infernal**

Eric Nawrocki

Sean Eddy's Lab



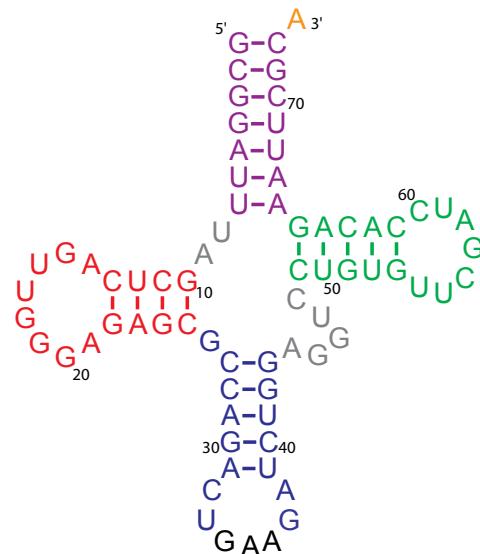
# Many functional RNAs adopt a conserved 3-dimensional structure

Three representations of a transfer RNA:

## Primary sequence

GC<sub>1</sub>GGAUUUAAGCUCAGUUGGG  
AGAGC<sub>2</sub>GCCAGACU<sub>3</sub>GAAGAUC  
UGGAGGUCCUGUGUUUCGAUC  
CACAGAAUUCGCA<sub>4</sub>

## Secondary structure

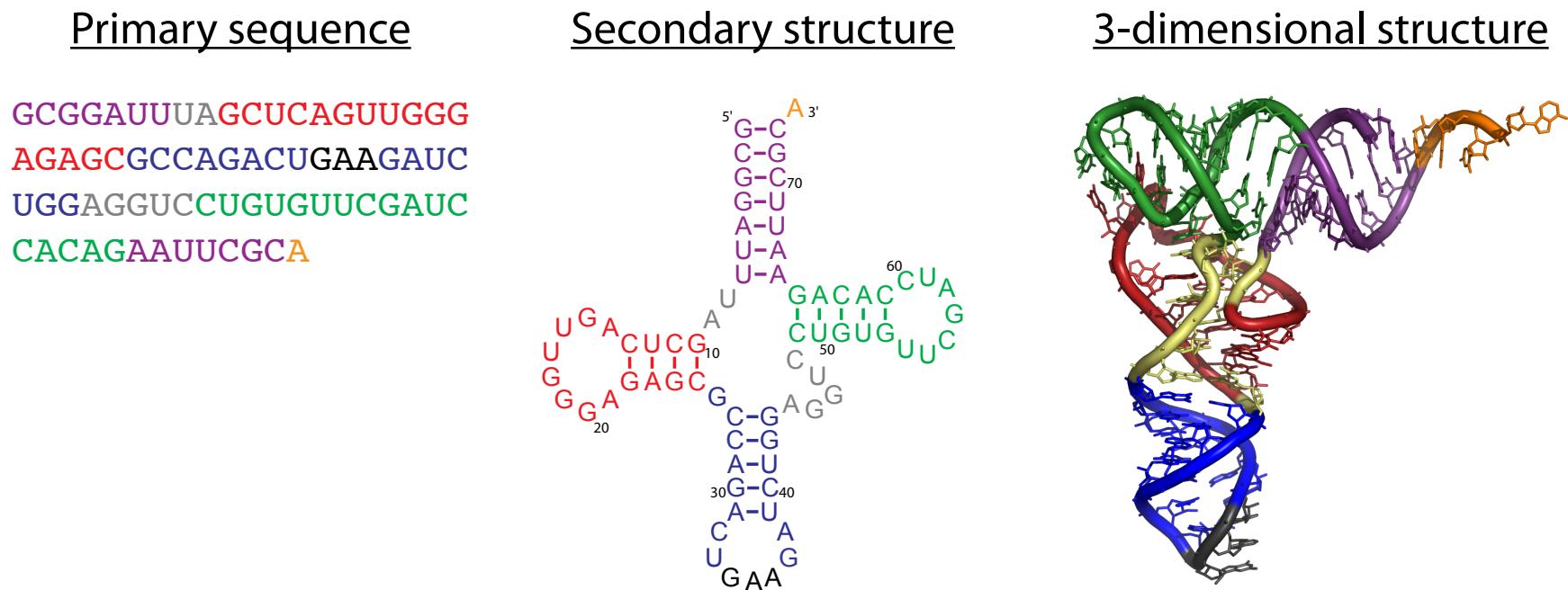


## 3-dimensional structure



# Many functional RNAs adopt a conserved 3-dimensional structure

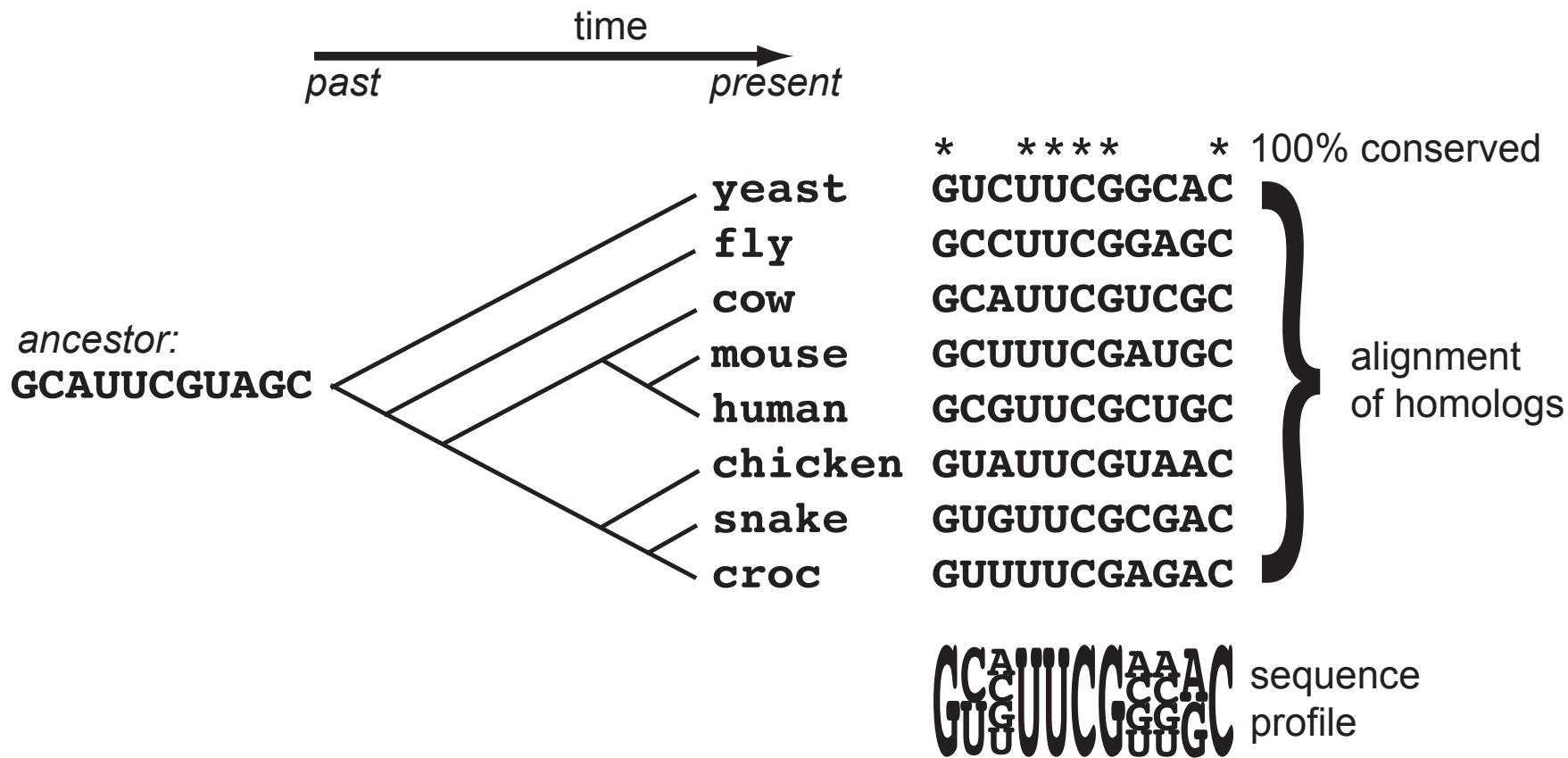
Three representations of a transfer RNA:



- BLAST: given a single sequence, search genomes for similar sequences.
- BLAST cannot take advantage of:
  - secondary structure
  - sequence conservation, which varies across the gene

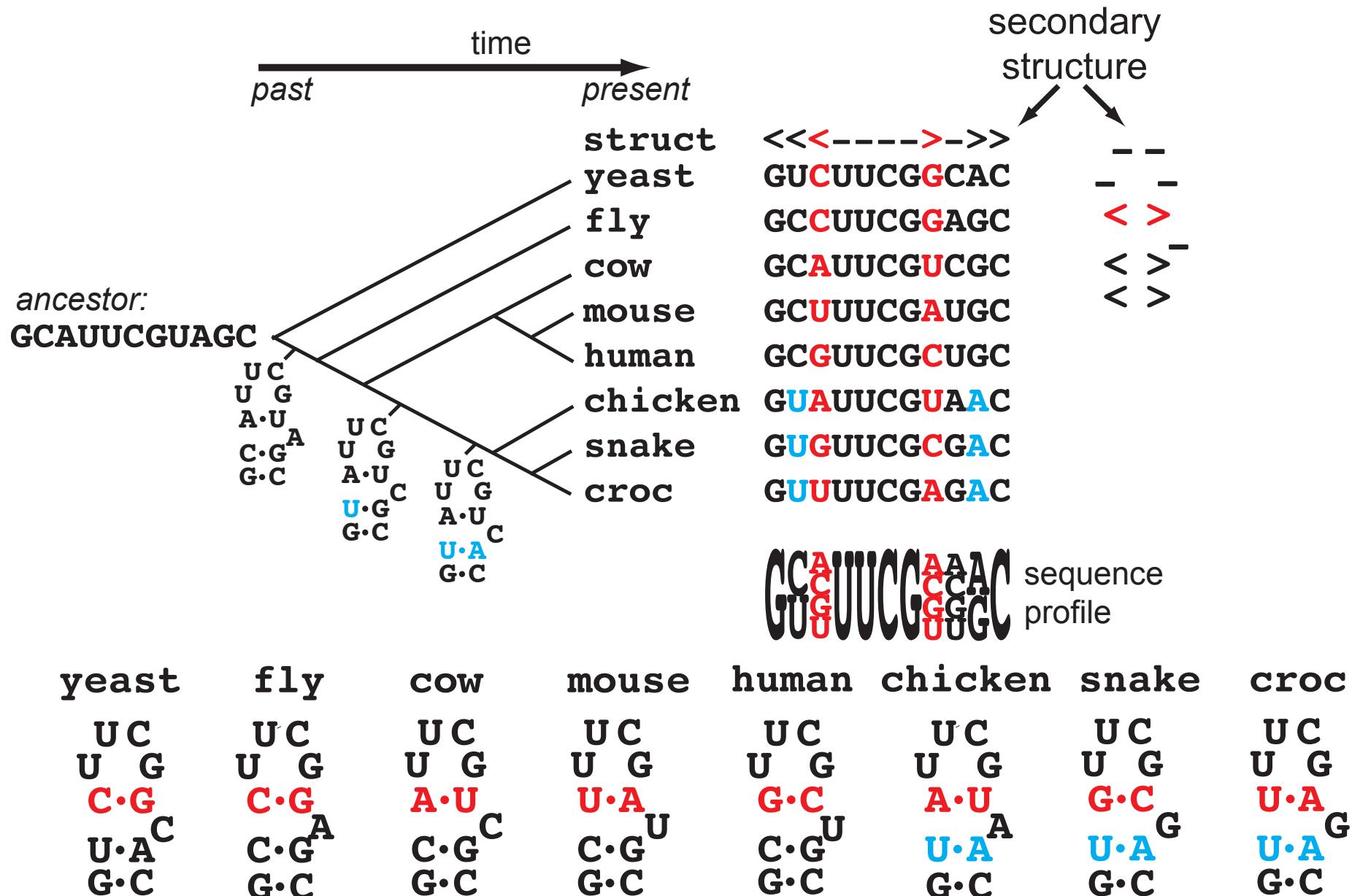
# Sequence conservation provides information for homology searches

Conservation levels vary across alignment columns.

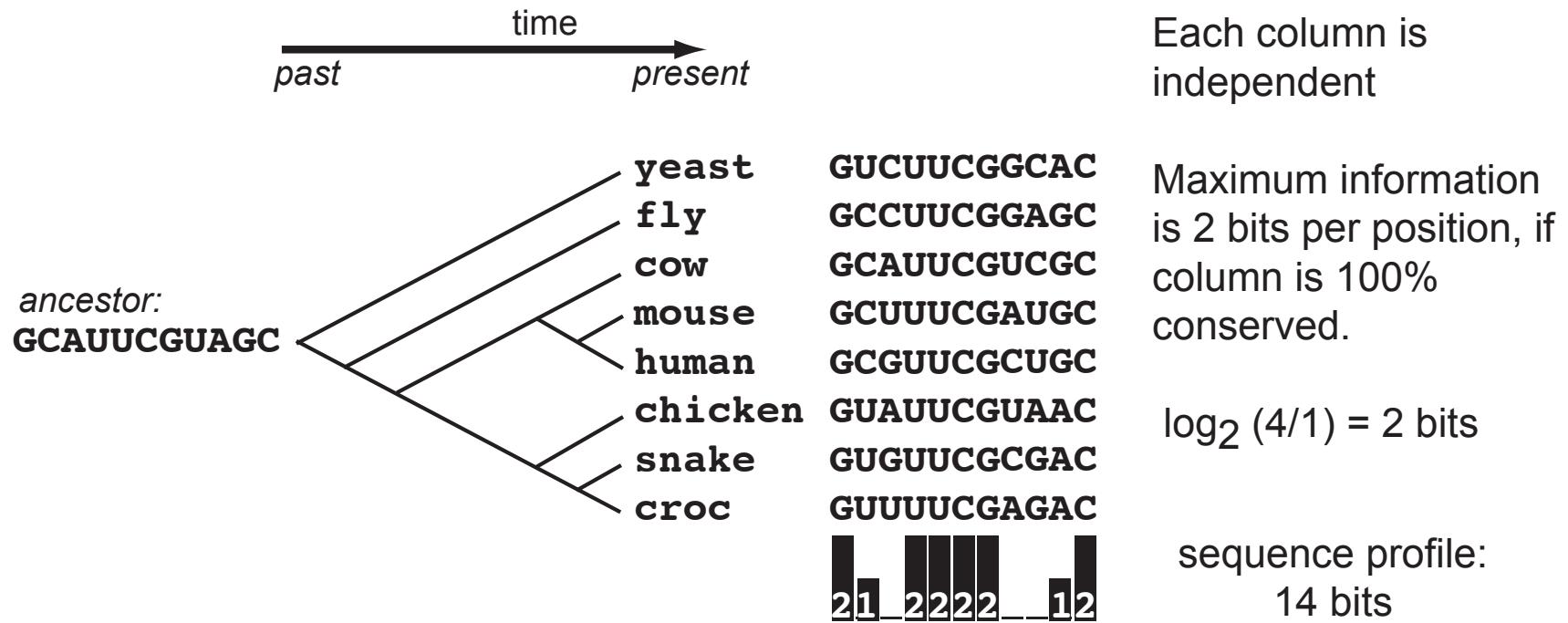


# Structure conservation provides additional information

Base-paired positions covary  
to maintain Watson-Crick complementarity.

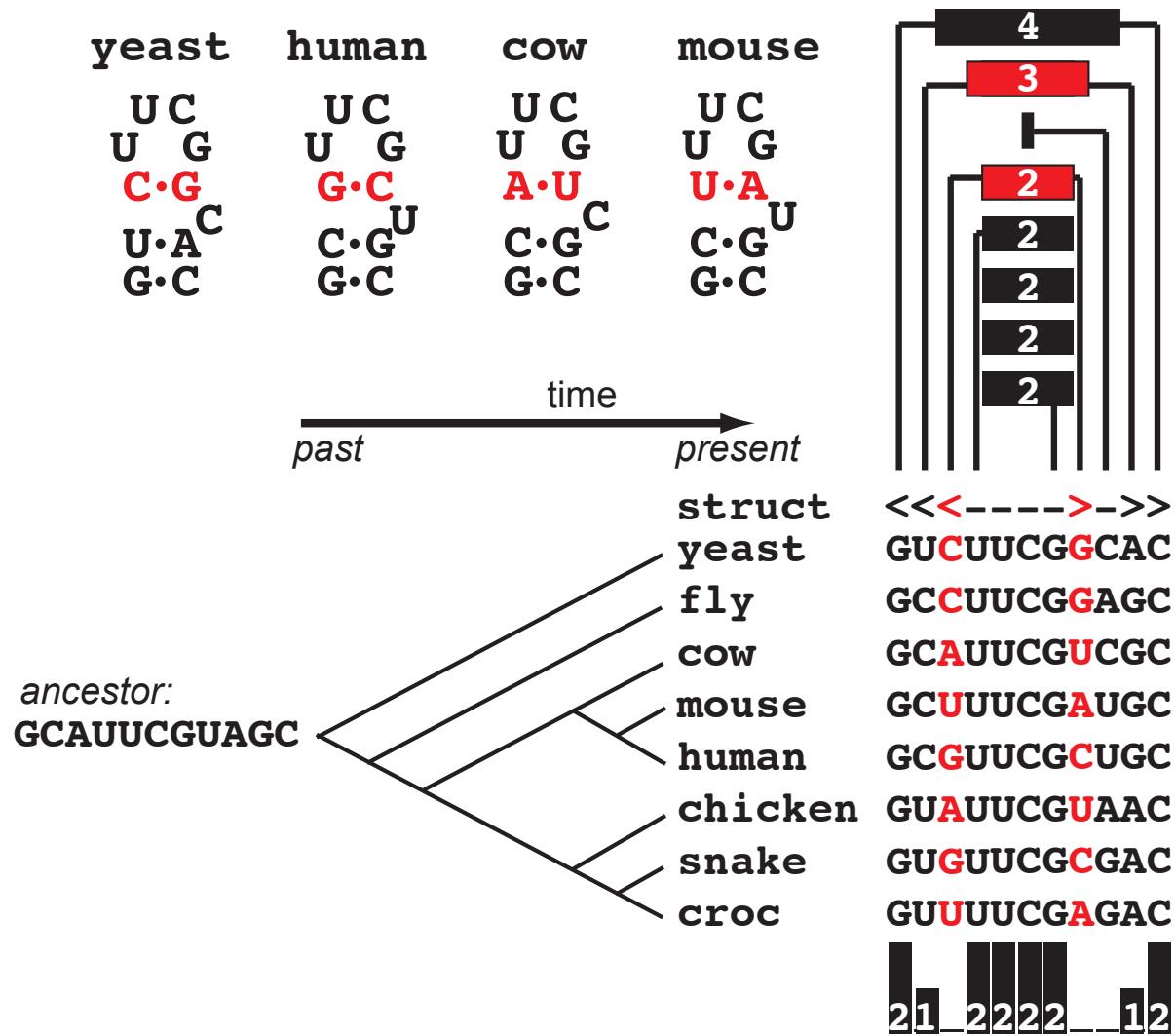


# Amount of information in a profile can be measured in bits



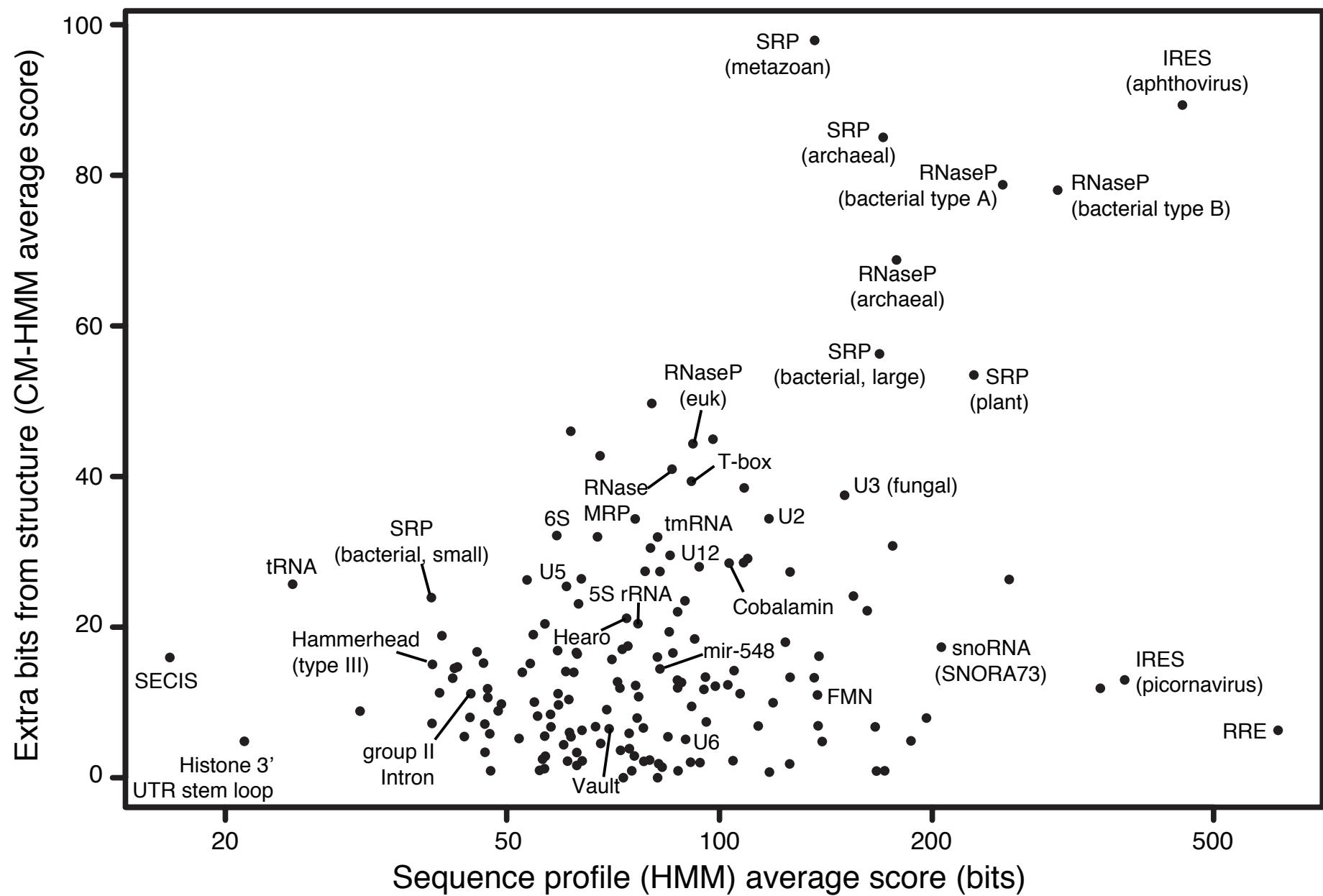
expect a match by chance: 1 in  $2^{14}$  nt  $= \sim 16$  Kb

**Structure contributes additional information from covariation**



expect a match by chance: 1 in  $2^{17}$  nt  $\approx$  130 Kb  
reducing expected false positives by  $2^3$  = 8-fold

# Levels of sequence and structure conservation in RNA families



# Eddy lab software for profile probabilistic models (since 1994)

	sequence profiles	sequence and structure profiles
models	profile HMMs	covariance models (CMs)
software	HMMER	Infernal
main use	proteins, repetitive DNA elements	structural RNAs
databases	Pfam and Dfam (14831 and 1132 entries)	Rfam (2450 families)
performance for RNAs	faster but less accurate	slower but more accurate



<http://hmmer.janelia.org>  
Eddy, SR. PLoS Comp. Biol.,  
7:e1002195, 2011.  
Eddy, SR. PLoS Comp. Biol.,  
4:e1000069, 2008.  
Eddy, SR. Bioinformatics,  
14:755-763, 1998.

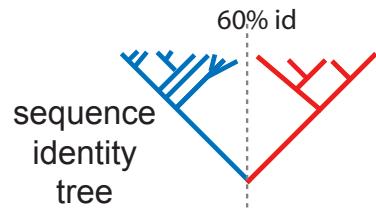
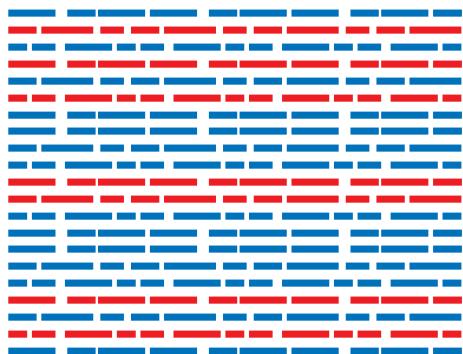


<http://infernal.janelia.org>  
Nawrocki EP, Eddy SR  
Bioinformatics,  
29:2933-2935, 2013.  
Eddy SR, Durbin R.  
Nucleic Acids Research,  
22:2079-2088, 1994.

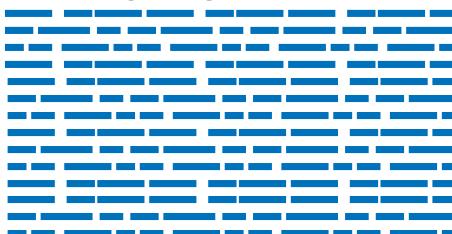
# Is the added complexity worth it?

## RMARK: a challenging internal RNA homology search benchmark

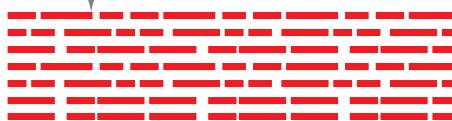
Rfam seed alignment:



training alignment

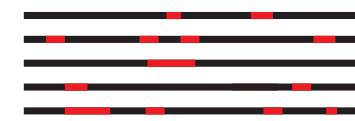


no train/test sequence pair is > 60% identical



test sequences

embed in  
pseudo-genome

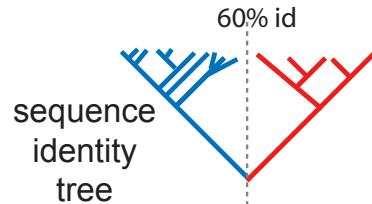
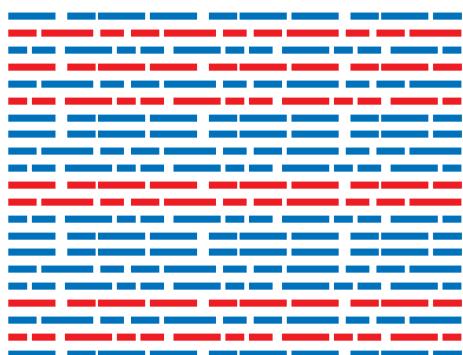


10 1Mb sequences  
with 780 embedded  
test seqs from 106 families

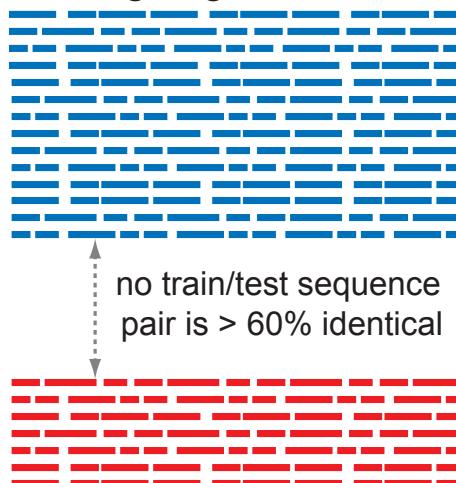
# Is the added complexity worth it?

## RMARK: a challenging internal RNA homology search benchmark

Rfam seed alignment:



training alignment



test sequences

profile  
(CM or HMM)

BLAST

search

embed in  
pseudo-genome



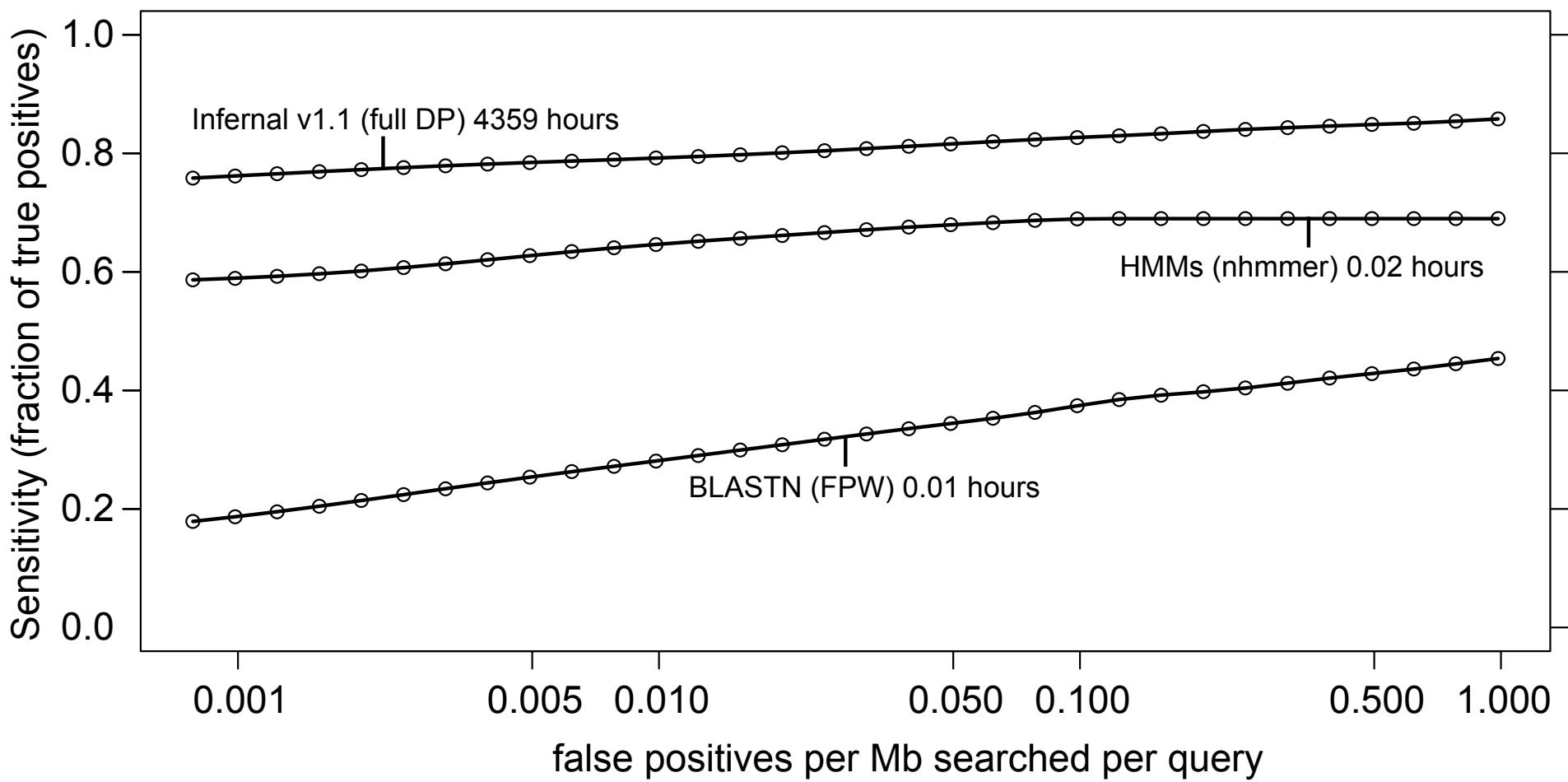
10 1Mb sequences  
with 780 embedded  
test seqs from 106 families

E=1E-40 132.53 bits rmark7 OLE 340023 339402 +  
...

E=0.0013 32.3 bits rmark3 6S 10135 10261 +  
E=0.0026 27.6 bits rmark6 tRNA 789278 789466 +  
E= 0.0061 28.3 bits rmark2 Cobalamin 32032 31787 -  
E=0.0231 25.4 bits rmark 6 FALSE 673200 673340 +  
E=0.0670 25.3 bits rmark6 tRNA 789278 789116 -  
...

E=103.3 16.4 bits rmark 4 FALSE 783222 782803 -

# Infernal outperforms primary-sequence based methods on our benchmark (and others\*, not shown)

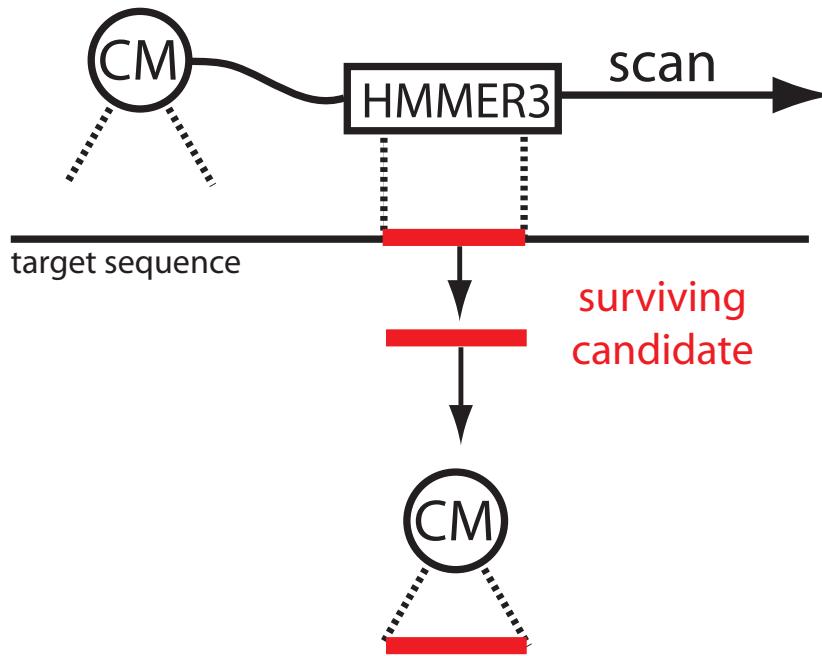


Nawrocki EP, Eddy SR. Bioinformatics, 29:2487-2489, 2013.

\*Freyhult EK, Bollback JP, Gardner PP. Genome Res. 2007 17: 117-125.

# Filter target database using profile HMMs\*

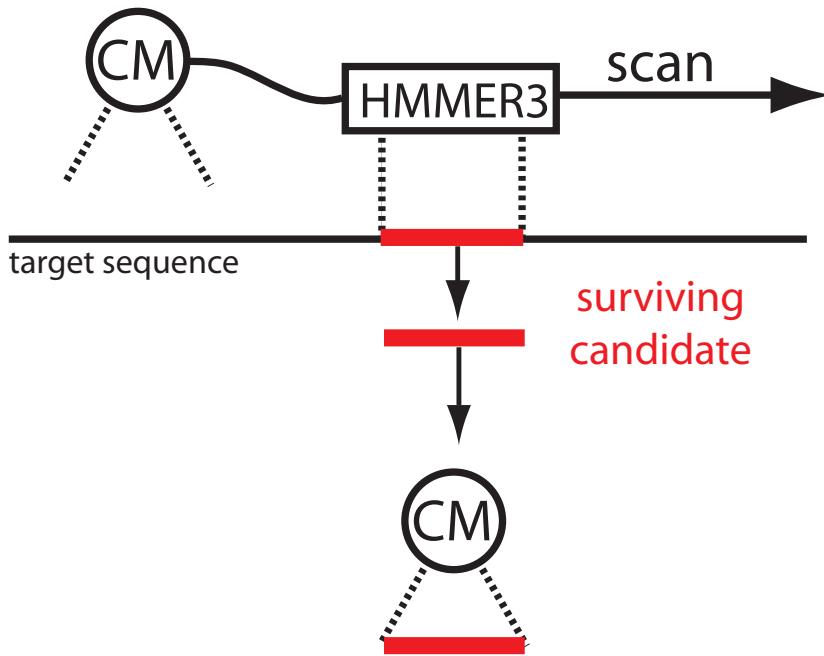
HMM filter first pass



surviving  
candidate

# Filter target database using profile HMMs\*

HMM filter first pass



- Even if we filter out 99% of the database (for up to 100X acceleration), searches will still be too slow.
- CM step needs to be accelerated.

# Accelerating CM alignment step 1: align sequence with HMM

## Accelerating CM alignment step 2: HMM posterior decoding to get confidence estimates

yeast GUGUUCGCUAC  
human -UCUUCGGCG-  
fly AGAUU-GUACU

1            5            11

**new sequence:**  
CUGUUCGCAAG

**probability "correct":**

1.0			
.8			
.6			
.4			
.2			
0			

C	C										
U		U									
G			G								
U				U							
U					U						
C						C					
G							G				
C								C			
A									A		
G											G

1            5            11

C											
U											
G											
U											
U											
C											
G											
C											
A											
G											

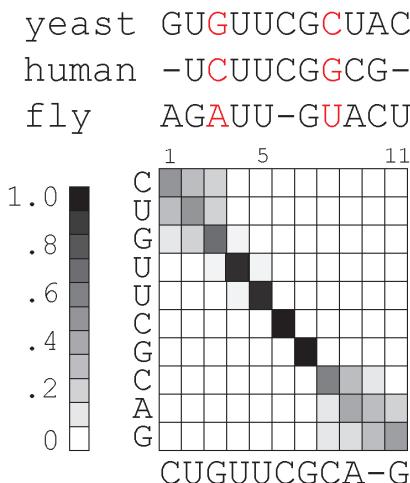
1            5            11

CUGUUCGCA-G

# Accelerating CM alignment step 3: use HMM alignment confidence to constrain CM alignment\*

**HMMs -**

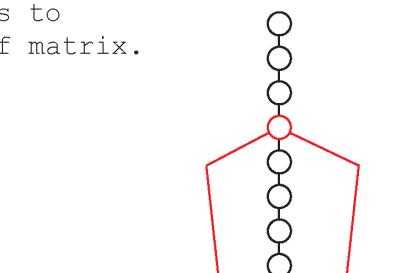
Each column of seed alignment corresponds to a column of matrix.



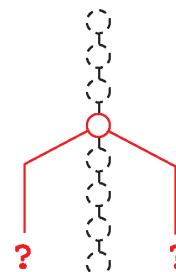
**CMs -**

Each column of seed alignment corresponds to a state.

yeast	human	fly
U C	U C	U
U G	U G	U G
<b>G•C</b>	<b>C•G</b>	<b>A•U</b>
U•A U	U•G C	G•C A
G•C		A•U



struct <<----->->  
 yeast GUGUUCG**C**UAC  
 human -UCUUCGG**G**CG-  
 fly AG**A**UU-G**U**ACU



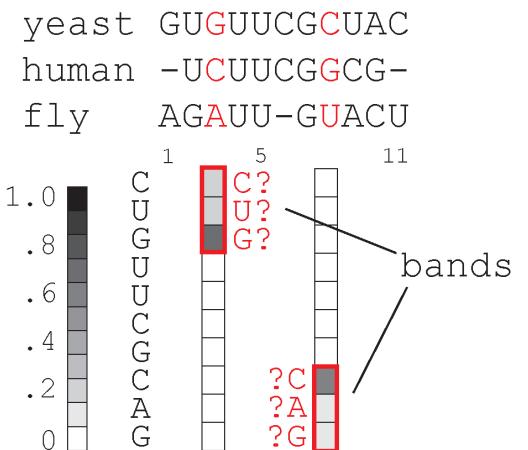
CUGUUCGCAG

45 possibilities

# Accelerating CM alignment step 3: use HMM alignment confidence to constrain CM alignment\*

**HMMs -**

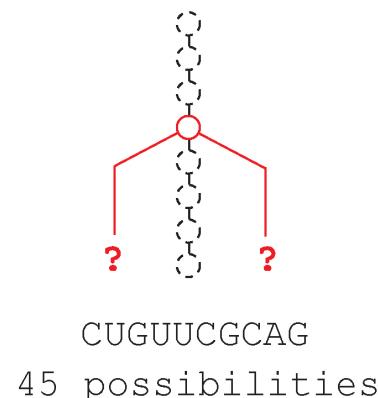
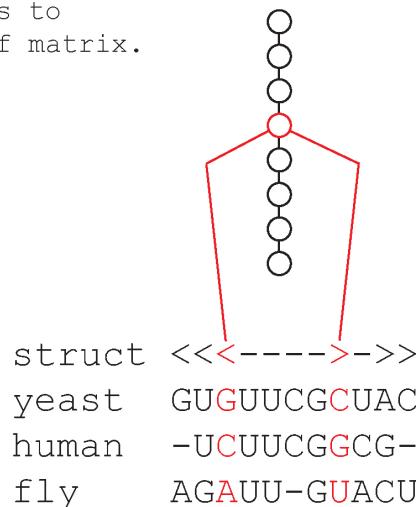
Each column of seed alignment corresponds to a column of matrix.



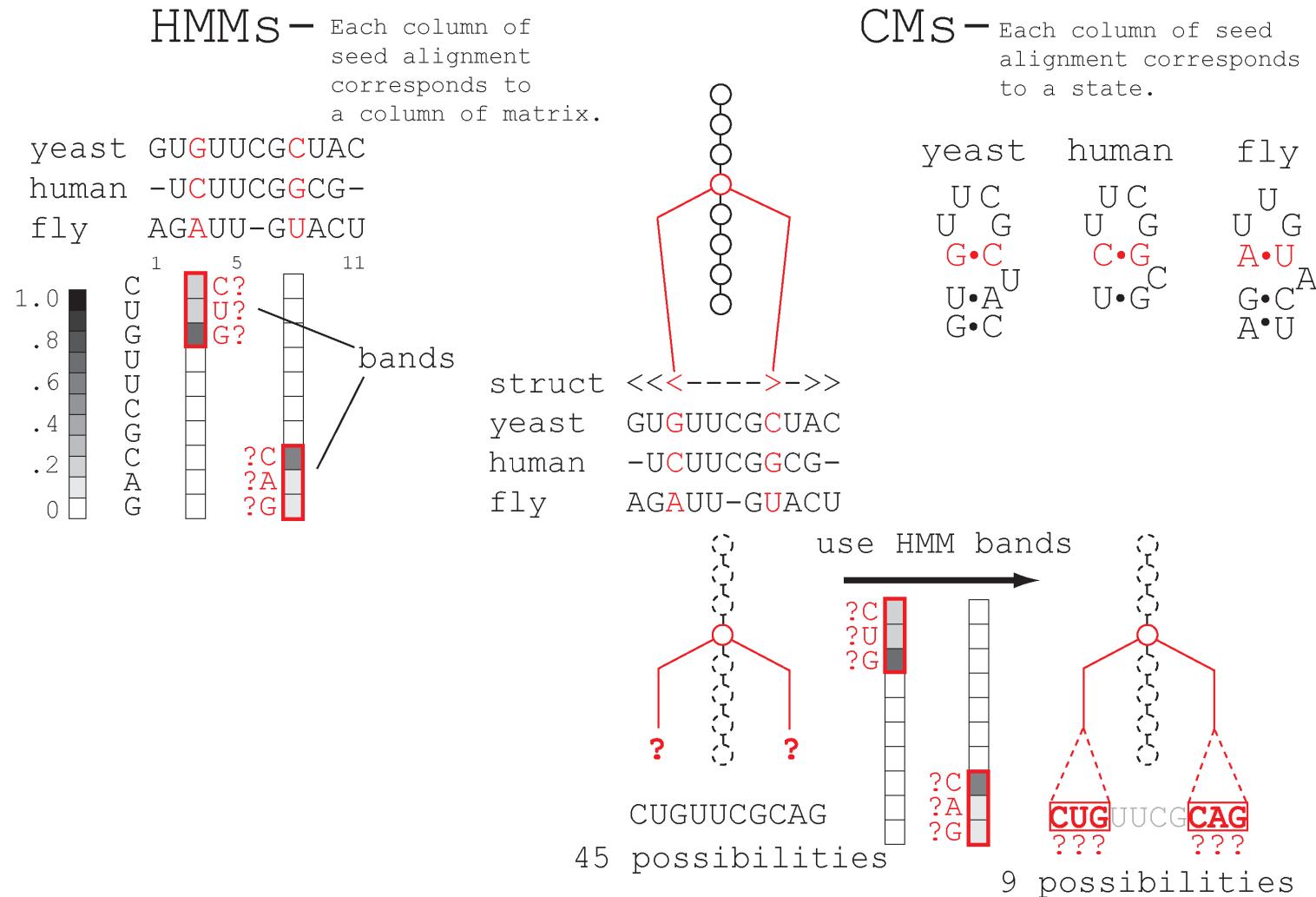
**CMs -**

Each column of seed alignment corresponds to a state.

yeast	human	fly
U C	U C	U
U G	U G	U G
<b>G•C</b>	<b>C•G</b>	<b>A•U</b>
U•A U	U•G C	G•C A
G•C		A•U

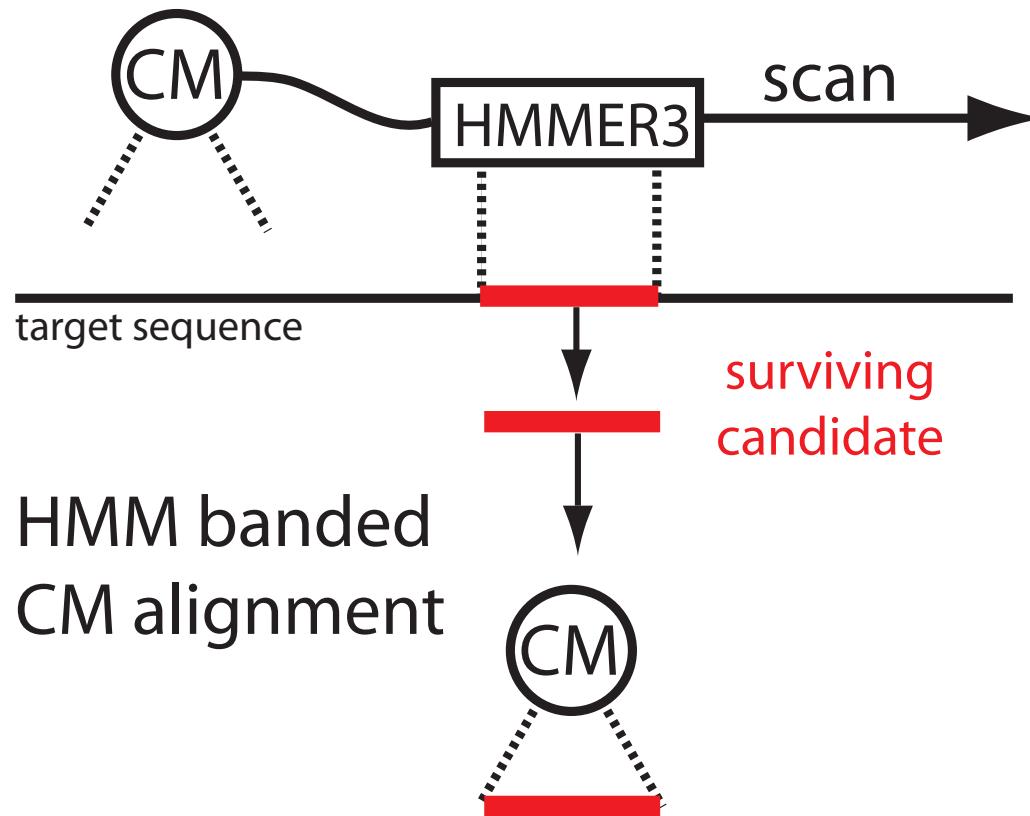


# Accelerating CM alignment step 3: use HMM alignment confidence to constrain CM alignment\*

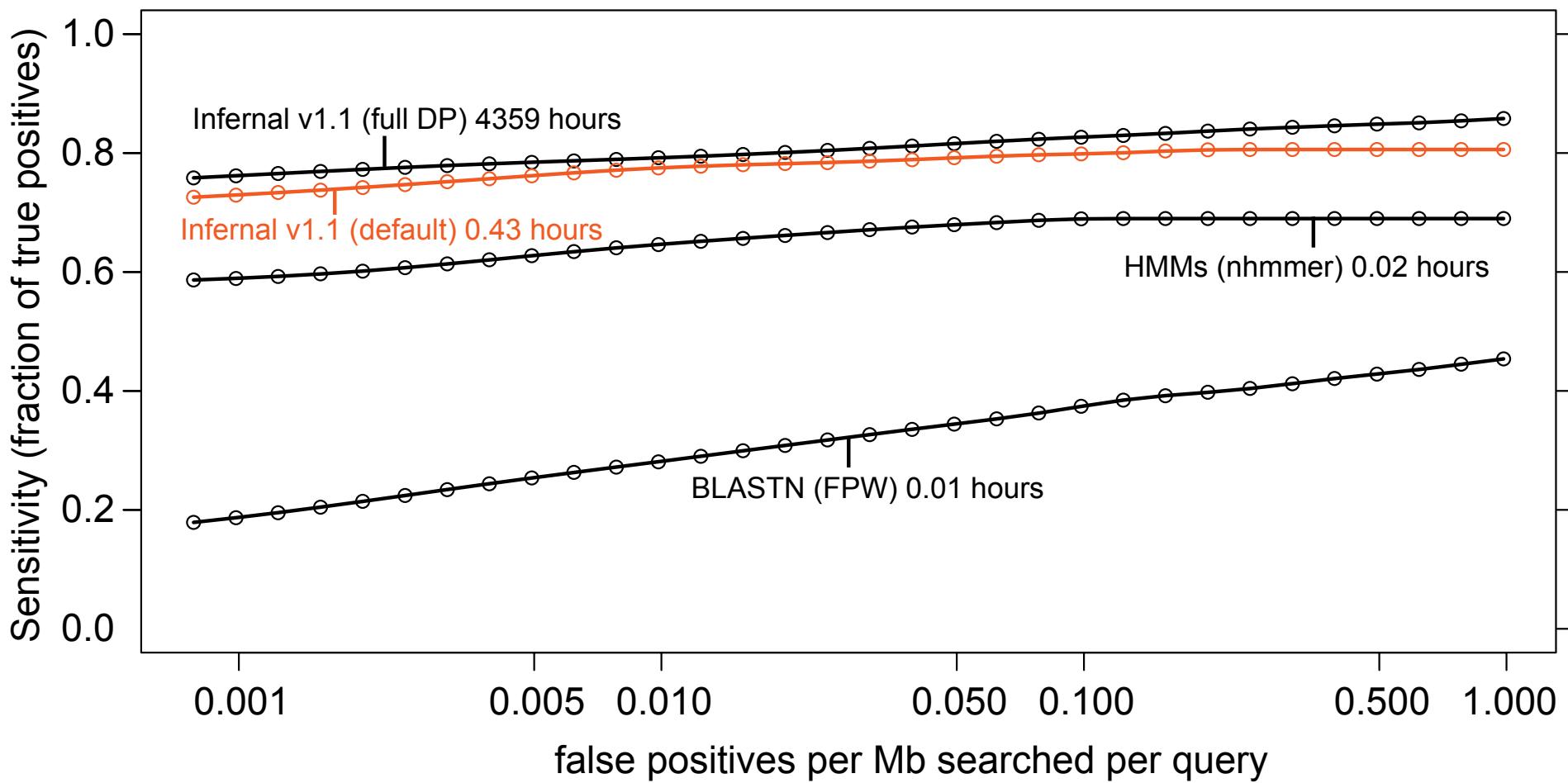


# Use HMMs as filters and to constrain CM alignment

## HMM filter first pass



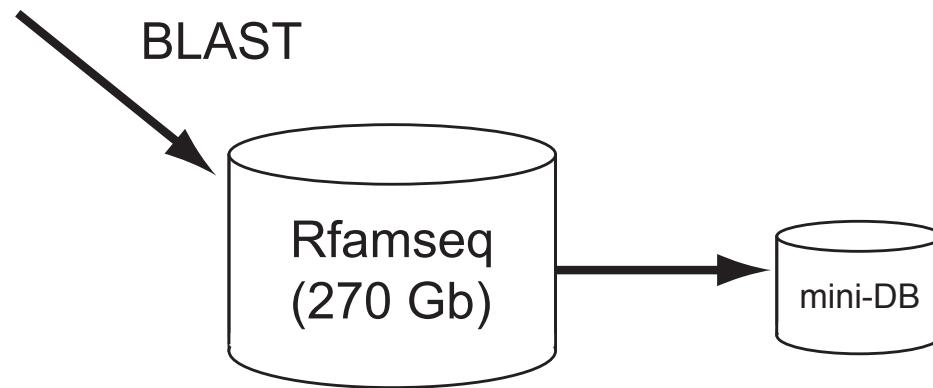
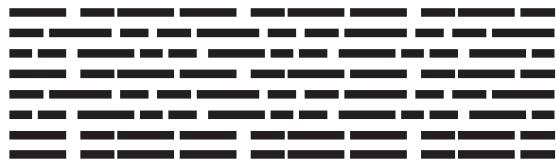
# HMM-based acceleration makes Infernal 10,000 times faster



Nawrocki EP, Eddy SR. Bioinformatics, 29:2487-2489, 2013.

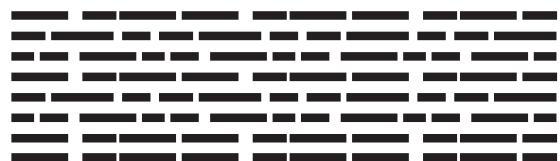
# Rfam used BLAST filters from 2003 to 2012

Rfam seed alignment:

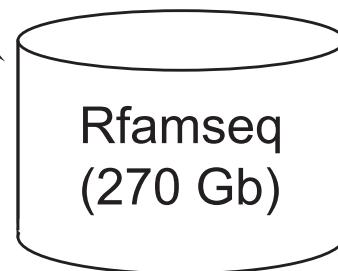


# Rfam used BLAST filters from 2003 to 2012

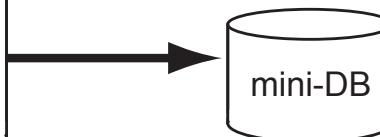
Rfam seed alignment:



BLAST



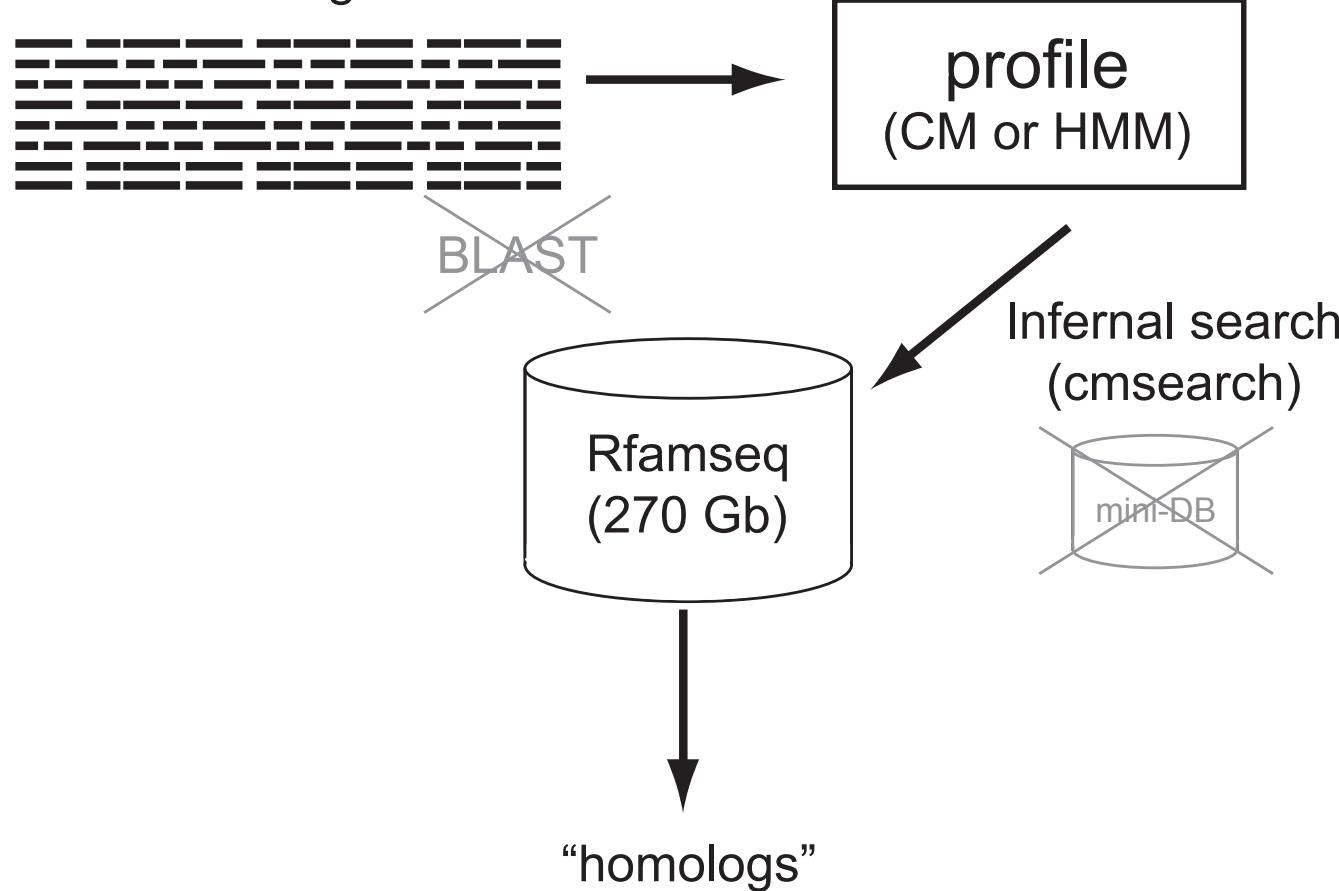
Infernal search  
(cmsearch)



“homologs”

# Rfam 12.0 (2014)\* first release without BLAST filtering

Rfam seed alignment:



# Rfam 12.0 (2014)\* first release without BLAST filtering

Search results against Rfamseq for 200 random families:

strategy	time (h)	# hits	# unique hits
Old (BLAST + Infernal 1.0)	4069.8	179,681	53
New (Infernal 1.1)	4222.2	201,814	22,312

# Acknowledgements

<b>Janelia</b>	<b>EBI (Rfam)</b>
<b>Sean Eddy</b>	<b>Alex Bateman</b>
Elena Rivas	<b>Rob Finn</b>
Travis Wheeler	<b>Sarah Burge</b>
<b>Tom Jones</b>	<b>Evan Floden</b>
Diana Kolbe	John Tate
Seolkyoung Jung	Jen Daub
Rob Finn	
Jody Clements	
Fred Davis	
Lee Henry	
Michael Farrar	

