

Modeling Structural RNA Families with Infernal

Eric Nawrocki

Sean Eddy's Lab



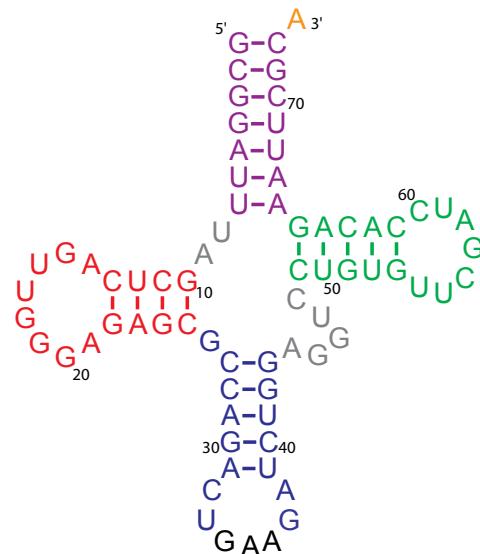
Many functional RNAs adopt a conserved 3-dimensional structure

Three representations of a transfer RNA:

Primary sequence

GC₁GGAUUUAAGCUCAGUUGGG
AGAGC₂GCCAGACU₃GAAGAUC
UGGAGGUCC₄UGUGUUCGAUC
CACAGAAUUCGCA₅

Secondary structure



3-dimensional structure



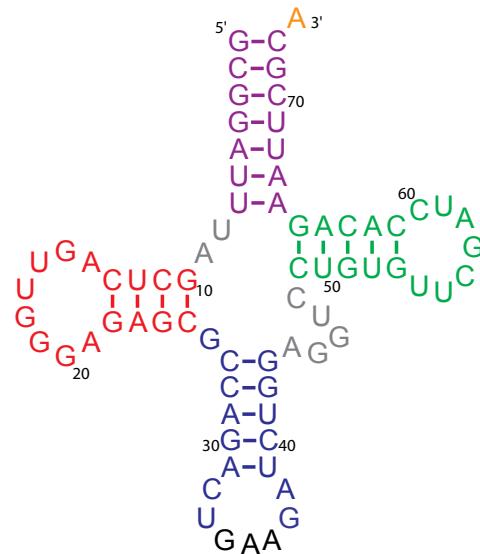
Many functional RNAs adopt a conserved 3-dimensional structure

Three representations of a transfer RNA:

Primary sequence

GC₁GGAUUUAGCUCAGUUGGG
AGAGCGCCAGACUGAAGAUC
UGGAGGUC₂CUGUGUUCGAUC
CACAGAAUUCGCA

Secondary structure



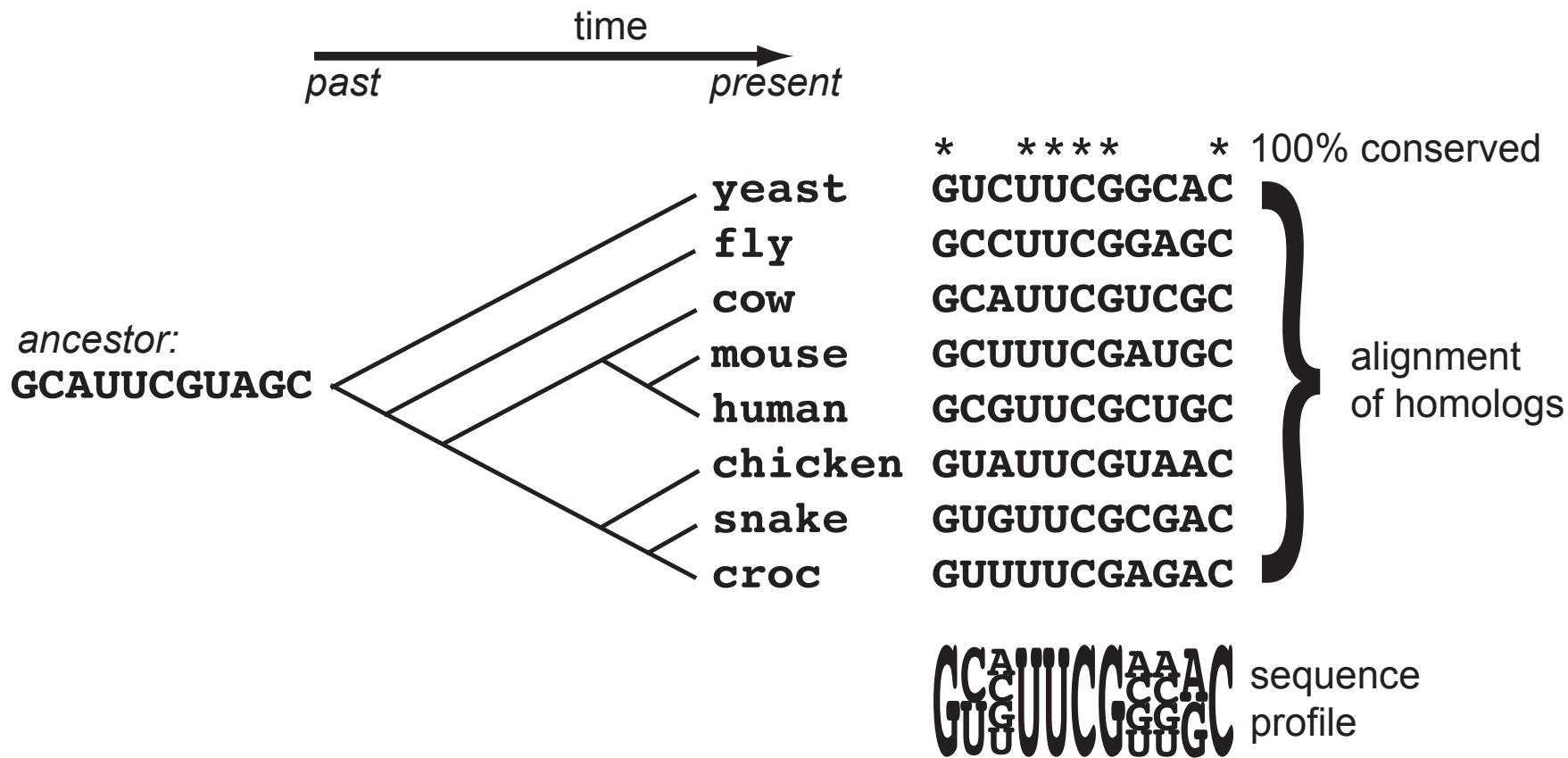
3-dimensional structure



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

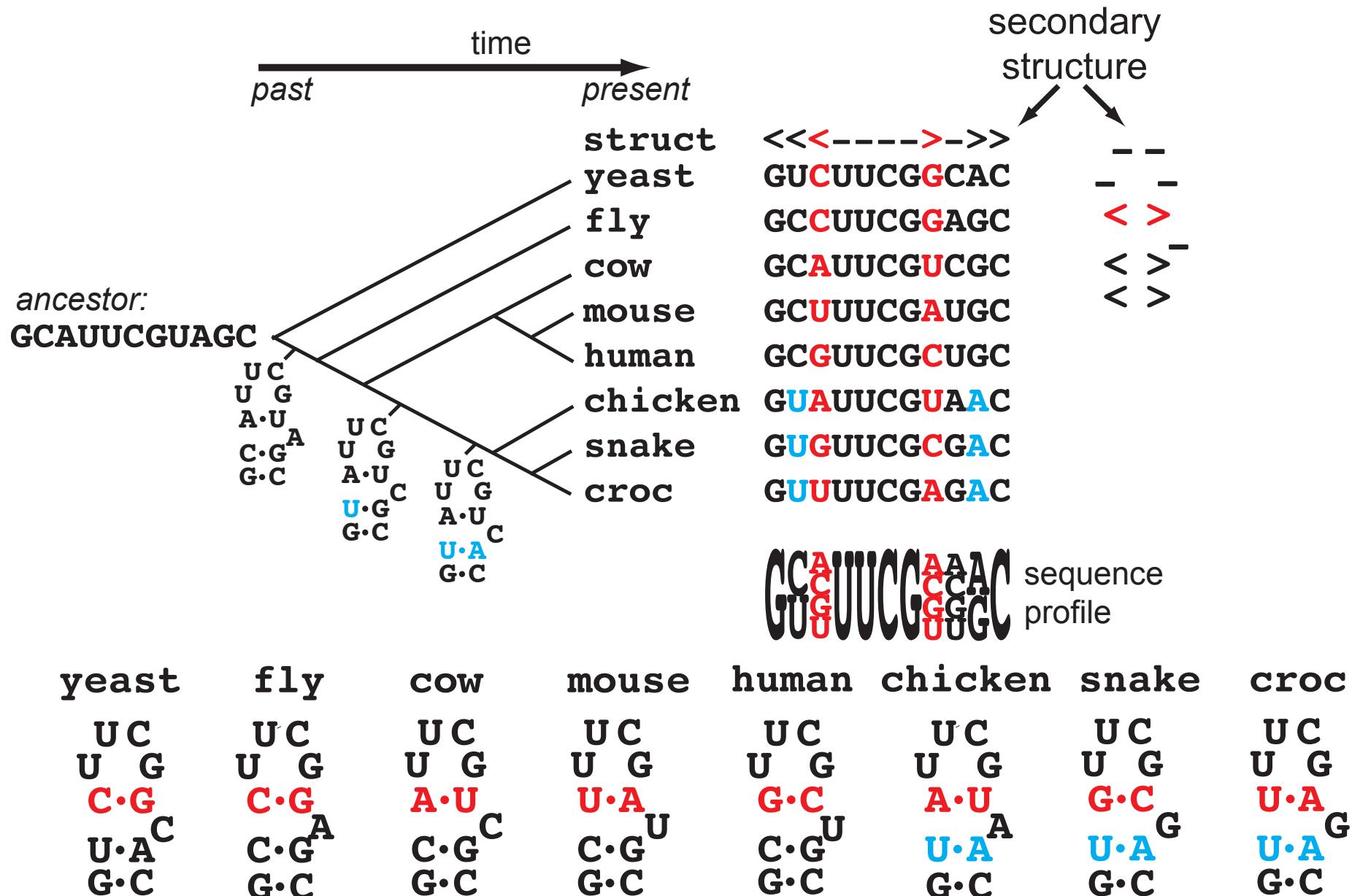
Sequence conservation provides information for homology searches

Conservation levels vary across alignment columns.

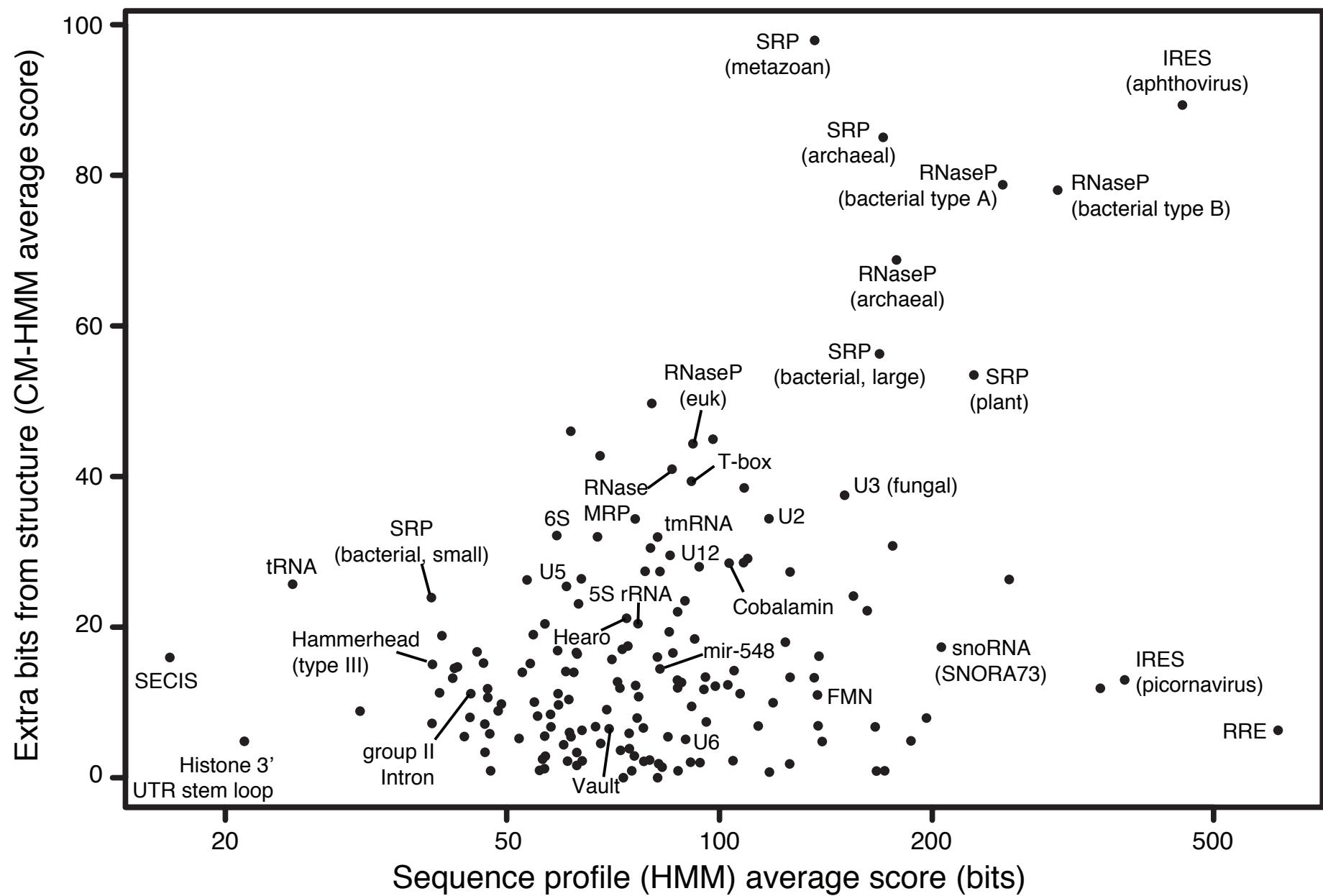


Structure conservation provides additional information

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



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.

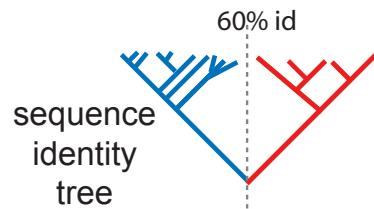
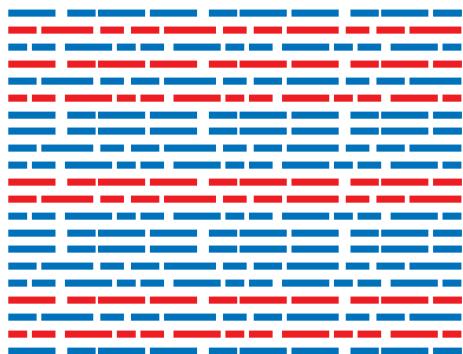


<http://infernal.janelia.org>
Nawrocki EP, Eddy SR
Bioinformatics,
29:2933-2935, 2013.
Eddy SR, Durbin R.
Nucleic Acids Research,
22:2079-2088, 1994.

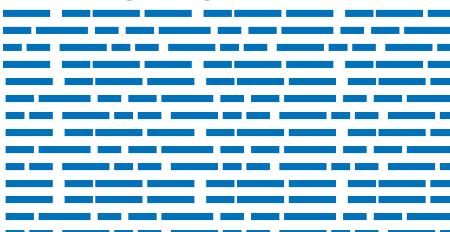
Is the added complexity worth it?

RMARK: a challenging internal RNA homology search benchmark

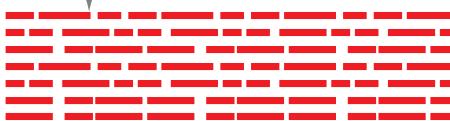
Rfam seed alignment:



training alignment

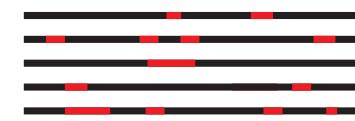


no train/test sequence pair is > 60% identical



test sequences

embed in
pseudo-genome

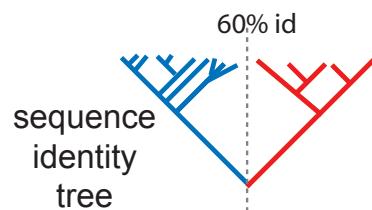
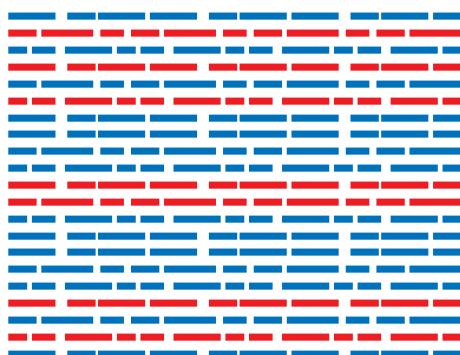


10 1Mb sequences
with 780 embedded
test seqs from 106 families

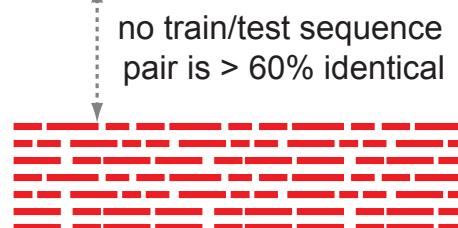
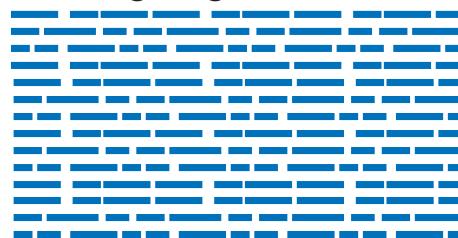
Is the added complexity worth it?

RMARK: a challenging internal RNA homology search benchmark

Rfam seed alignment:



training alignment



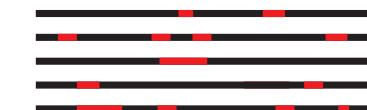
test sequences

profile
(CM or HMM)

BLAST

search

embed in
pseudo-genome



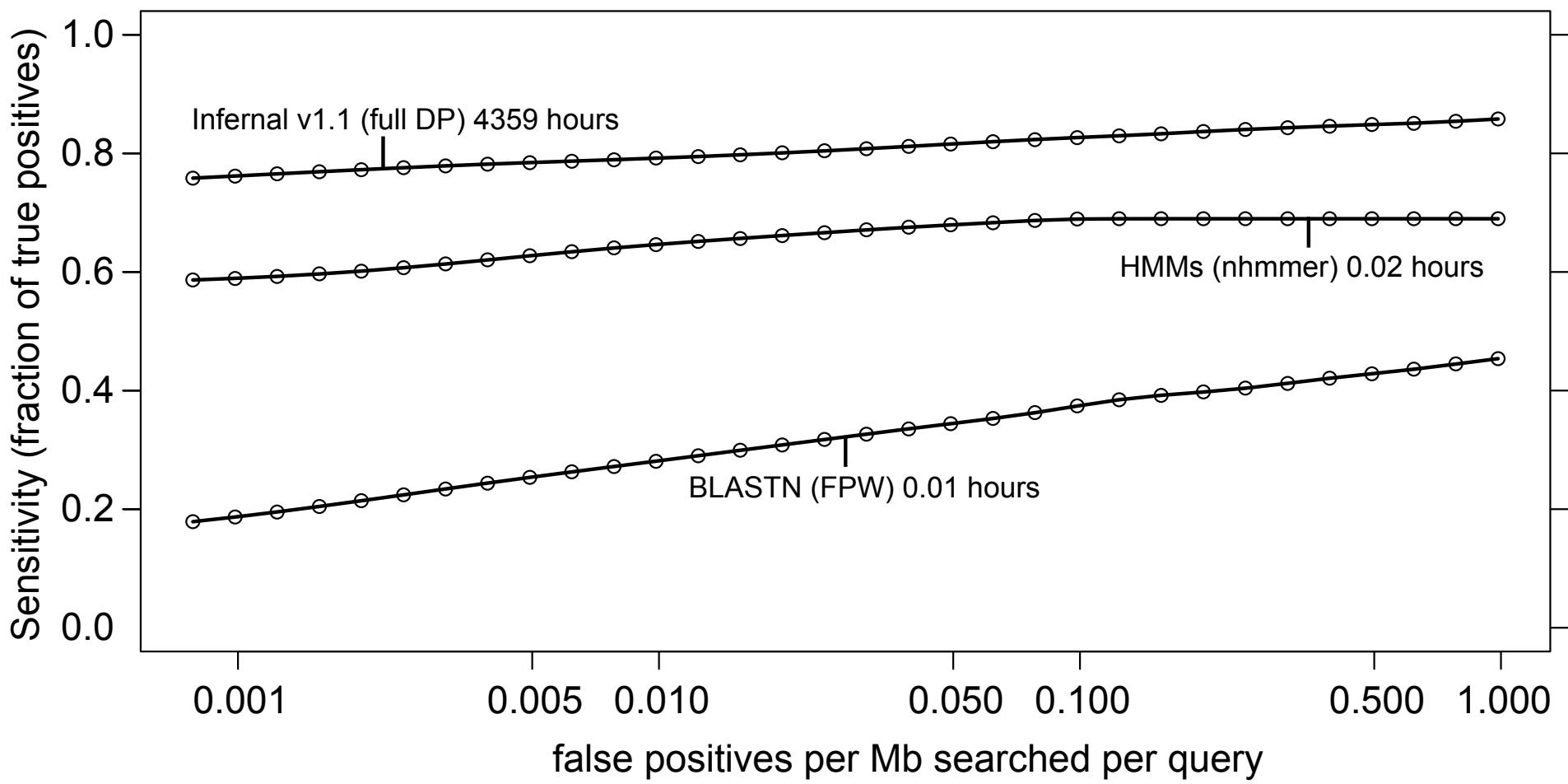
10 1Mb sequences
with 780 embedded
test seqs from 106 families

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 -
...

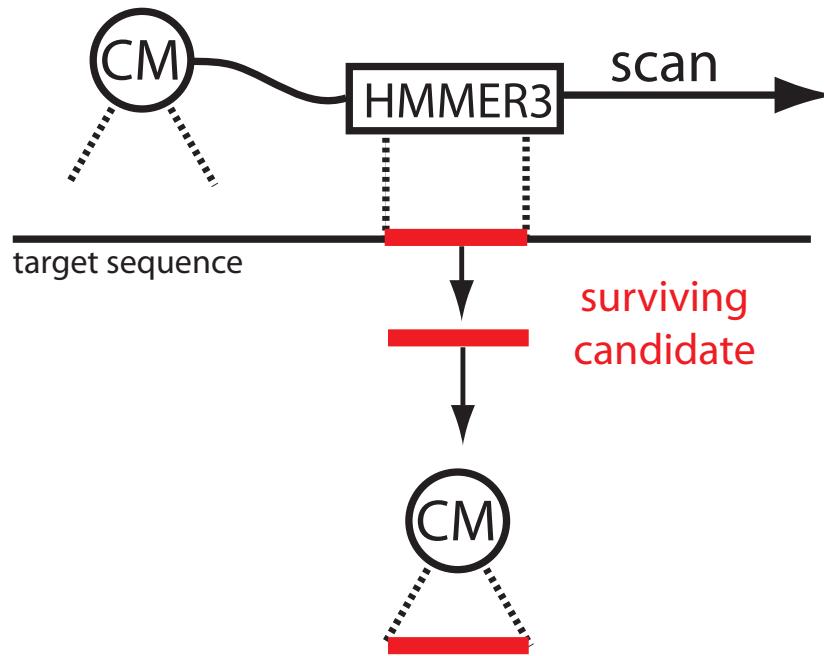
Infernal outperforms primary-sequence based methods on our benchmark (and others*, not shown)



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

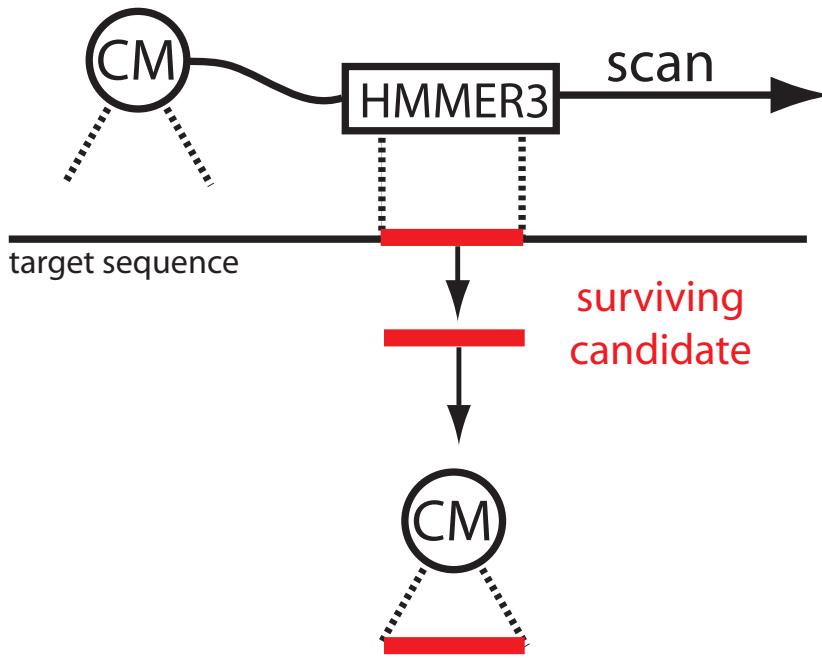
Filter target database using profile HMMs*

HMM filter first pass



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.

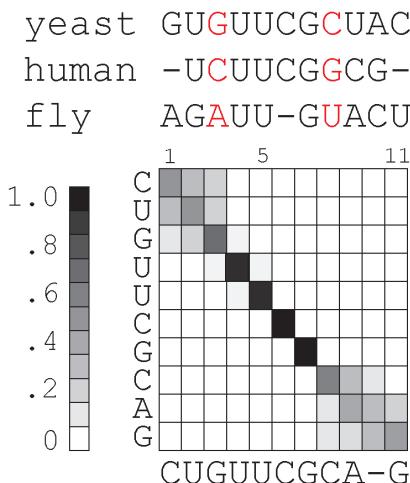
Accelerating CM alignment step 1: align sequence with HMM

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

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

HMMs -

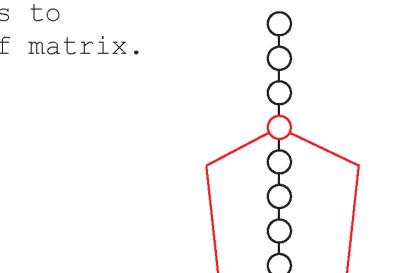
Each column of seed alignment corresponds to a column of matrix.



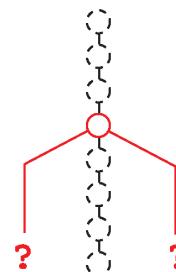
CMs -

Each column of seed alignment corresponds to a state.

yeast	human	fly
U C	U C	U
U G	U G	U G
G•C	C•G	A•U
U•A U	U•G C	G•C A
G•C		A•U



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



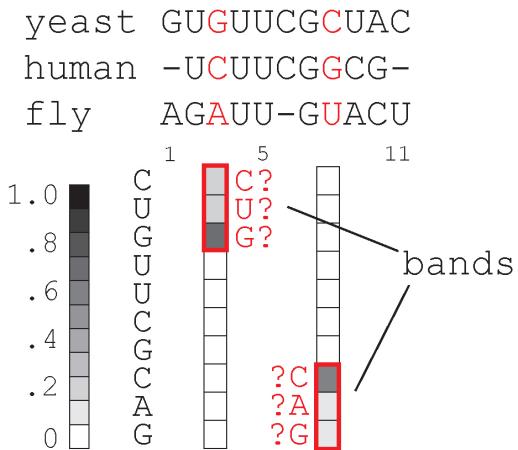
CUGUUCGCAG

45 possibilities

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

HMMs -

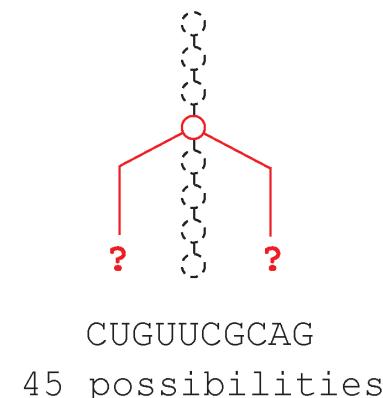
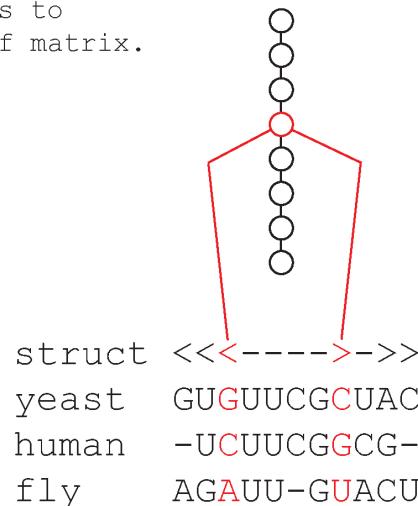
Each column of seed alignment corresponds to a column of matrix.



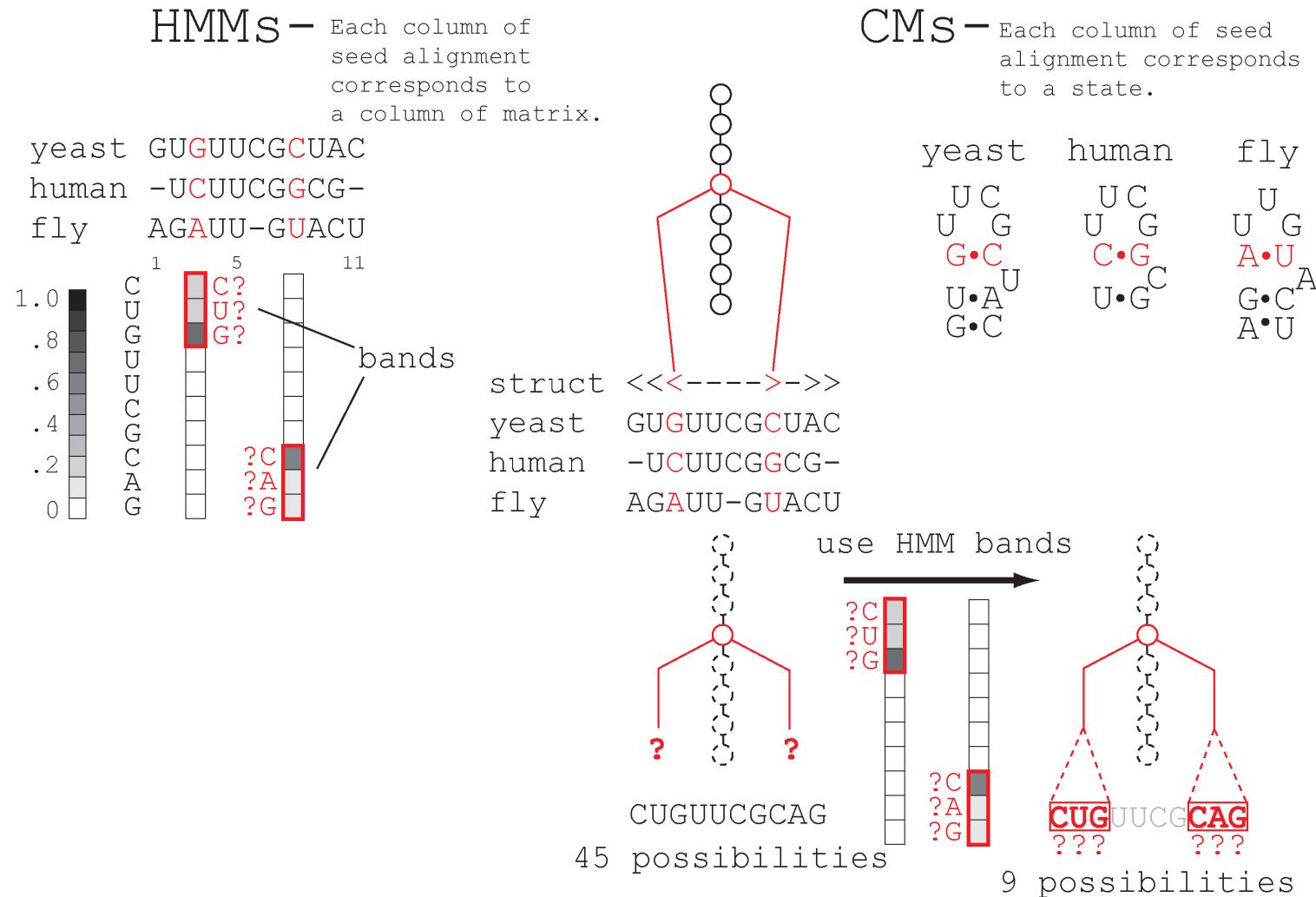
CMs -

Each column of seed alignment corresponds to a state.

yeast	human	fly
U C	U C	U
U G	U G	U G
G•C	C•G	A•U
U•A U	U•G C	G•C A
G•C		A•U

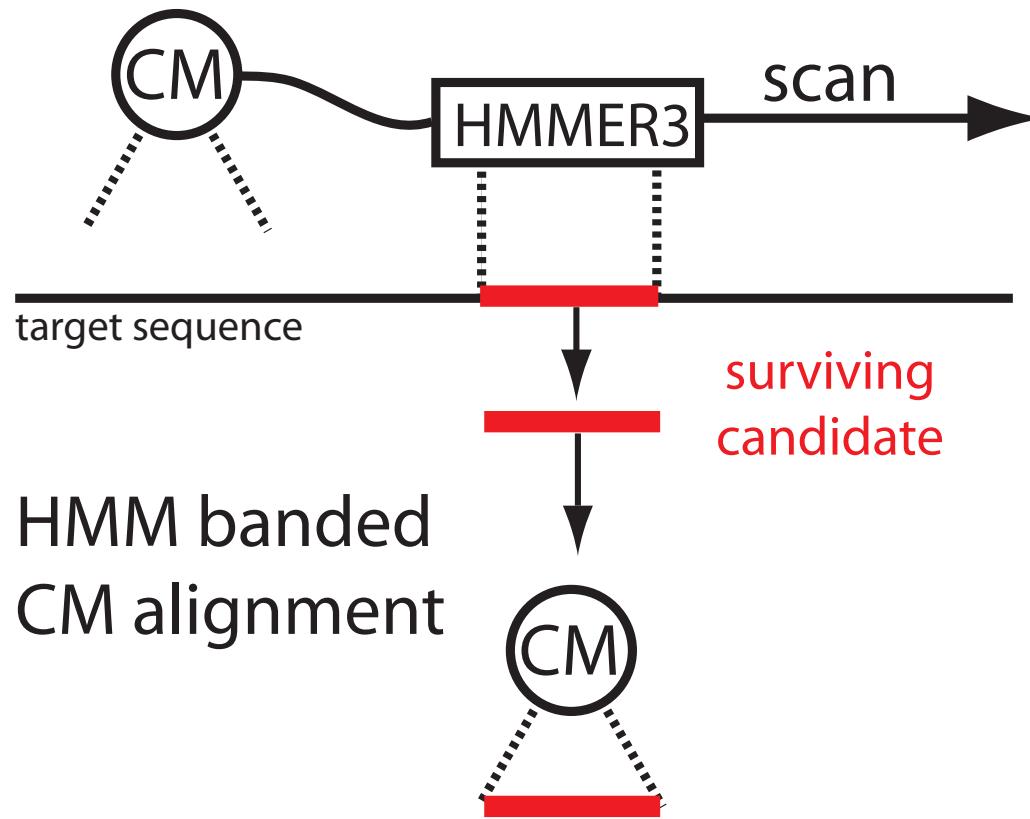


Accelerating CM alignment step 3: use HMM alignment confidence to constrain CM alignment*

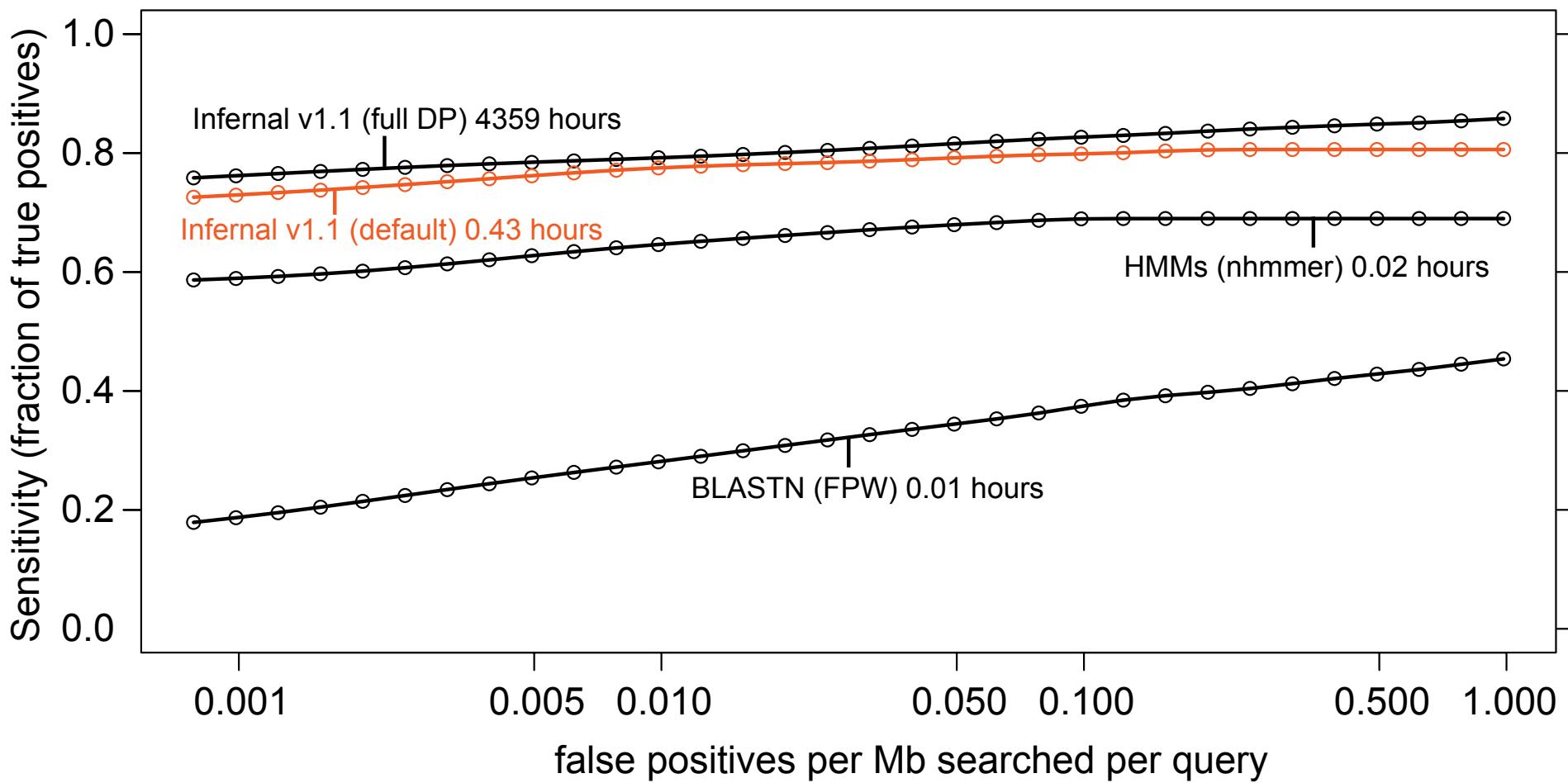


Use HMMs as filters and to constrain CM alignment

HMM filter first pass



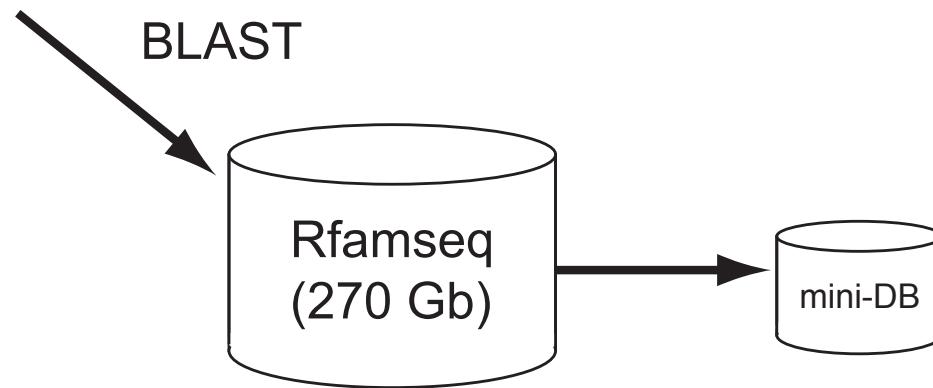
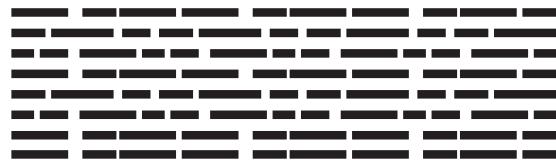
HMM-based acceleration makes Infernal 10,000 times faster



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

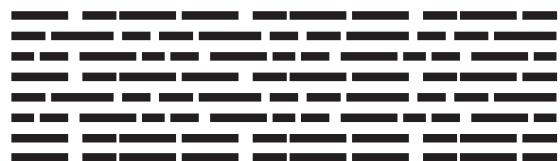
Rfam used BLAST filters from 2003 to 2012

Rfam seed alignment:

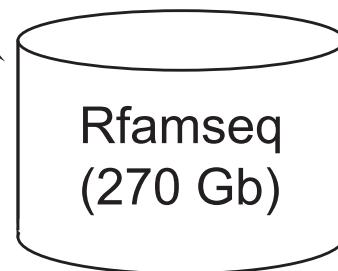


Rfam used BLAST filters from 2003 to 2012

Rfam seed alignment:



BLAST



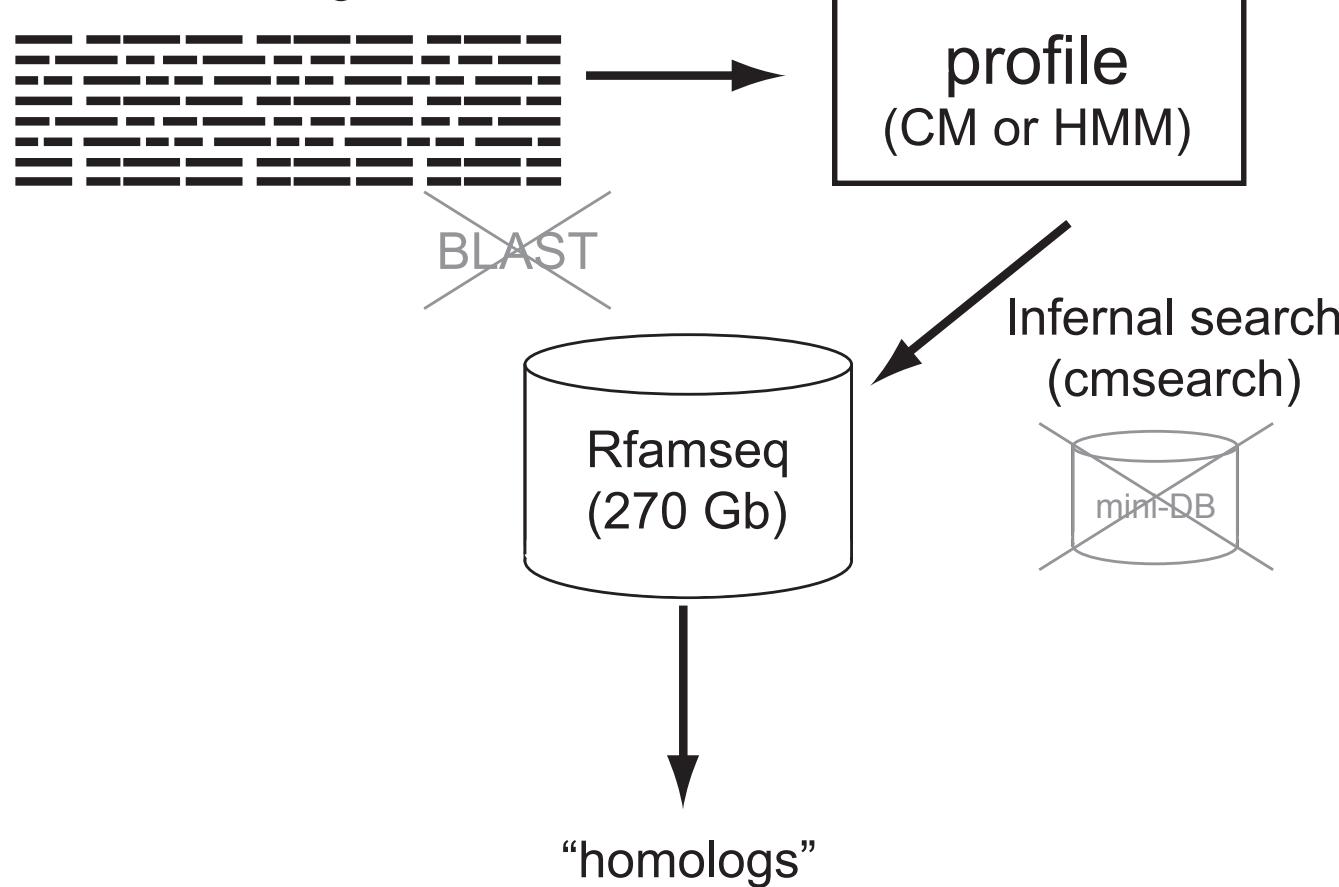
Infernal search
(cmsearch)



“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

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