

Structural RNA and viral sequence analysis

Eric Nawrocki

Intramural Research Program
National Library of Medicine
National Institutes of Health



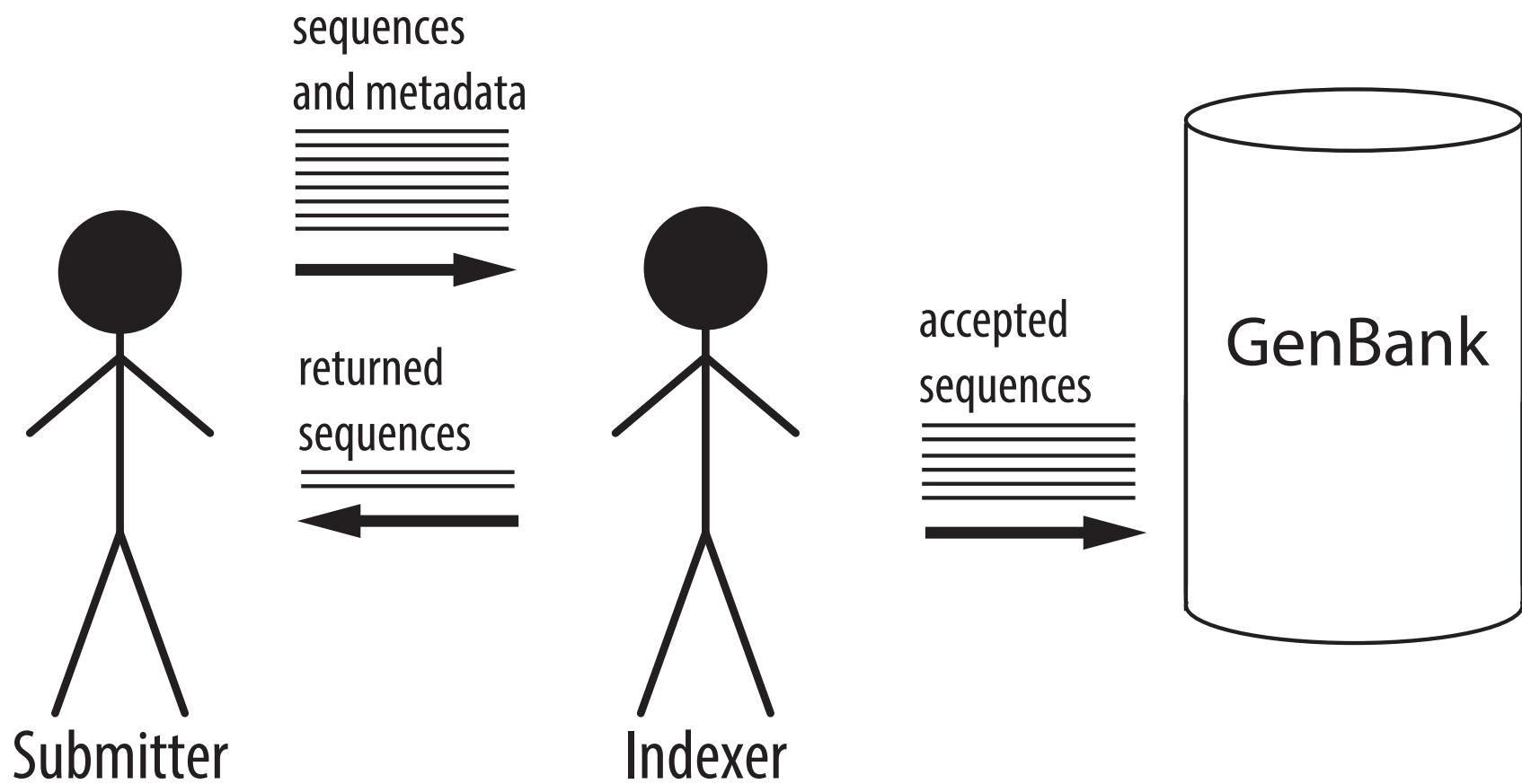
Two main areas of my research:

- 1. Viral sequence analysis tools, since 2015**

- 2. Structural RNA analysis tools, since 2004**

add logos?

GenBank indexers handle incoming sequence submissions



SOFTWARE

Open Access

VADR: validation and annotation of virus sequence submissions to GenBank

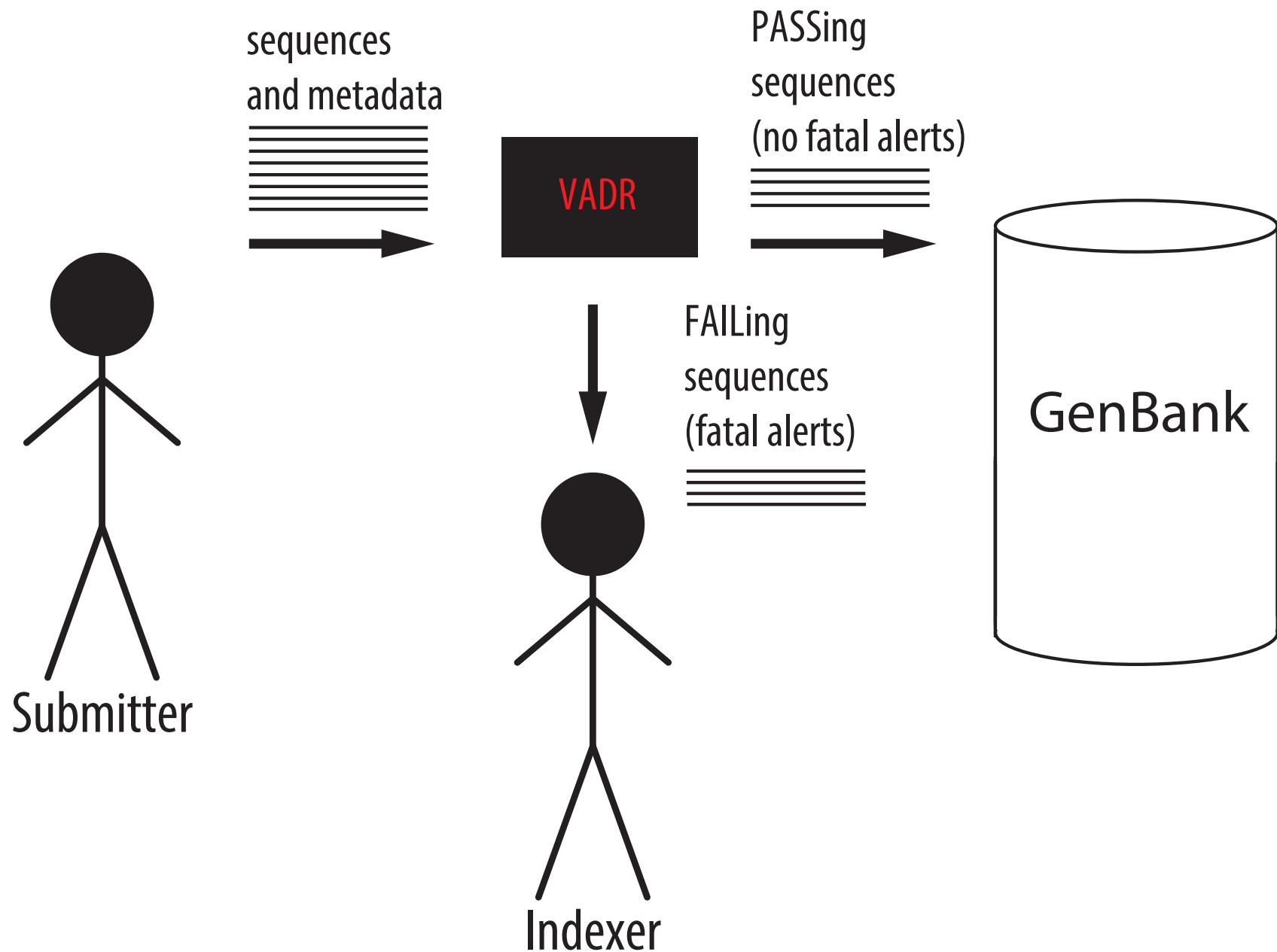


Alejandro A. Schäffer^{1,2}, Eneida L. Hatcher², Linda Yankie², Lara Shonkwiler^{2,3}, J. Rodney Brister², Ilene Karsch-Mizrachi² and Eric P. Nawrocki^{2*} 

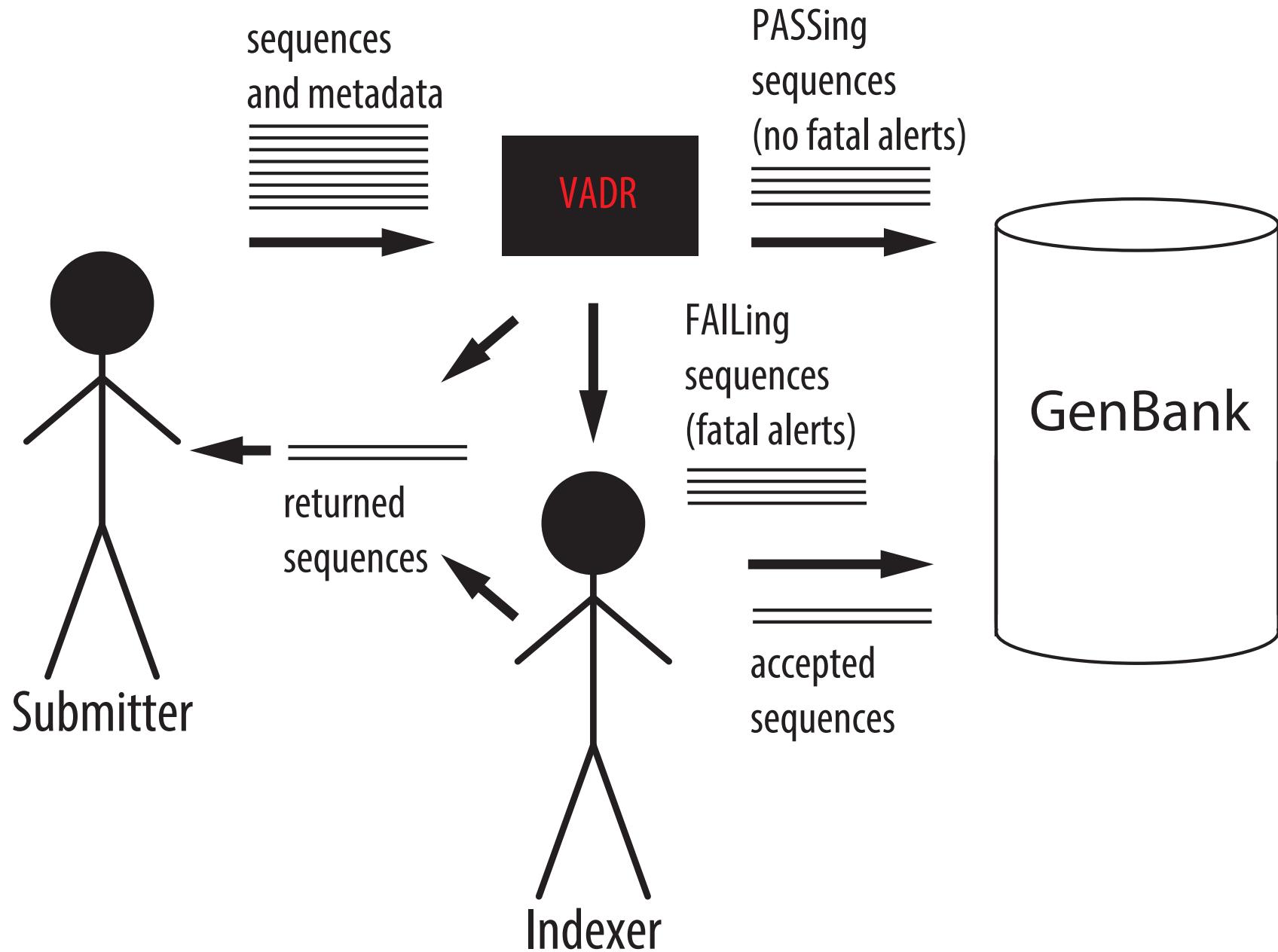
- general tool for reference-based annotation of viral sequences
- used for Norovirus and Dengue virus submissions since 2018
- used for SARS-CoV-2 submissions since March 2020
- currently also used manually for RSV, MpoX, and some Influenza submissions

VADR assists GenBank indexers:

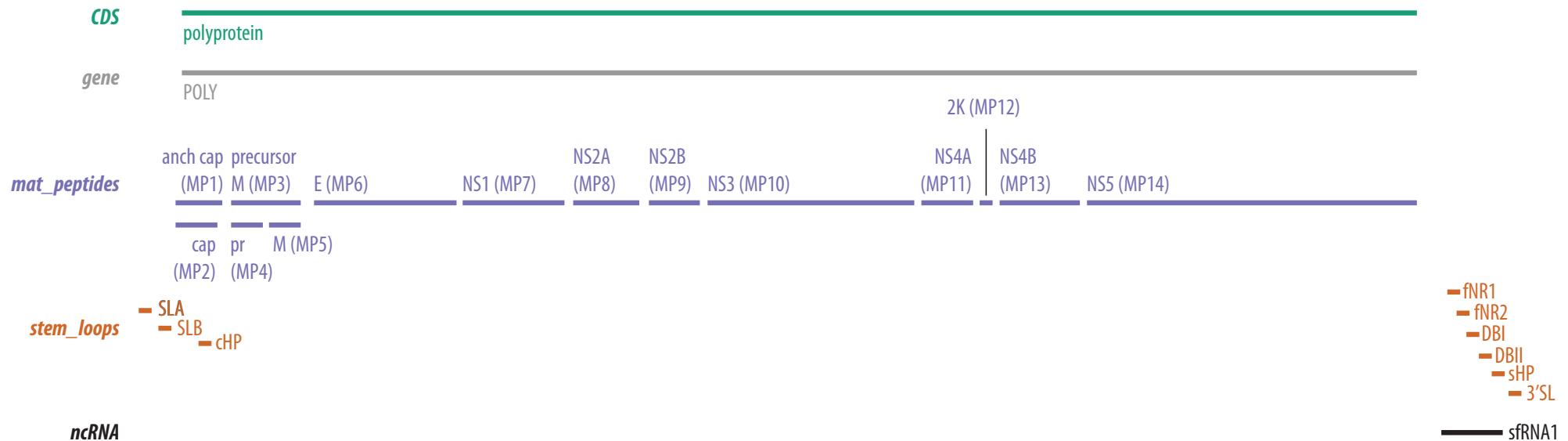
Each sequence **PASSes** or **FAILs**



**Indexers decide fate of some FAILing sequences
but some are sent directly back to submitter with error reports**



VADR builds a reference model of a RefSeq and its features



NC_001477 MODEL



Group: Dengue; Subgroup: 1

VADR validates and annotates each input sequence using its best-matching model

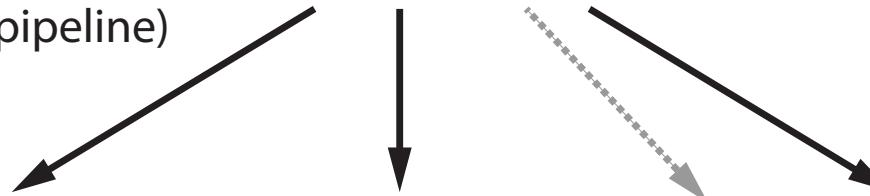
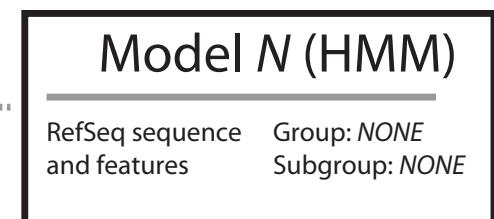
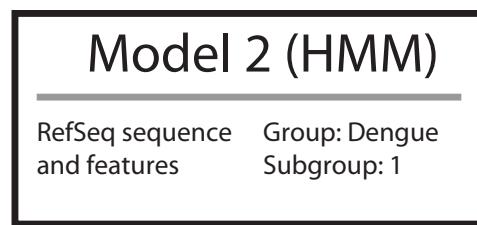
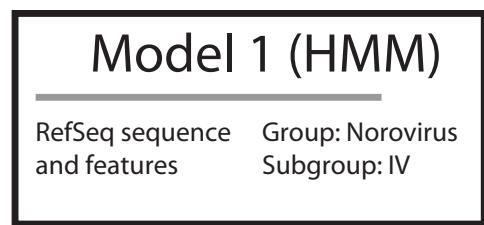
- Each sequence S proceeds through 4 stages:
 1. **Classification**
 2. **Coverage determination**
 3. **Alignment**
 4. **Protein validation**

Different types of alerts are identified and reported at each stage

Stage 1: Classification

Score each sequence
with all models
(HMMER3 shortened pipeline)

input sequences:



low HMM score

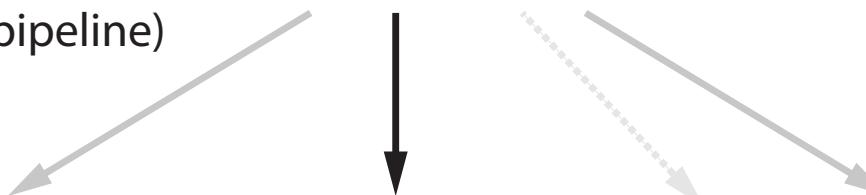
highest HMM score

low HMM score

Stage 1: Classification

Score each sequence
with all models
(HMMER3 shortened pipeline)

input sequences:



Model 1 (HMM)

RefSeq sequence
and features Group: Norovirus
Subgroup: IV

Model 2 (HMM)

RefSeq sequence
and features Group: Dengue
Subgroup: 1

Model N (HMM)

RefSeq sequence
and features Group: NONE
Subgroup: NONE

low HMM score

highest HMM score

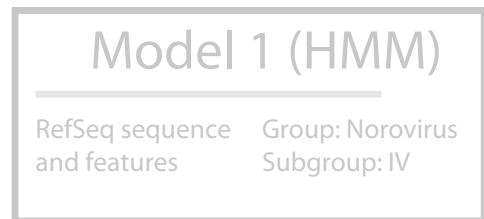
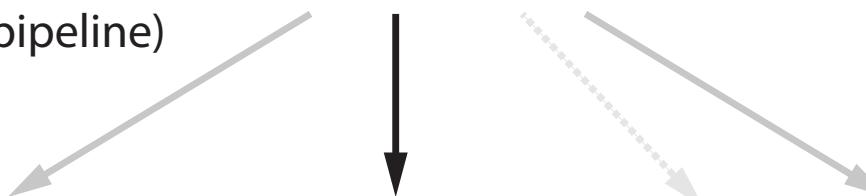
low HMM score

***best-matching model
used in remaining stages***

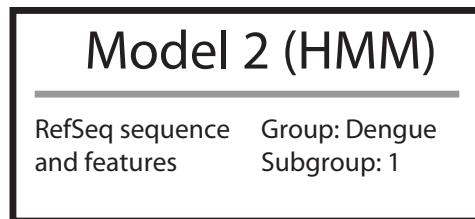
Stage 1: Classification

Score each sequence
with all models
(HMMER3 shortened pipeline)

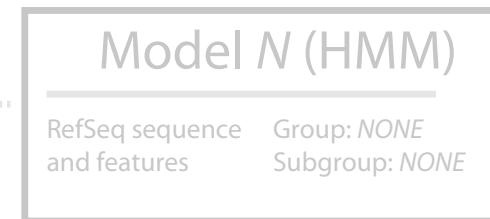
input sequences:



low HMM score



highest HMM score



low HMM score

***best-matching model
used in remaining stages***

code	S/F	error message	description
Fatal alerts detected in the classification stage			
noannotn*	S	NO_ANNOTATION	no significant similarity detected
revcompl*	S	REVCOMPLEM	sequence appears to be reverse complemented
incsbgrp	S	INCORRECT_SPECIFIED_SUBGROUP	score difference too large between best overall model and best specified subgroup model
incgroup	S	INCORRECT_SPECIFIED_GROUP	score difference too large between best overall model and best specified group model
Non-fatal alerts detected in the classification stage			
qstsbgp	S	QUESTIONABLE_SPECIFIED_SUBGROUP	best overall model is not from specified subgroup
qstgroup	S	QUESTIONABLE_SPECIFIED_GROUP	best overall model is not from specified group
indfclas	S	INDEFINITE_CLASSIFICATION	low score difference between best overall model and second best model (not in best model's subgroup)
lowscore	S	LOW_SCORE	score to homology model below low threshold

Stage 2: Coverage determination

Search each sequence with best-matching model (HMMER3 full pipeline)

input sequences that match best to NC_001477:

S1 _____
S2 _____
S3 _____
S4 _____



NC_001477 MODEL



Group: Dengue; Subgroup: 1



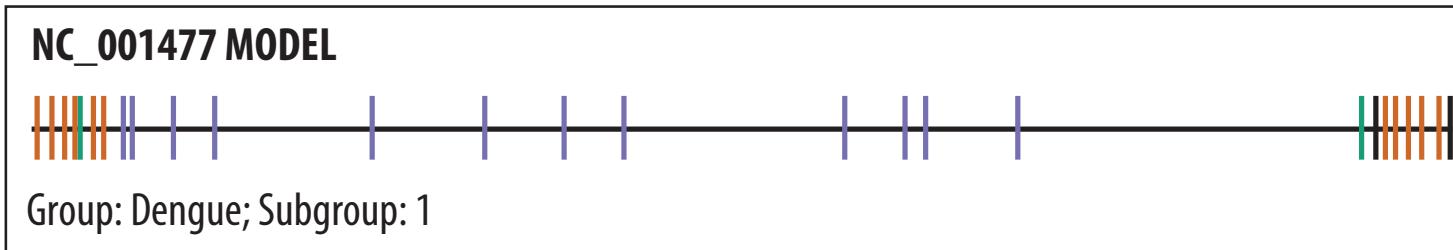
NC_001477 full length sequence
S1 (expected)
NC_001477 partial or truncated sequence
S2 (expected)

Stage 2: Coverage determination

Search each sequence with best-matching model (HMMER3 full pipeline)

input sequences that match best to NC_001477:

S1 _____
S2 _____
S3 _____
S4 _____

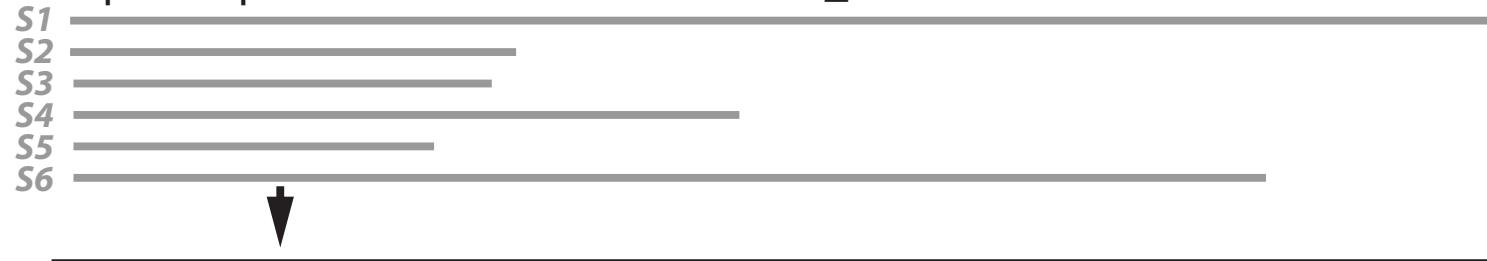


code	S/F	error message	description
Fatal alerts detected in the coverage stage			
lowcovrg	S	LOW_COVERAGE	low sequence fraction with significant similarity to homology model
dupregin	S	DUPLICATE_REGIONS	similarity to a model region occurs more than once
discontn	S	DISCONTINUOUS_SIMILARITY	not all hits are in the same order in the sequence and the homology model
indfstrn	S	INDEFINITE_STRAND	significant similarity detected on both strands
lowsim5s	S	LOW_SIMILARITY_START	significant similarity not detected at 5' end of the sequence
lowsim3s	S	LOW_SIMILARITY_END	significant similarity not detected at 3' end of the sequence
lowsimis	S	LOW_SIMILARITY	internal region without significant similarity
Non-fatal alerts detected in the coverage stage			
biasdseq	S	BIASED_SEQUENCE	high fraction of score attributed to biased sequence composition

Stage 3: Alignment and feature mapping

Align each sequence to its best-matching model (Infernal's cmalign)

input sequences that match best to NC_001477:



NC_001477 MODEL



Group: Dengue; Subgroup: 1



Stage 3: Alignment and feature mapping

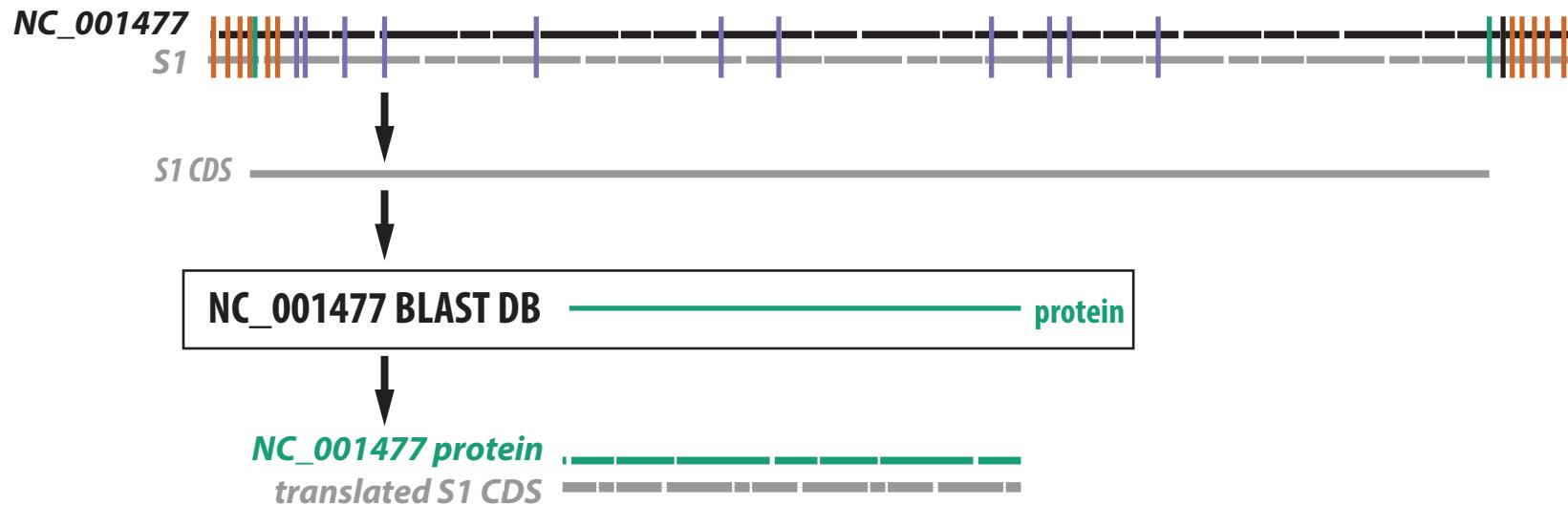
Align each sequence to its best-matching model (Infernal's cmalign)



code	S/F	error message	description
Fatal alerts detected in the annotation stage			
unexdivg*	S	UNEXPECTED_DIVERGENCE	sequence is too divergent to confidently assign nucleotide-based annotation
noftrann*	S	NO_FEATURES_ANNOTATED	sequence similarity to homology model does not overlap with any features
mutstart	F	MUTATION_AT_START	expected start codon could not be identified
mutendcd	F	MUTATION_AT_END	expected stop codon could not be identified, predicted CDS stop by homology is invalid
mutendns	F	MUTATION_AT_END	expected stop codon could not be identified, no in-frame stop codon exists 3' of predicted valid start codon
mutendex	F	MUTATION_AT_END	expected stop codon could not be identified, first in-frame stop codon exists 3' of predicted stop position
unexleng	F	UNEXPECTED_LENGTH	length of complete coding (CDS or mat_peptide) feature is not a multiple of 3
cdsstoppn	F	CDS_HAS_STOP_CODON	in-frame stop codon exists 5' of stop position predicted by homology to reference
peptrans	F	PEPTIDE_TRANSLATION_PROBLEM	mat_peptide may not be translated because its parent CDS has a problem
pepadjcy	F	PEPTIDE_ADJACENCY_PROBLEM	predictions of two mat_peptides expected to be adjacent are not adjacent
indfantn	F	INDEFINITE_ANNOTATION	nucleotide-based search identifies CDS not identified in protein-based search
indf5gap	F	INDEFINITE_ANNOTATION_START	alignment to homology model is a gap at 5' boundary
indf5loc	F	INDEFINITE_ANNOTATION_START	alignment to homology model has low confidence at 5' boundary
indf3gap	F	INDEFINITE_ANNOTATION_END	alignment to homology model is a gap at 3' boundary
indf3loc	F	INDEFINITE_ANNOTATION_END	alignment to homology model has low confidence at 3' boundary
lowsim5f	F	LOW FEATURE SIMILARITY_START	region within annotated feature at 5' end of sequence lacks significant similarity
lowsim3f	F	LOW FEATURE SIMILARITY_END	region within annotated feature at 3' end of sequence lacks significant similarity
lowsimif	F	LOW FEATURE SIMILARITY	region within annotated feature lacks significant similarity

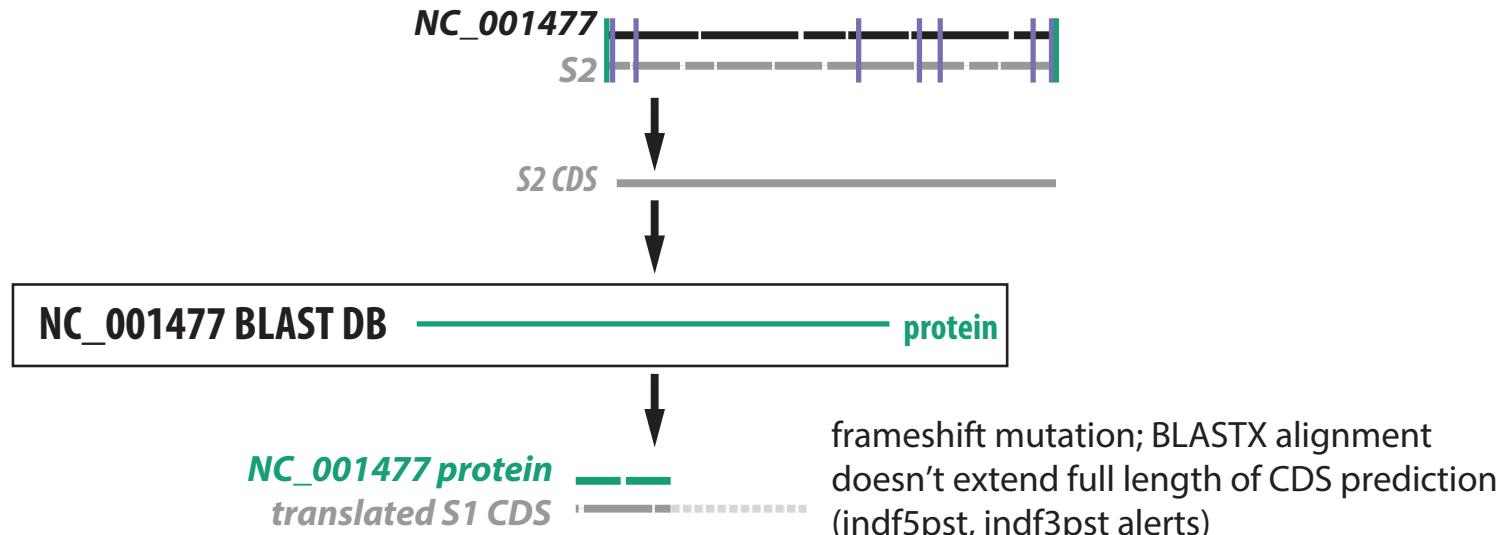
Stage 4: Protein validation (Alejandro Schäffer)

Compare each predicted CDS to model (RefSeq) proteins with BLASTX



Stage 4: Protein validation (Alejandro Schäffer)

Compare each predicted CDS to model (RefSeq) proteins with BLASTX



code	S/F	error message	description
Fatal alerts detected in the protein validation stage			
cdsstopp	F	CDS_HAS_STOP_CODON	stop codon in protein-based alignment
indf5antp	F	INDEFINITE_ANNOTATION	protein-based search identifies CDS not identified in nucleotide-based search
indf5plg	F	INDEFINITE_ANNOTATION_START	protein-based alignment extends past nucleotide-based alignment at 5' end
indf5pst	F	INDEFINITE_ANNOTATION_START	protein-based alignment does not extend close enough to nucleotide-based alignment 5' endpoint
indf3plg	F	INDEFINITE_ANNOTATION_END	protein-based alignment extends past nucleotide-based alignment at 3' end
indf3pst	F	INDEFINITE_ANNOTATION_END	protein-based alignment does not extend close enough to nucleotide-based alignment 3' endpoint
indfstrp	F	INDEFINITE_STRAND	strand mismatch between protein-based and nucleotide-based predictions
insertnp	F	INSERTION_OF_NT	too large of an insertion in protein-based alignment
deletinp	F	DELETION_OF_NT	too large of a deletion in protein-based alignment

VADR used for Norovirus and Dengue virus sequences since 2018

	Norovirus	Dengue virus
length	7.6Kb	10.7Kb
# seqs	44,936	113,211
% seqs full length	5.1%	8.4%
% Ns	0.5%	0.2%
% seqs with stretch of \geq 50 Ns	1.0%	0.4%
average % identity	81.6%	94.4%

VADR v1.0 performance

seconds per sequence	42.4	92.6
required RAM	8Gb	8Gb
total running time, CPU days	1.1	10.2

SARS-CoV-2 sequence submissions have increased since early 2020

month	year	#new seqs	#cumulative seqs
Jan	2020	32	32
Feb	2020	58	90
Mar	2020	332	422
Apr	2020	1541	1963
May	2020	2974	4937
Jun	2020	3394	8331
Jul	2020	3604	11,935
Aug	2020	3818	15,753
Sep	2020	6731	22,484
Oct	2020	11,939	34,423
Nov	2020	4274	38,697
Dec	2020	4530	43,227
Jan	2021	8775	52,002
Feb	2021	26,078	78,080
Mar	2021	42,607	120,687
Apr	2021	97,095	217,782
May	2021	104,729	322,511
Jun	2021	46,187	368,698
Jul	2021	43,336	412,034
Aug	2021	141,958	553,992
Sep	2021	267,562	821,554
Oct	2021	239,296	1,060,850
Nov	2021	267,270	1,328,120
Dec	2021	288,771	1,616,891
Jan	2022	258,522	1,875,413
Feb	2022	230,185	2,105,598

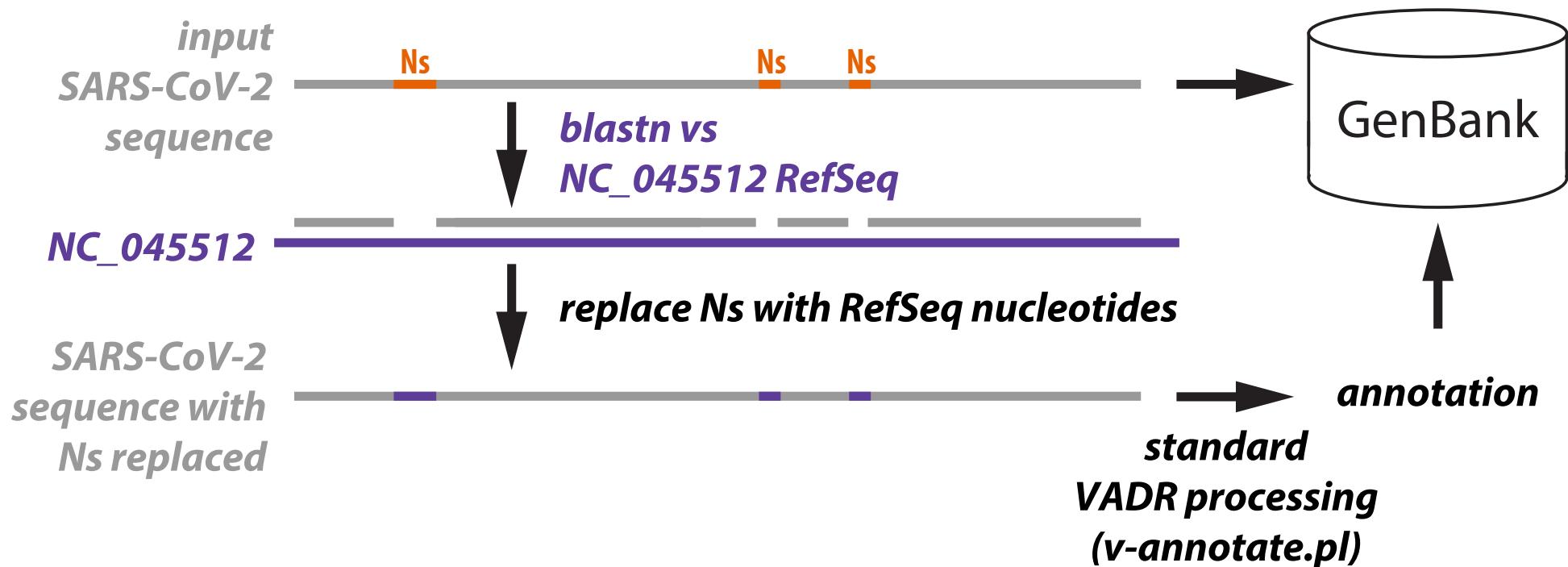
SARS-CoV-2 sequences differ from Norovirus and Dengue virus in several ways that impact VADR processing

	Norovirus	Dengue virus	SARS-CoV-2
length	7.6Kb	10.7Kb	29.9Kb
# seqs	44,936	113,211	1,616,891
% seqs full length	5.1%	8.4%	99.7%
% Ns	0.5%	0.2%	1.4%
% seqs with stretch of \geq 50 Ns	1.0%	0.4%	38.7%
average % identity	81.6%	94.4%	99.4%

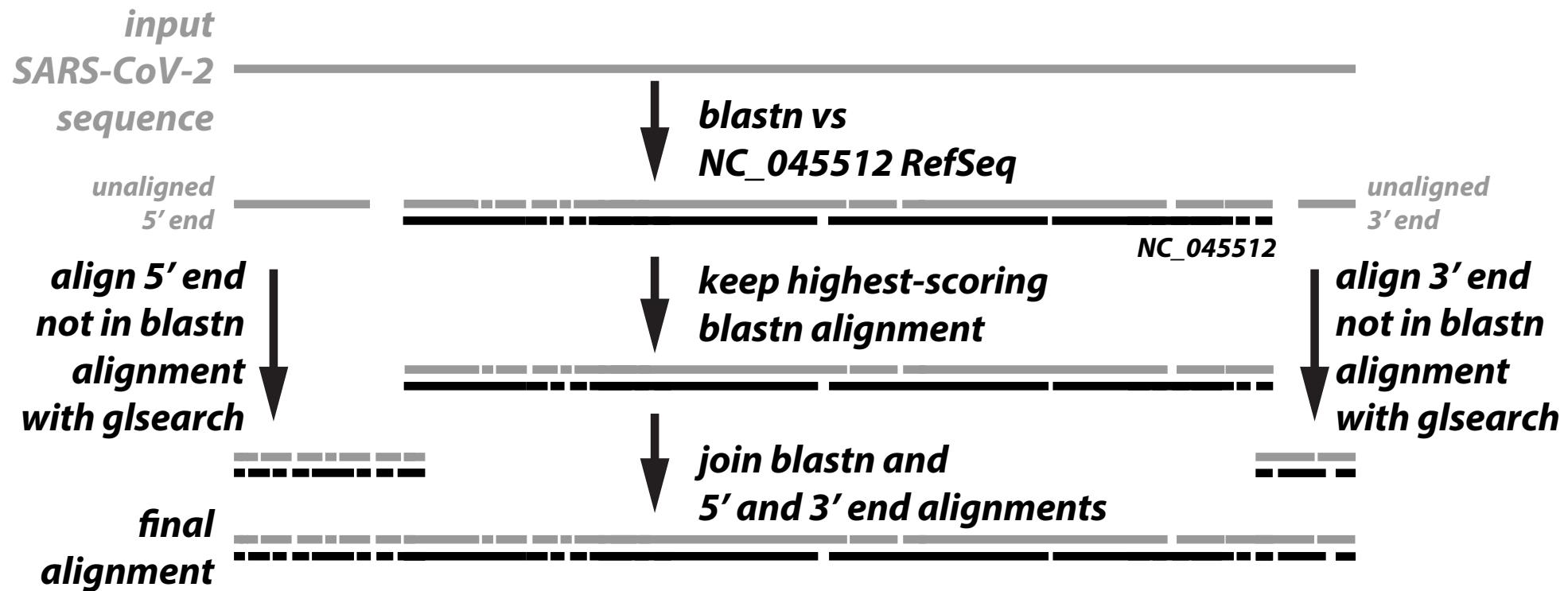
VADR v1.0 performance

seconds per sequence	42.4	92.6	331.8
required RAM	8Gb	8Gb	64Gb
total running time, CPU days	1.1	10.2	6187.6

Replacing Ns with expected nucleotides allows many 'good' sequences to pass

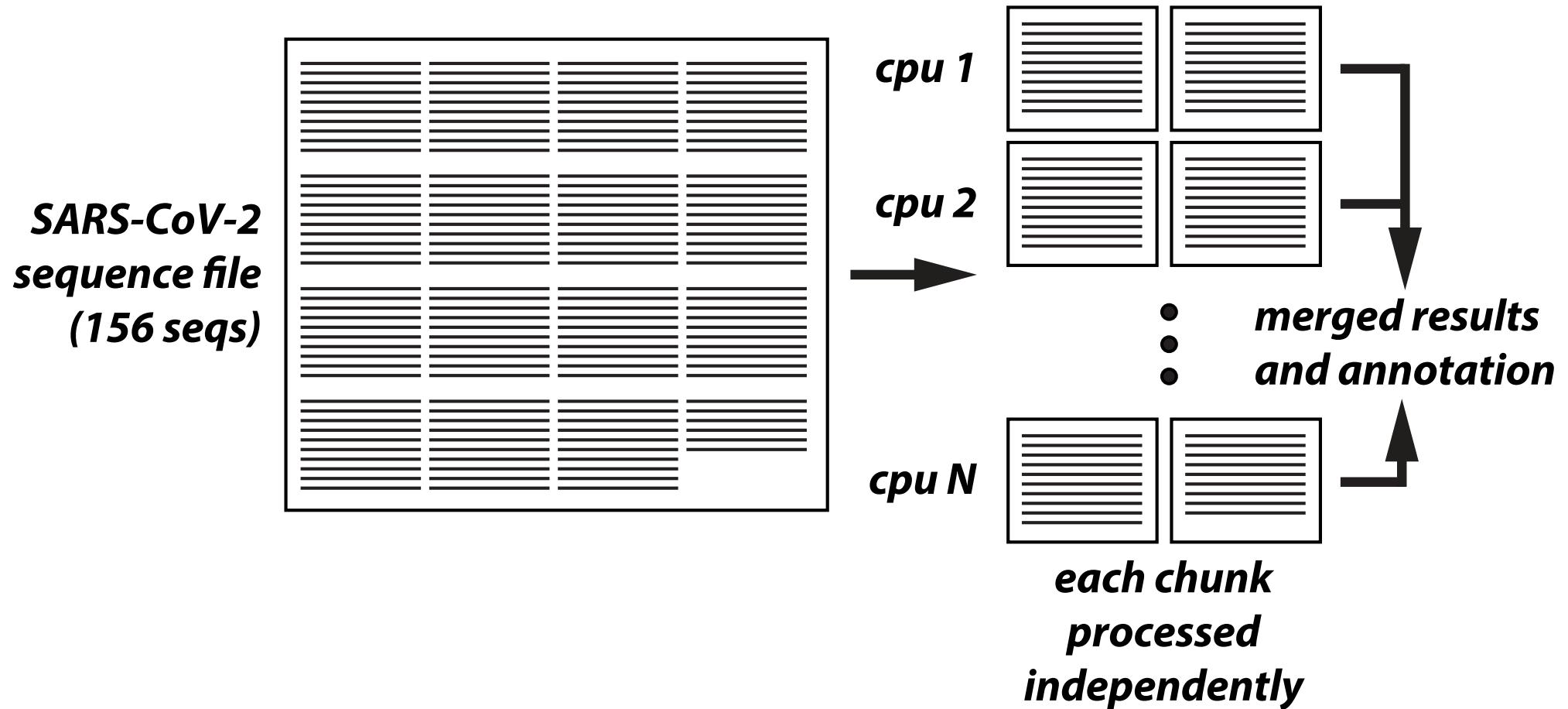


Seeded alignment using blastn makes alignment stage faster



Using glsearch instead of cmalign reduces memory requirement

- lower memory requirement (2Gb max) allows for multi-threading



VADR is now 1000-fold faster in practice for SARS-CoV-2 processing

VADR version	seeded alignment?	N replacement?	glsearch?	# cpus	required RAM	secs per seq	hours per 100K seqs	speedup vs v1.0
v1.0	-	-	-	1	64 Gb	329.91	9164.3	-
v1.4.1	+	+	+	8	16 Gb	0.33	9.3	986.8

VADR is now fast enough to handle hundreds of thousands of sequences per month

month	year	#new seqs	#cumulative seqs
Jan	2020	32	32
Feb	2020	58	90
Mar	2020	332	422
Apr	2020	1541	1963
May	2020	2974	4937
Jun	2020	3394	8331
Jul	2020	3604	11,935
Aug	2020	3818	15,753
Sep	2020	6731	22,484
Oct	2020	11,939	34,423
Nov	2020	4274	38,697
Dec	2020	4530	43,227
Jan	2021	8775	52,002
Feb	2021	26,078	78,080
Mar	2021	42,607	120,687
Apr	2021	97,095	217,782
May	2021	104,729	322,511
Jun	2021	46,187	368,698
Jul	2021	43,336	412,034
Aug	2021	141,958	553,992
Sep	2021	267,562	821,554
Oct	2021	239,296	1,060,850
Nov	2021	267,270	1,328,120
Dec	2021	288,771	1,616,891
Jan	2022	258,522	1,875,413
Feb	2022	230,185	2,105,598

Besides getting faster, VADR has changed in other ways (work with Linda Yankie and Vince Calhoun and GenBank team)

- 13 releases between March 2020 and January 2022
- 3 additional models (all eventually dropped):
 - B.1.1.7 (alpha)
 - B.1.525
 - 28254-deletion
- allow some alerts for non-essential ORFs without failing sequence
(they become a `misc_feature` instead)

Faster SARS-CoV-2 sequence validation and annotation for GenBank using VADR

Eric P. Nawrocki *

National Center for Biotechnology Information, U.S. National Library of Medicine, National Institutes of Health, Bethesda, MD, 20894, USA

Received September 08, 2022; Revised November 28, 2022; Editorial Decision December 14, 2022; Accepted January 03, 2023

Additional VADR models and development

	length	num models	new feature(s)	author
RSV	15Kb	2	alignment-based models	EPN
COX-1	1.5Kb	86	protein-coding gene	EPN
Mpox	197Kb	1	minimap alignment	EPN
Influenza	1-2Kb	70	segmented virus	EPN
Zika	11Kb	?	?	EBD

Additional VADR models and development

	length	num models	new feature(s)	author
RSV	15Kb	2	alignment-based models	EPN
COX-1	1.5Kb	86	protein-coding gene	EPN
Mpox	197Kb	1	minimap alignment	EPN
Influenza	1-2Kb	70	segmented virus	EPN
Zika	11Kb	?	?	EBD

Database, 2024, baae091

DOI: <https://doi.org/10.1093/database/baae091>

Original article



Influenza sequence validation and annotation using VADR

Vincent C. Calhoun, Eneida L. Hatcher, Linda Yankie, Eric P. Nawrocki *

National Center for Biotechnology Information, U.S. National Library of Medicine, National Center for Biotechnology Information, 8600 Rockville Pike, Bethesda, MD 20894, United States

Additional VADR models and development

- Alex Greninger's lab at Univ of Washington:
 - sequences a wide variety of human pathogenic viruses
 - previously developed the VAPiD software tool for validating and annotating viral sequences
 - now a collaborator that builds VADR models

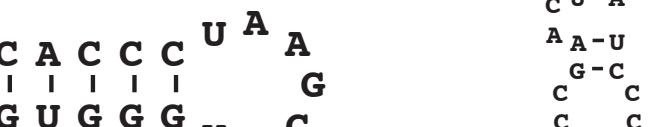
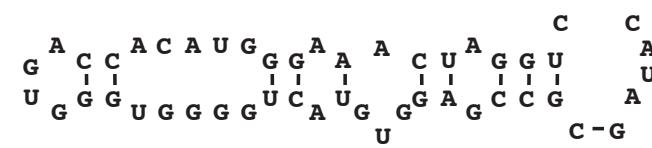
Additional VADR models and development

	length	num models	new feature(s)	author
RSV	15Kb	2	alignment-based models	EPN
COX-1	1.5Kb	86	protein-coding gene	EPN
Mpox	197Kb	1	minimap alignment	EPN
Influenza	1-2Kb	70	segmented virus	EPN
Zika	11Kb	?	?	EBD
Herpes Simplex Virus (HSV)	150Kb	2	-	Greninger Lab
Human meta-pneumovirus (HMPV)	13Kb	6	-	Greninger Lab

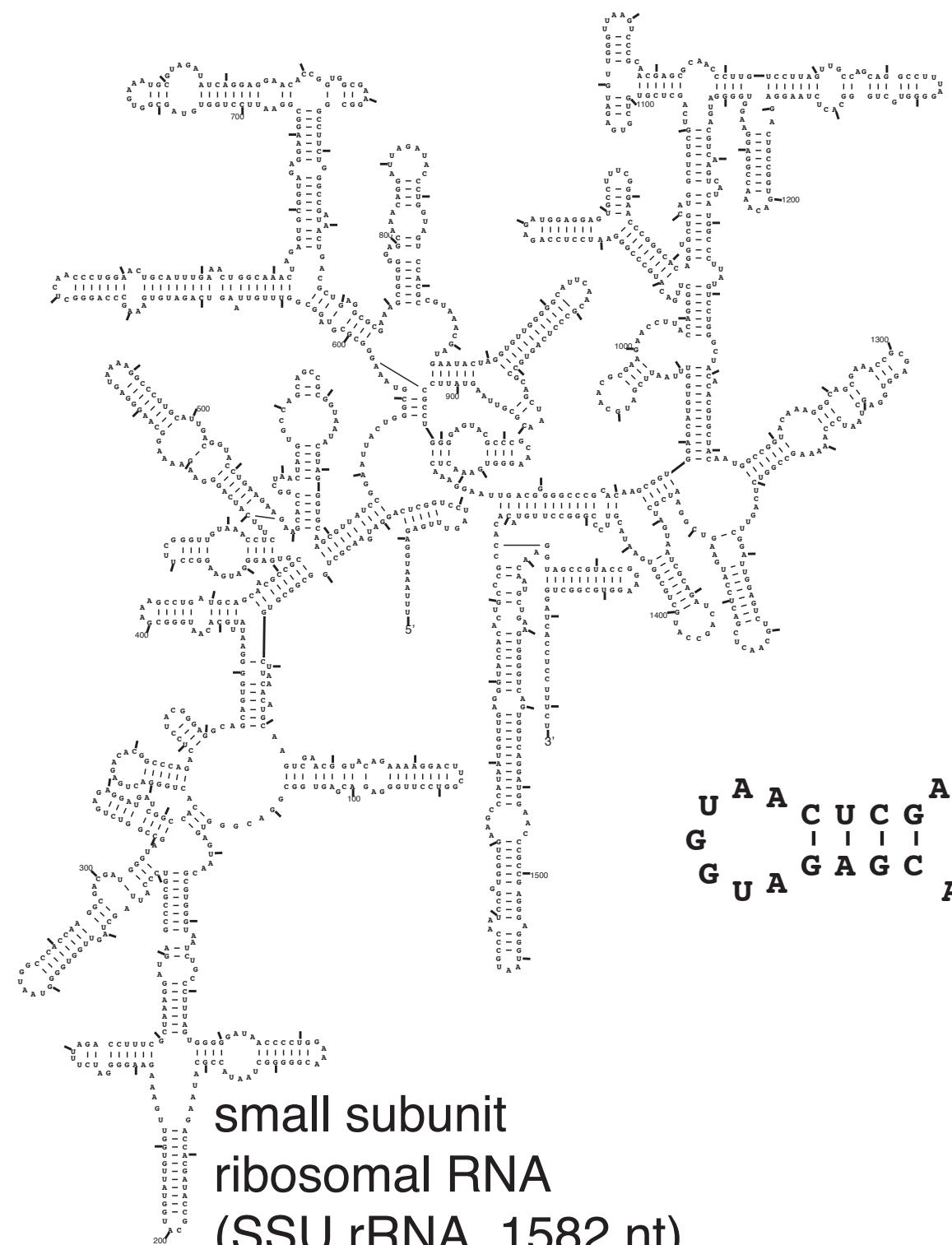
Future directions for VADR

- NCBI Virus FY2025 goals:
 - VADR web server
 - Replacement of FLAN with VADR
- Development of and application of models for additional viruses (our group and Greninger lab)

5S ribosomal RNA (119 nt)



**small subunit
ribosomal RNA
(SSU rRNA, 1582 nt)**



Functional RNAs play many vital roles in the cell

	key RNAs involved	archaea	bacteria	eukarya	viruses
translation	ribosomal RNAs	x	x	x	
	transfer RNAs	x	x	x	
	RNase P RNA	x	x	x	
	snoRNAs	x		x	
	SRP RNA	x	x	x	
	tmRNA		x		
	RNaseMRP			x	
gene expression	riboswitches	?	x	?	
	microRNAs			x	x
	6S RNA		x	x	
splicing	U1, U2, U4, U5, U6			x	
other	tracrRNA	x	x		
	telomerase RNA			x	
	group I introns	x	x	x	x
	sfRNAs				x
	many more...				x

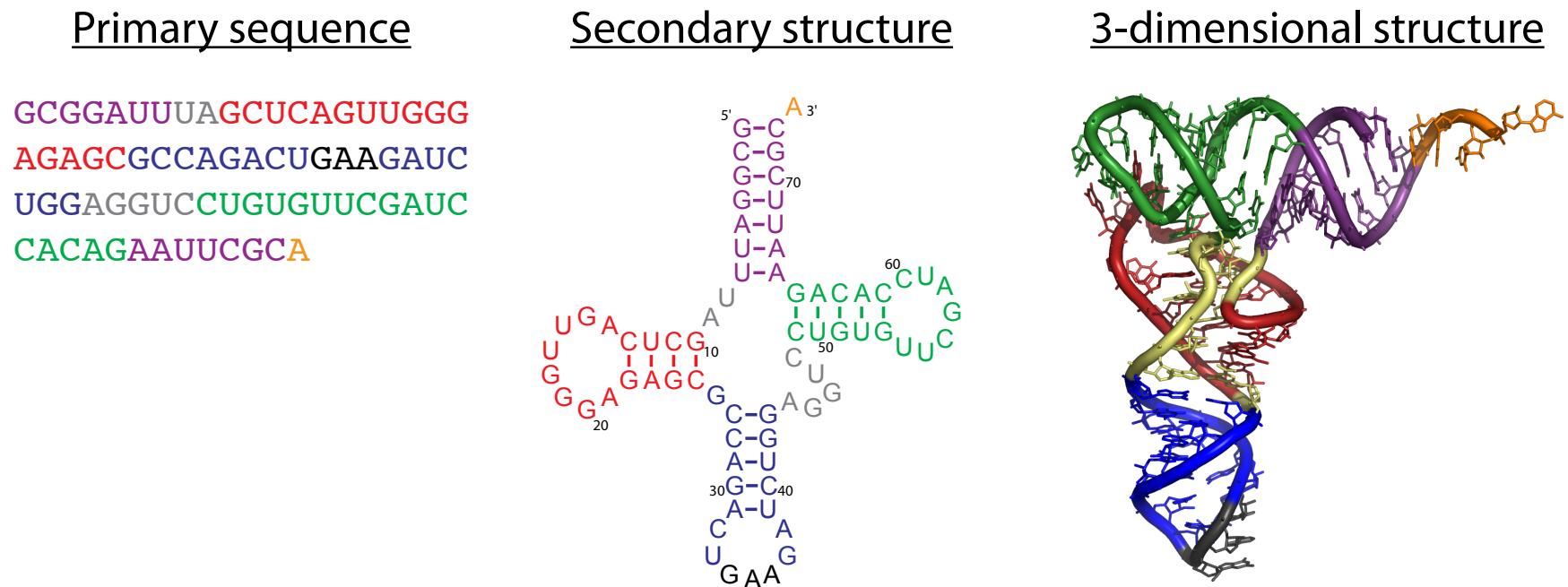
Functional RNAs play many vital roles in the cell

	key RNAs involved	archaea	bacteria	eukarya	viruses
translation	ribosomal RNAs	x	x	x	
	transfer RNAs	x	x	x	
	RNase P RNA	x	x	x	
	snoRNAs	x		x	
	SRP RNA	x	x	x	
	tmRNA		x		
	RNaseMRP			x	
gene expression	riboswitches	?	x	?	
	microRNAs			x	x
	6S RNA		x	x	
splicing	U1, U2, U4, U5, U6			x	
other	tracrRNA	x	x		
	telomerase RNA			x	
	group I introns	x	x	x	x
	sfRNAs				x
	many more...				x



database of more than 4100 non-coding RNA families
each represented by a secondary structure, alignment, and covariance model.

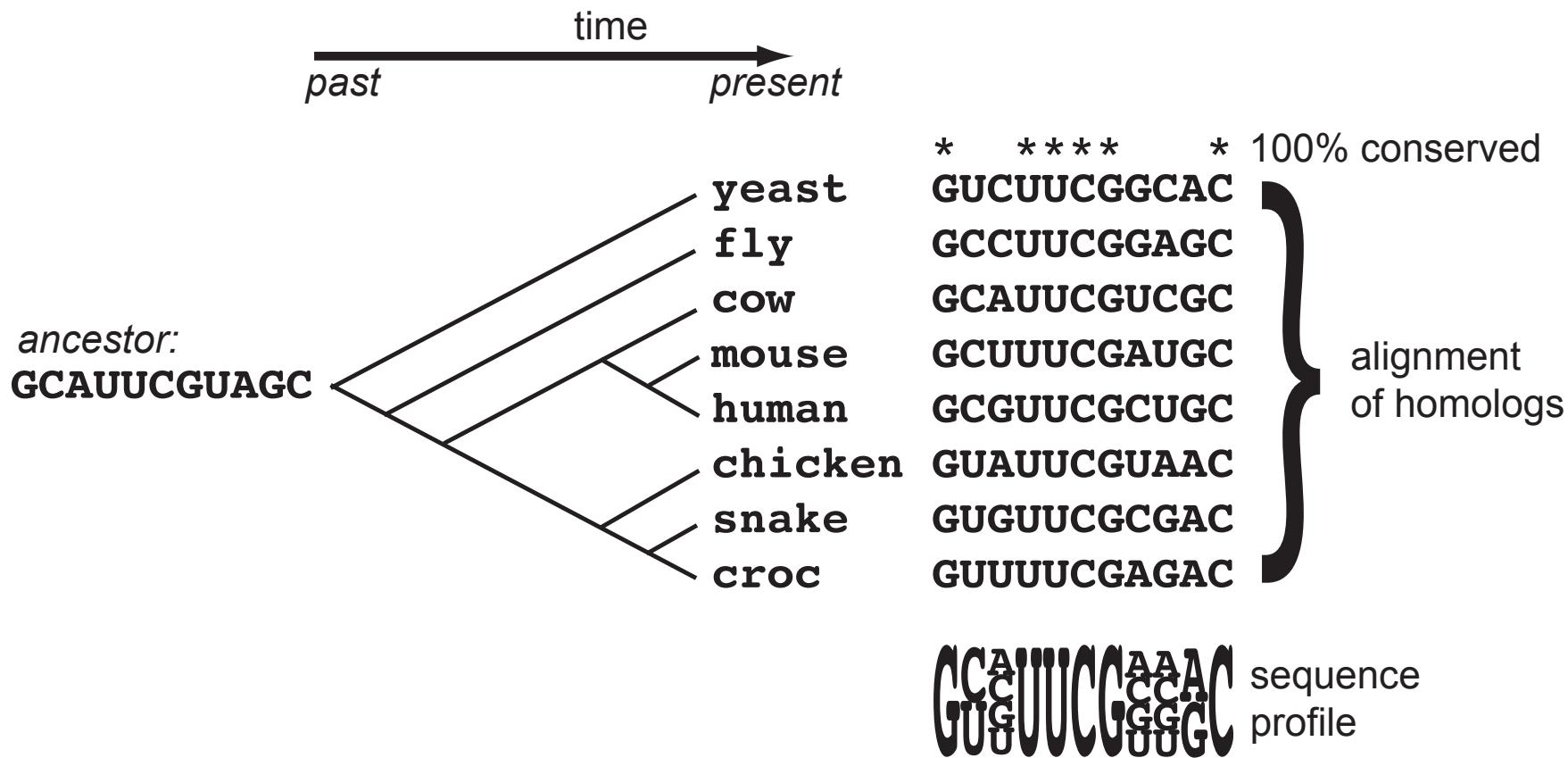
Many functional RNAs adopt a conserved 3-dimensional structure



- BLAST: given a single sequence, search genomes for similar sequences.
- Structural RNAs are difficult to find
 - short (~ 100 nt) and evolve rapidly at sequence level
 - lack open reading frames
 - small, 4 letter alphabet
- BLAST cannot take advantage of:
 - sequence conservation, which varies across the gene
 - secondary structure

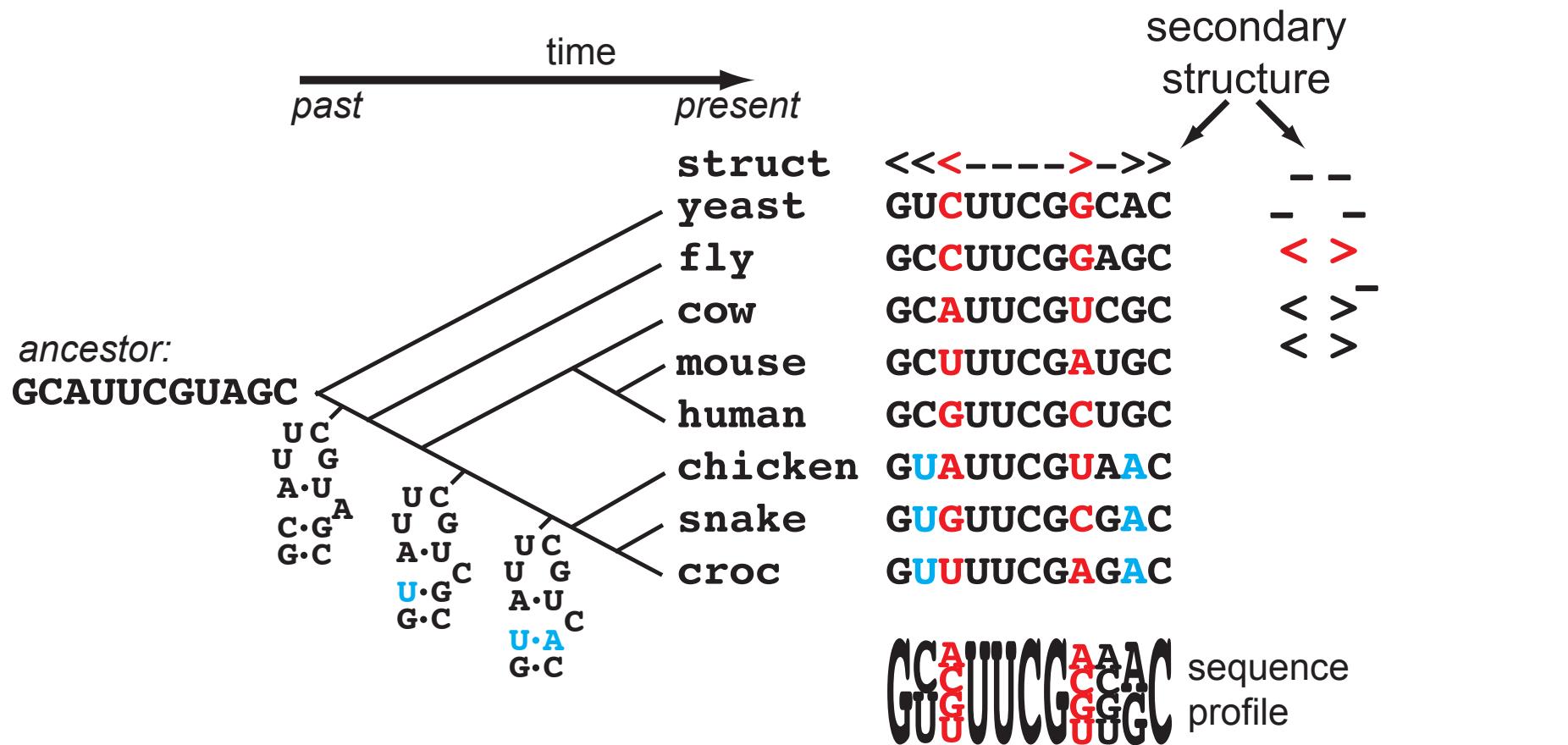
Sequence conservation provides information for homology searches

Conservation levels vary across alignment columns.



Structure conservation provides additional information

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



profile HMMs and covariance models

	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 (23794 and 4150 entries)	Rfam (4178 families)
performance for RNAs	faster but less accurate	slower but more accurate

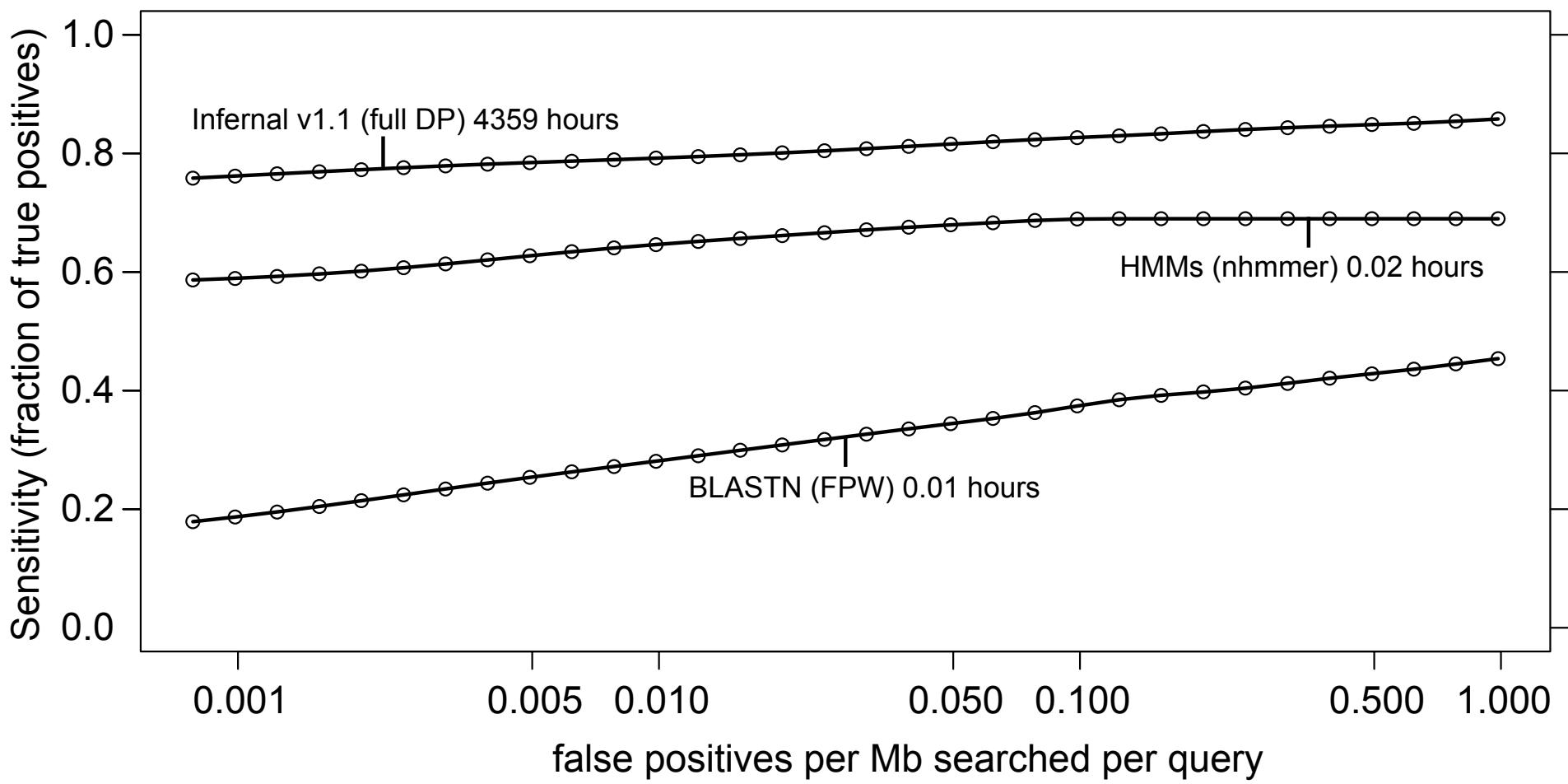


<http://hmmer.org>
Potter et al. NAR
46:W200-204
Wheeler, TJ, Eddy SR.
Bioinformatics, 29:2487-89, 2013.
Eddy, SR. PLoS Comp. Biol.,
7:e1002195, 2011.
Eddy, SR. Bioinformatics,
14:755-763, 1998.



<http://eddylab.org/infernal/>
Nawrocki EP, Eddy SR.
Bioinformatics, 29:
2487-2489, 2013.
Eddy SR, Durbin R.
Nucleic Acids Research,
22:2079-2088, 1994.

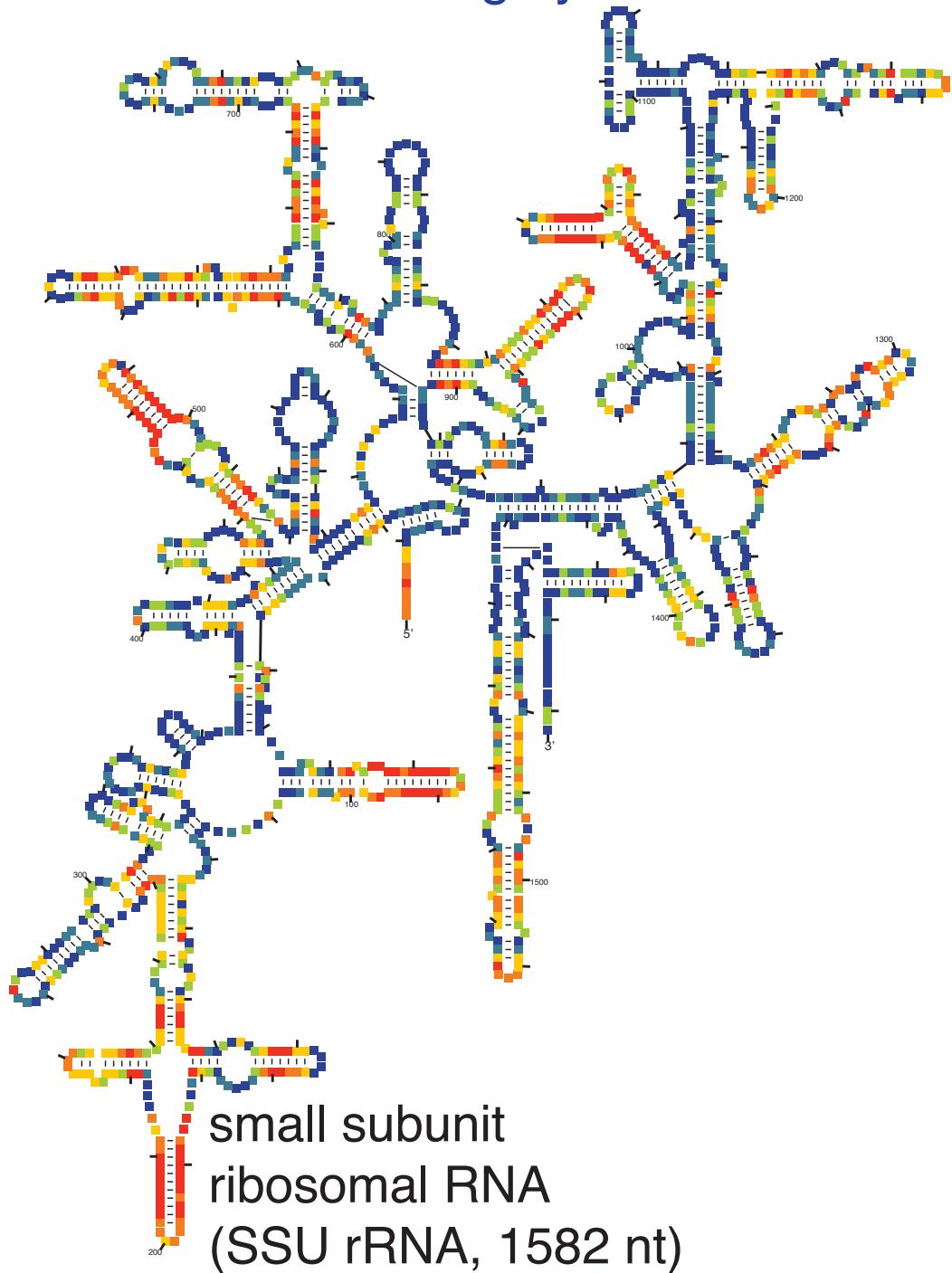
Infernal outperforms primary-sequence based methods on our benchmark (and others*, not shown)



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

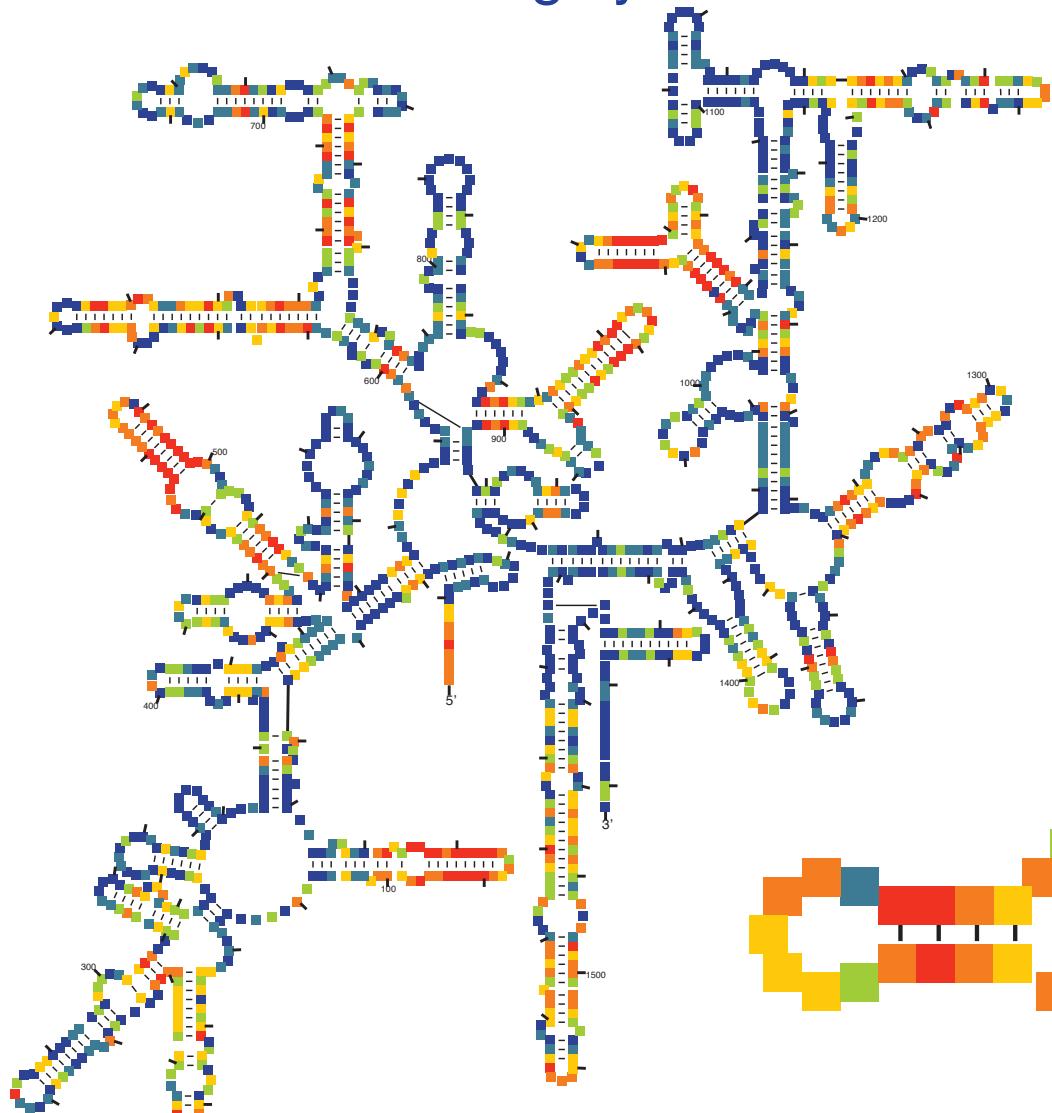
Sequence conservation per position

blue:highly conserved red: highly variable

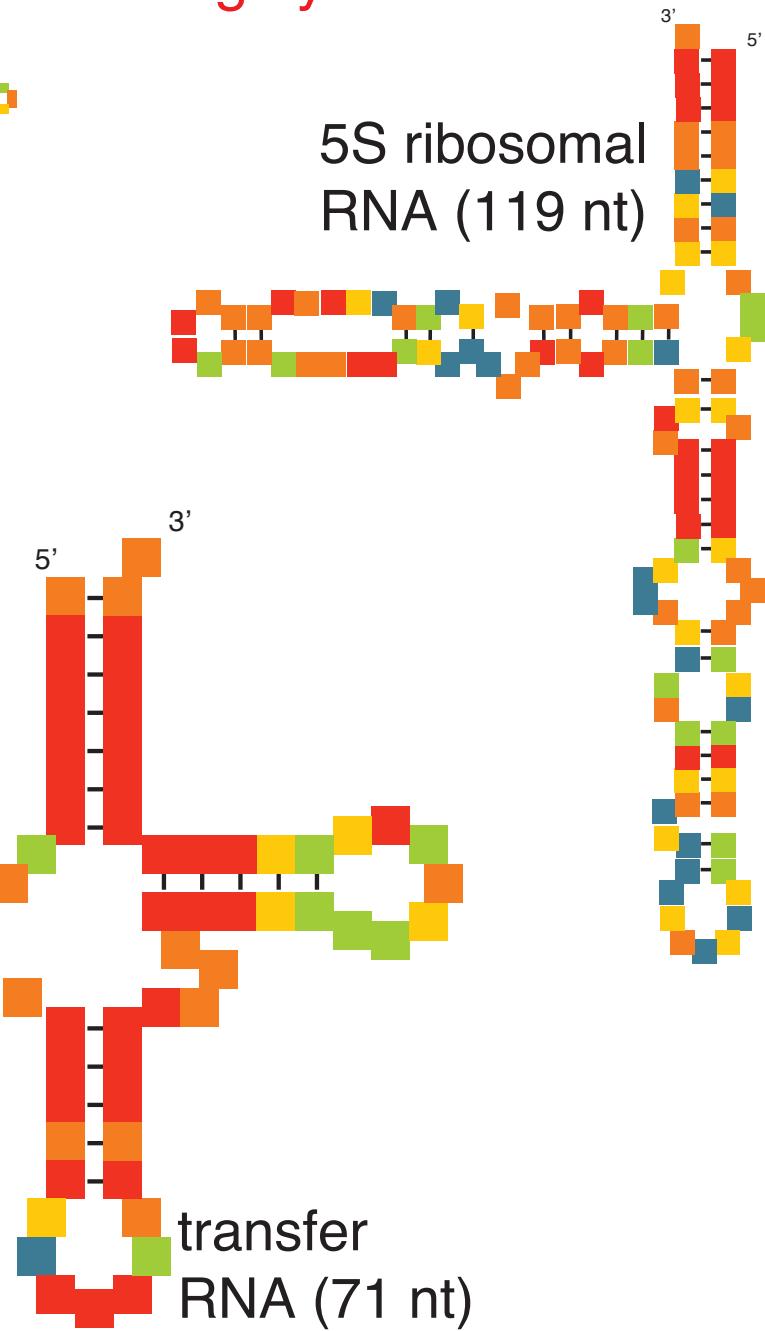


Sequence conservation per position

blue:highly conserved red: highly variable



small subunit
ribosomal RNA
(SSU rRNA, 1582 nt)

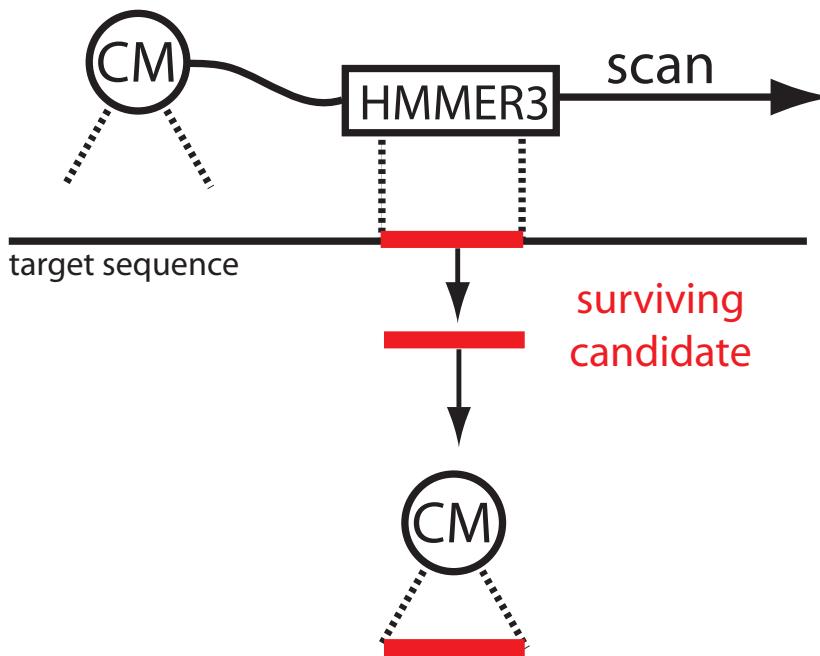


5S ribosomal
RNA (119 nt)

transfer
RNA (71 nt)

Filter target database using profile HMMs*

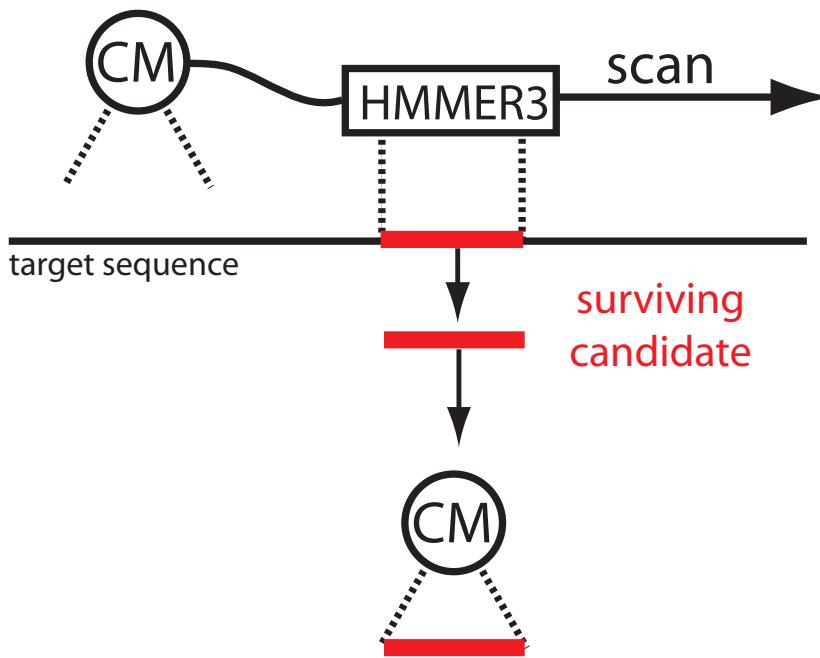
HMM filter first pass



*Weinberg, Ruzzo, RECOMB, 243-251, 2004; Weinberg, Ruzzo, Bioinformatics, 22(1) 35-39 2006.

Filter target database using profile HMMs*

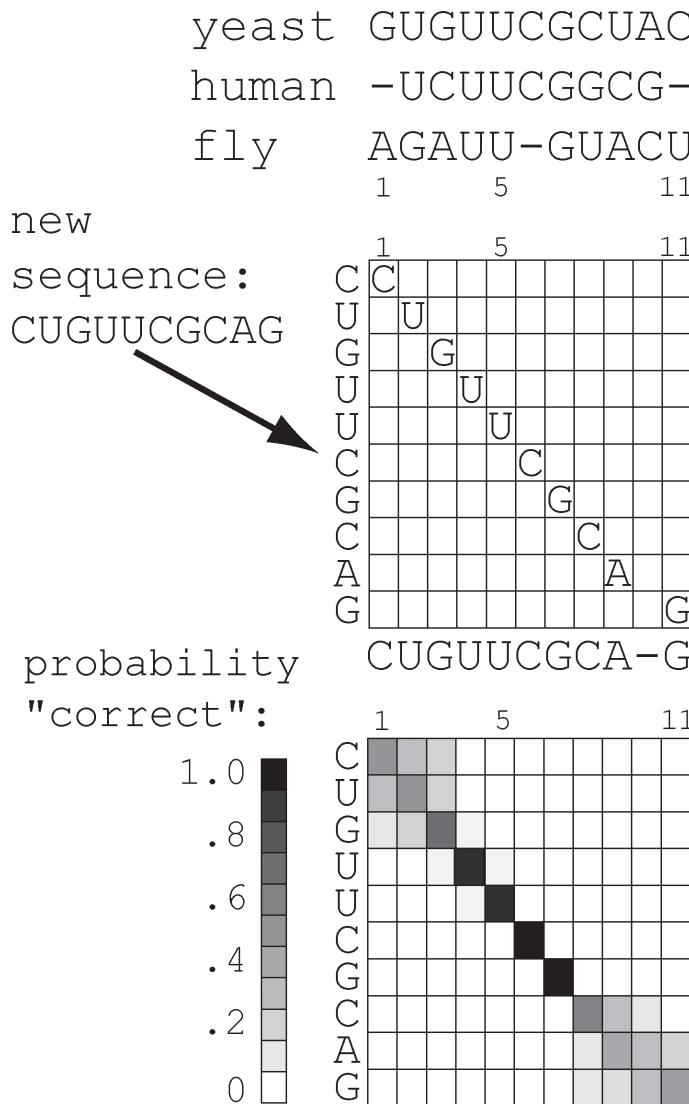
HMM filter first pass



- 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.

*Weinberg, Ruzzo, RECOMB, 243-251, 2004; Weinberg, Ruzzo, Bioinformatics, 22(1) 35-39 2006.

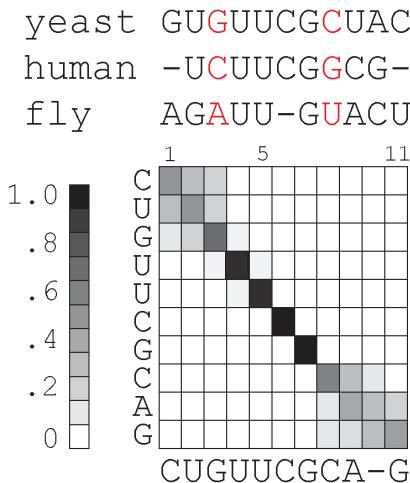
Accelerating CM alignment step 1: HMM posterior decoding to get confidence estimates



Accelerating CM alignment step 2: use HMM alignment confidence to constrain CM alignment

HMMs -

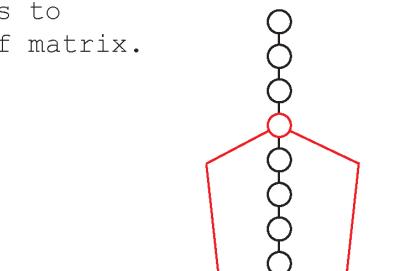
Each column of seed alignment corresponds to a column of matrix.



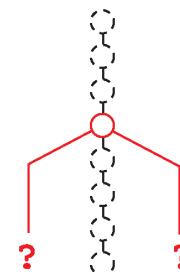
CMs -

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 AG**A**UU-G**U**ACU



CUGUUCGCAG
 45 possibilities

Accelerating CM alignment step 2: use HMM alignment confidence to constrain CM alignment

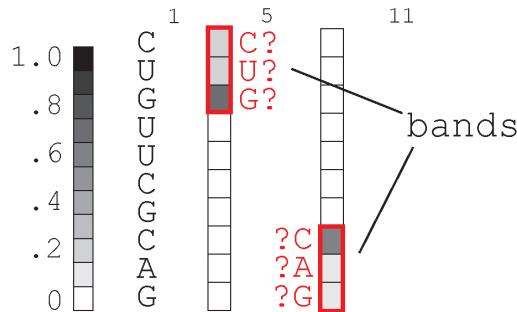
HMMs -

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

yeast GUGUUCGCUAC

human -UCUUCGGCG-

fly AGAUU-GUACU



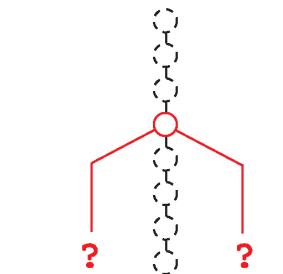
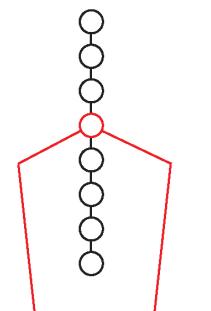
struct <<----->->>
yeast GUGUUCGCUAC
human -UCUUCGGCG-
fly AGAUU-GUACU

CMs -

Each column of seed alignment corresponds to a state.

yeast human fly

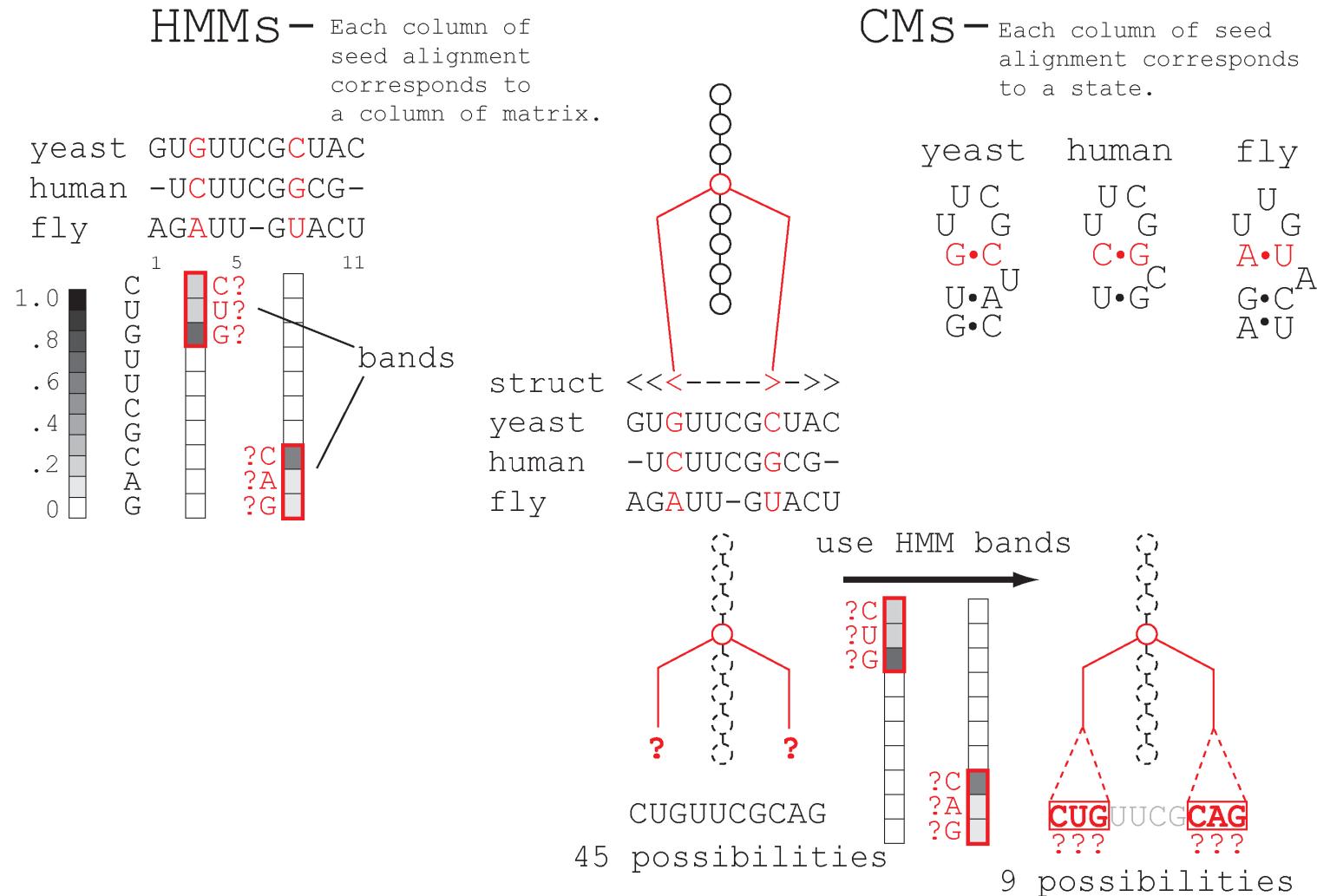
U C	U C	U
U G	U G	U G
G•C	C•G	A•U
U•A U	U•G C	G•C A
G•C		A•U



CUGUUCGCAG

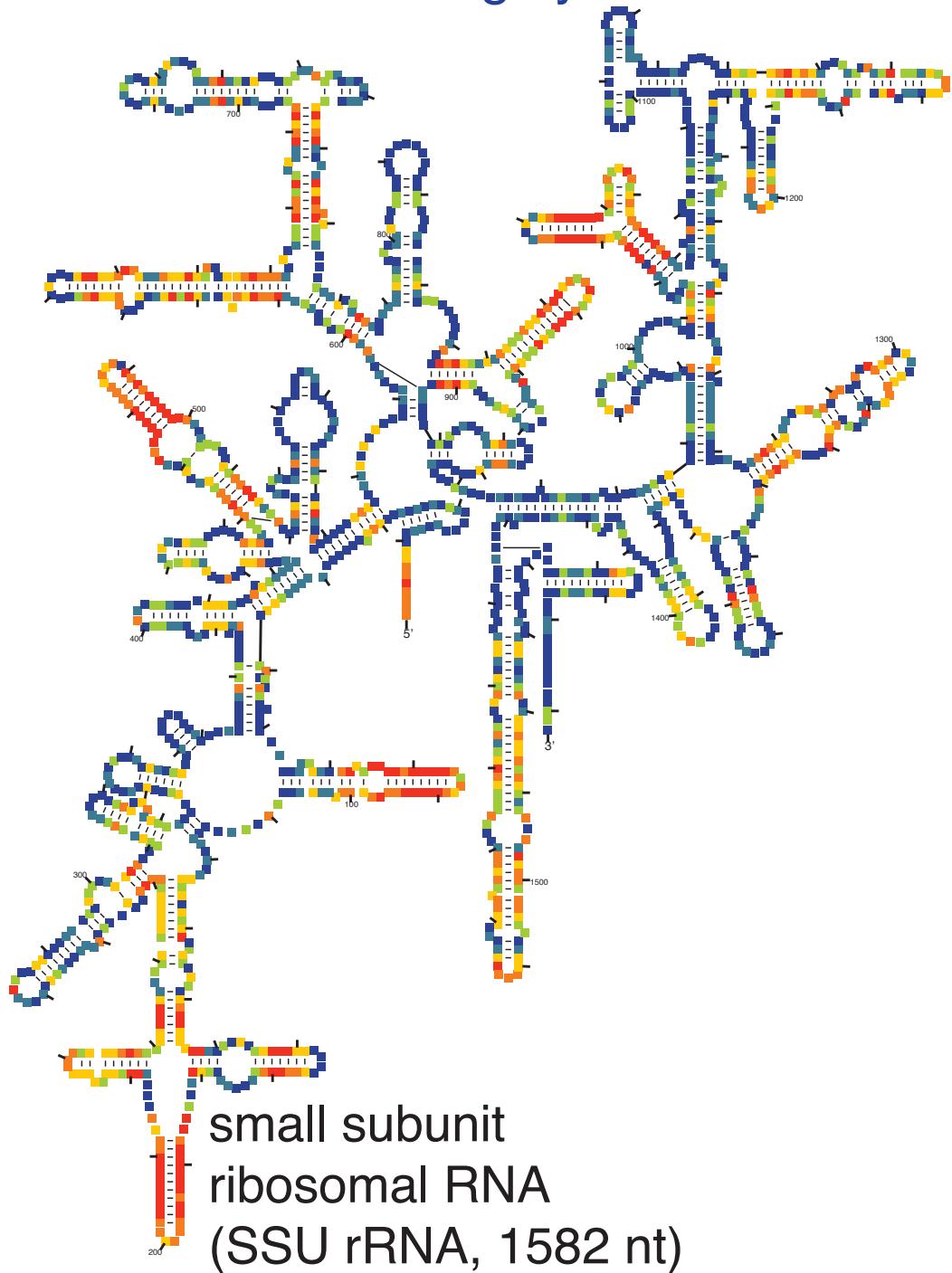
45 possibilities

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



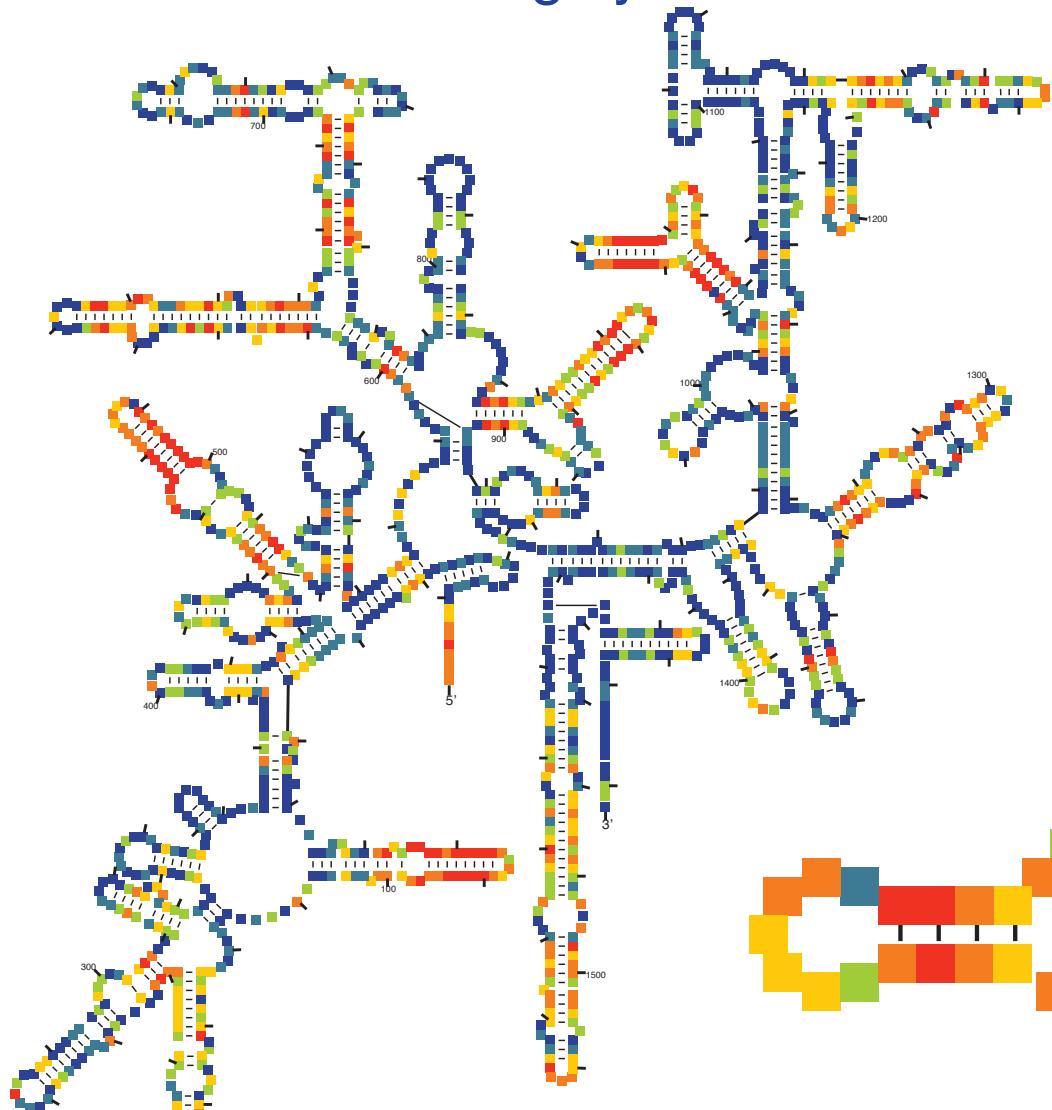
Sequence conservation per position

blue:highly conserved red: highly variable

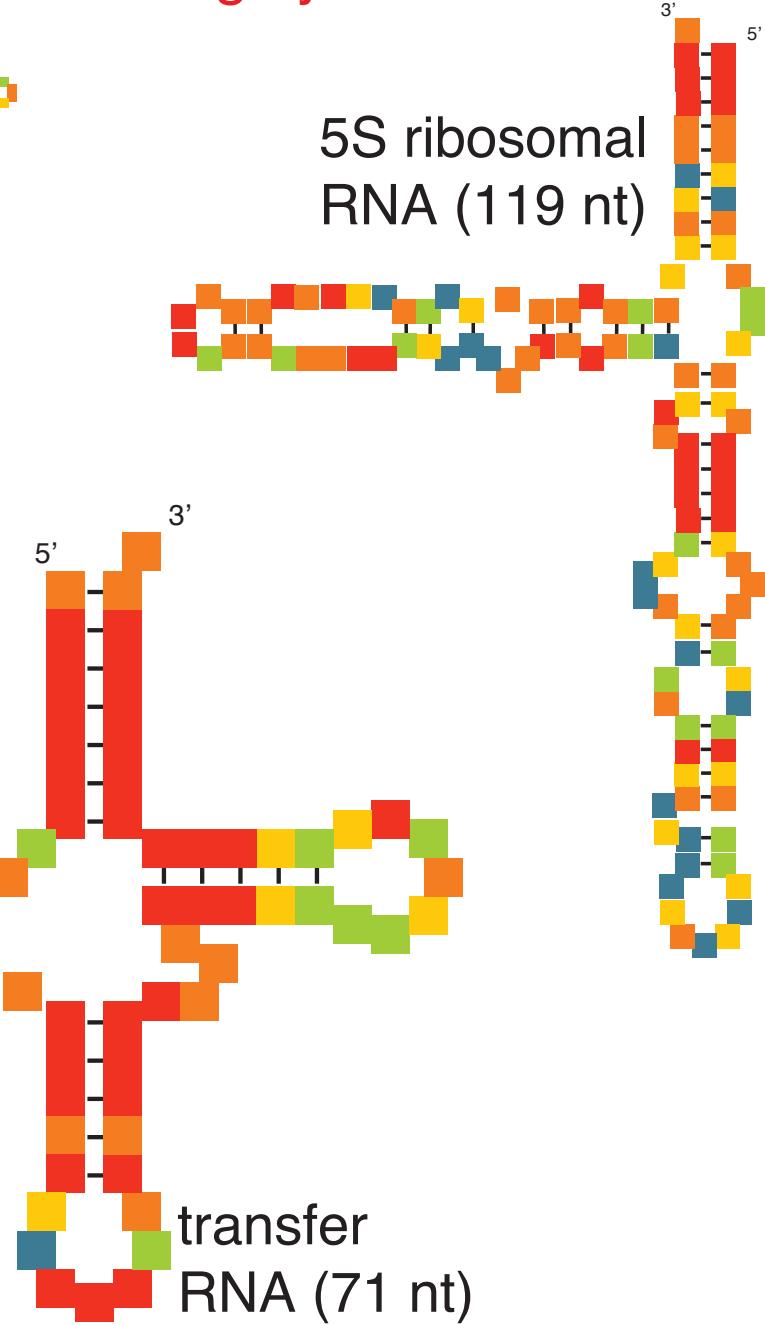


Sequence conservation per position

blue:highly conserved red: highly variable



small subunit
ribosomal RNA
(SSU rRNA, 1582 nt)

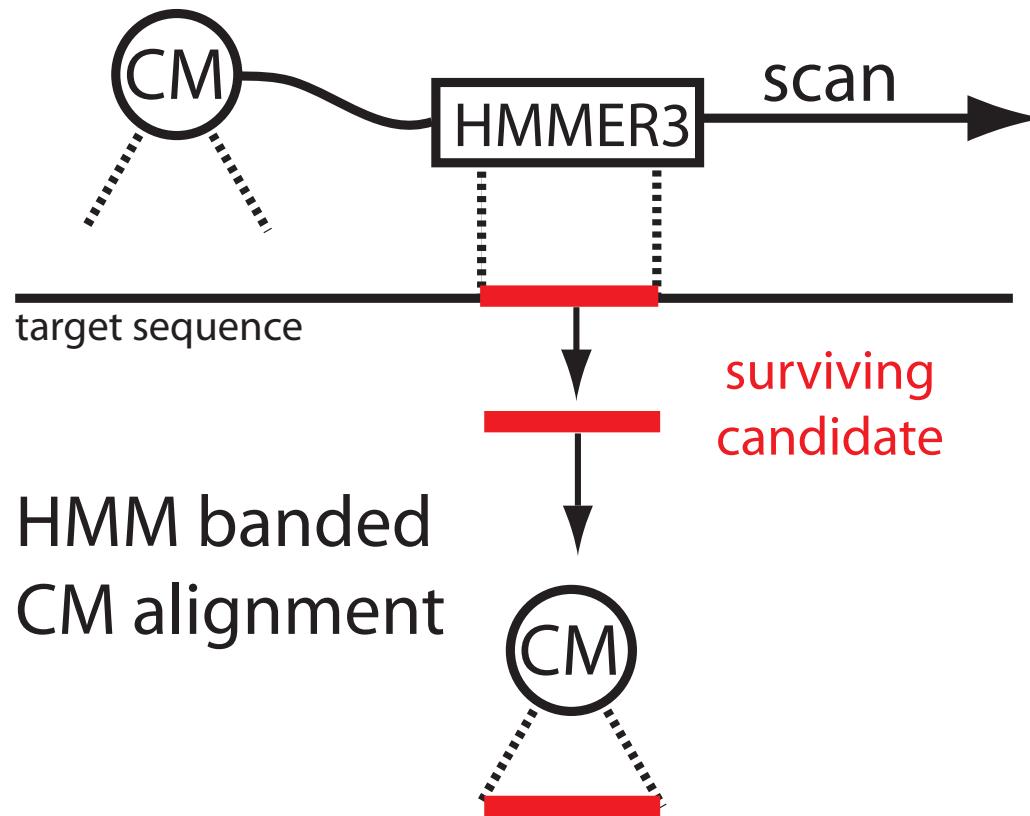


5S ribosomal
RNA (119 nt)

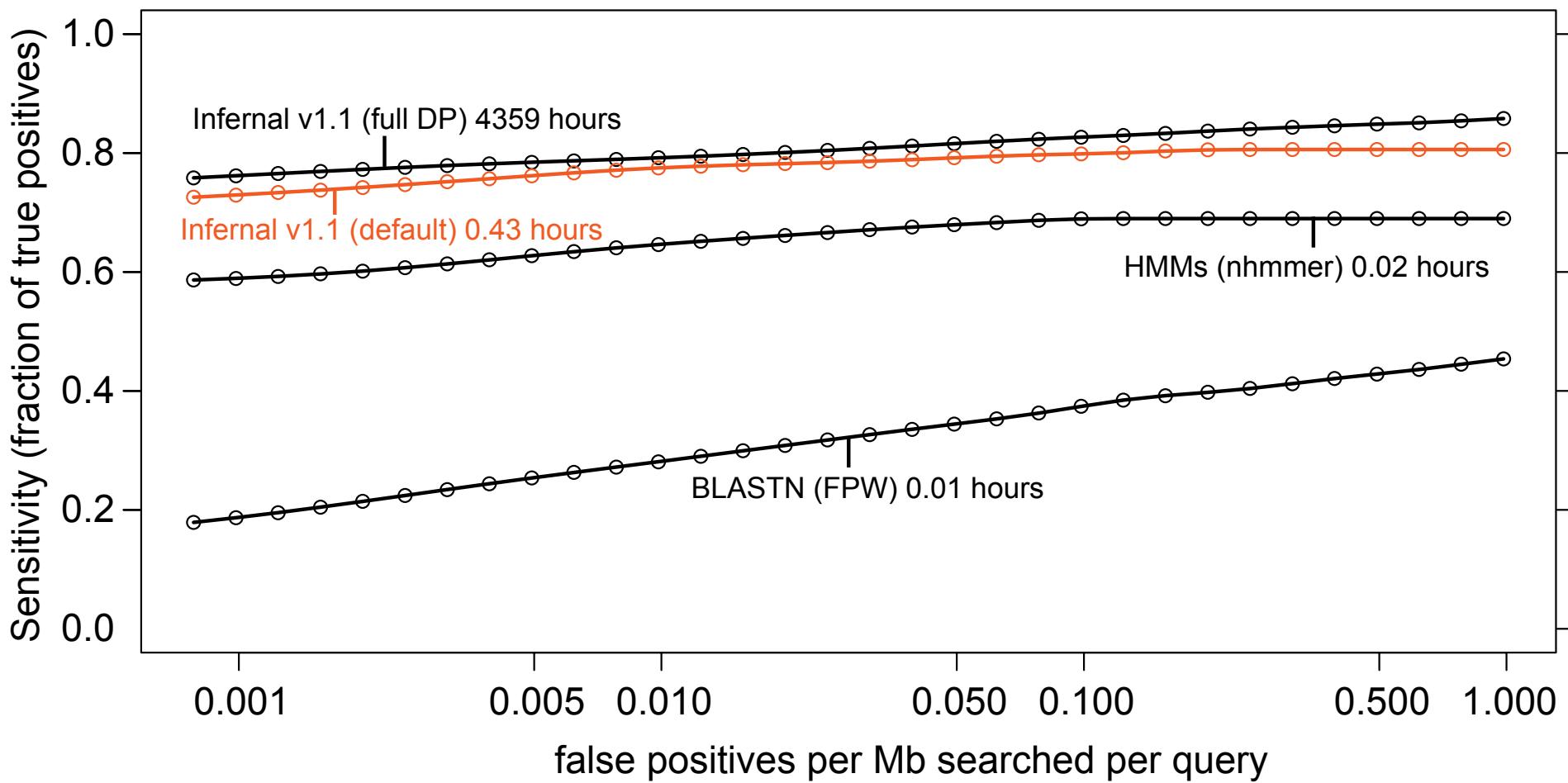
transfer
RNA (71 nt)

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.

Practical structural RNA genome annotation

- Faster Infernal integrated into:
 - NCBI prokaryotic genome annotation pipeline PGAP
 - NCBI eukaryotic genome annotation (Françoise Thibaud-Nissen)

6614–6624 *Nucleic Acids Research*, 2016, Vol. 44, No. 14
doi: 10.1093/nar/gkw569

Published online 24 June 2016

NCBI prokaryotic genome annotation pipeline

Tatiana Tatusova^{1,†}, Michael DiCuccio^{1,†}, Azat Badretdin¹, Vyacheslav Chetvernin¹, Eric P. Nawrocki¹, Leonid Zaslavsky¹, Alexandre Lomsadze², Kim D. Pruitt¹, Mark Borodovsky^{2,3,*‡} and James Ostell^{1,‡}

¹National Center for Biotechnology Information, U.S. National Library of Medicine, Bethesda, MD 20894, USA,

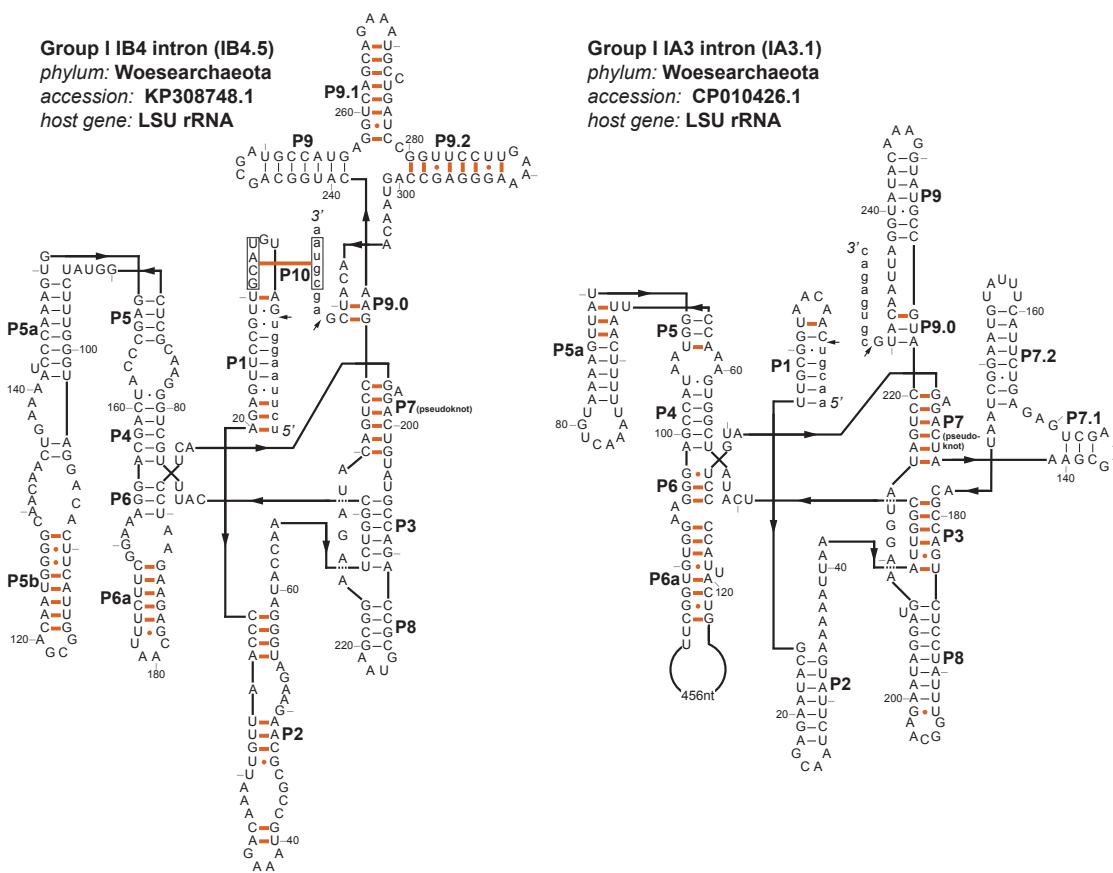
²Wallace H. Coulter Department of Biomedical Engineering, Georgia Tech, Atlanta, GA 30332, USA and ³School of Computational Science and Engineering, Georgia Tech, Atlanta, GA 30332, USA

Group I introns are widespread in archaea

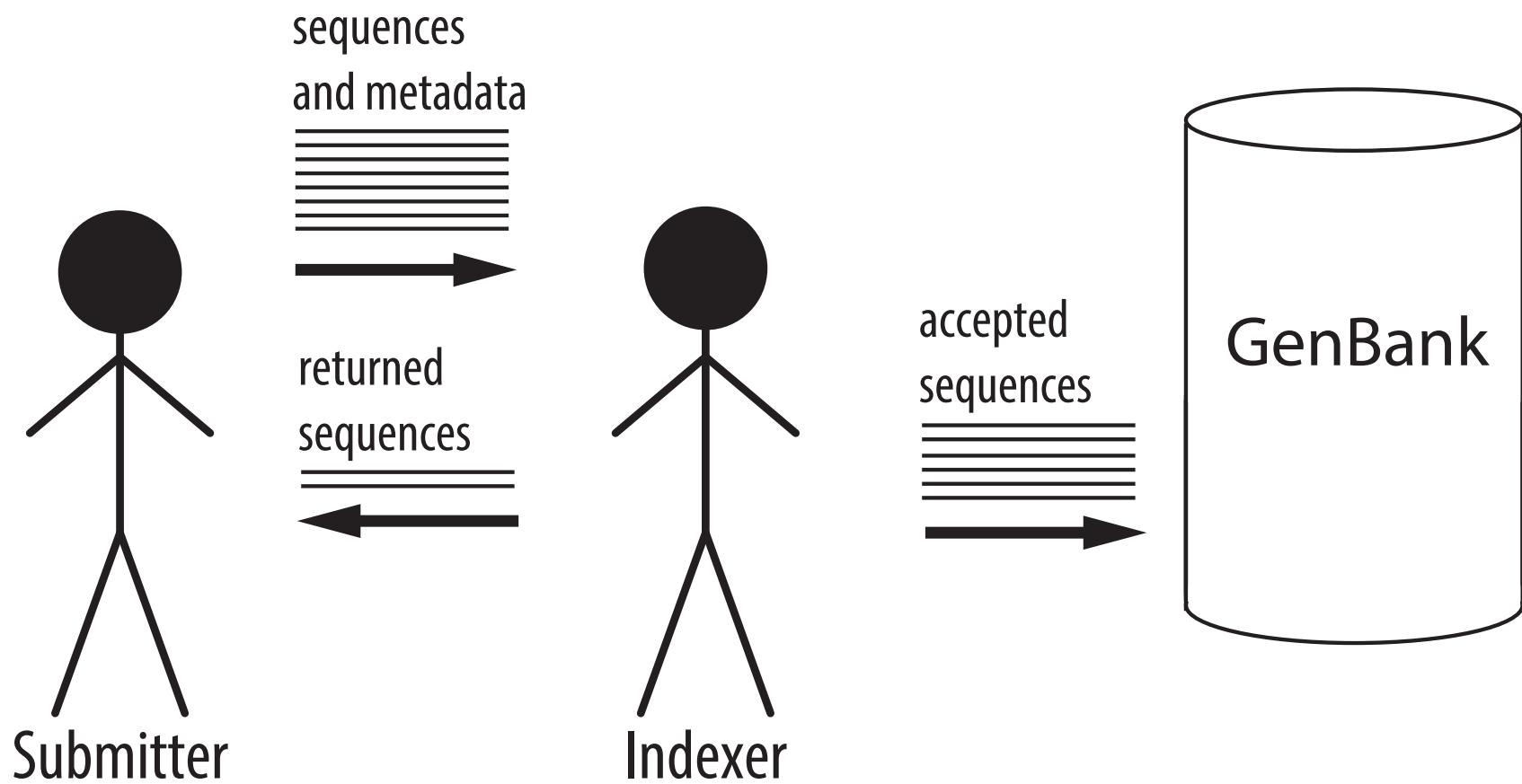
Eric P. Nawrocki^{1,*}, Thomas A. Jones^{2,3} and Sean R. Eddy^{2,3,4,*}

¹National Center for Biotechnology Information, U.S. National Library of Medicine, Bethesda, MD 20894, USA,

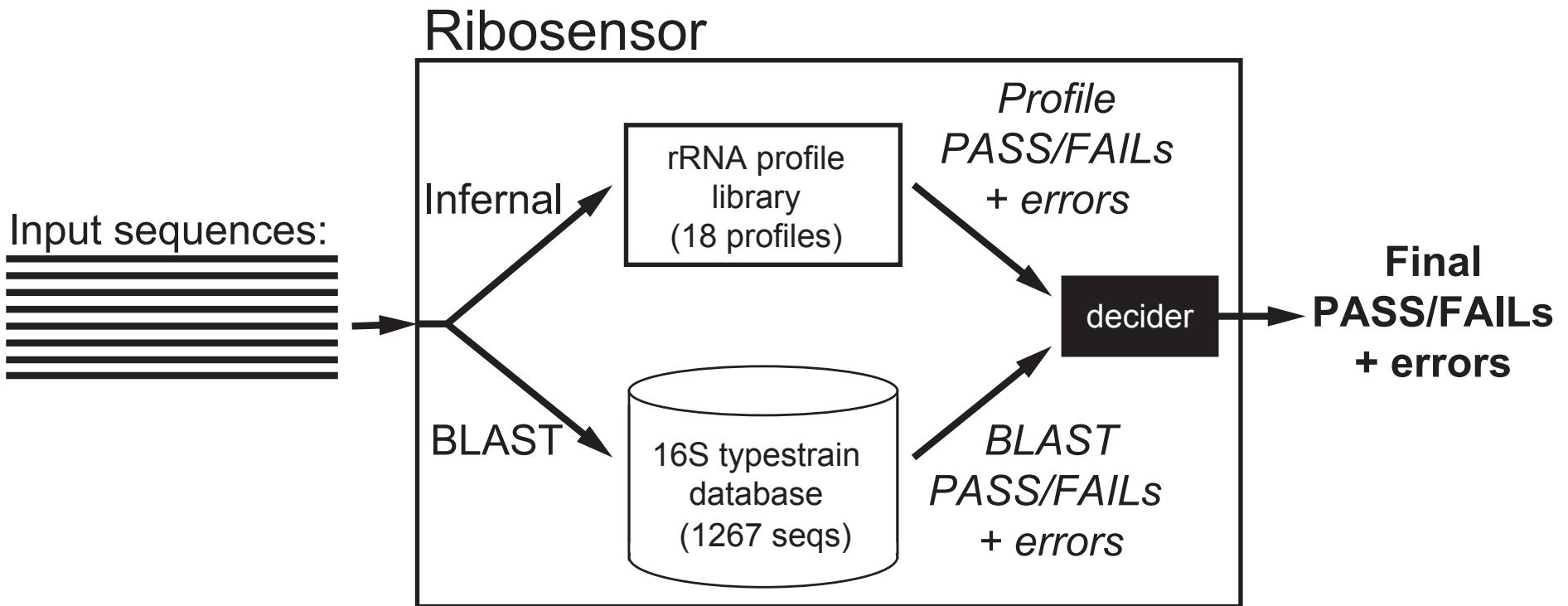
²Howard Hughes Medical Institute, Harvard University, Cambridge, USA, ³Department of Molecular and Cellular Biology, Harvard University, Cambridge, USA and ⁴School of Engineering and Applied Sciences, Harvard University, Cambridge, USA



GenBank indexers handle incoming sequence submissions



Ribosensor: a tool for evaluating ribosomal RNA datasets using profiles and BLAST*



- Profile-based analysis:
 - 18 ribosomal RNA models (15 SSU rRNA, 3 LSU rRNA); 8 from Rfam
 - Detects unexpected features ("UnacceptableModel", "DuplicatedRegions", etc.)
- Profile and BLAST results considered together to determine PASS/FAIL

*Alejandro Schäffer developed the BLAST-based scheme

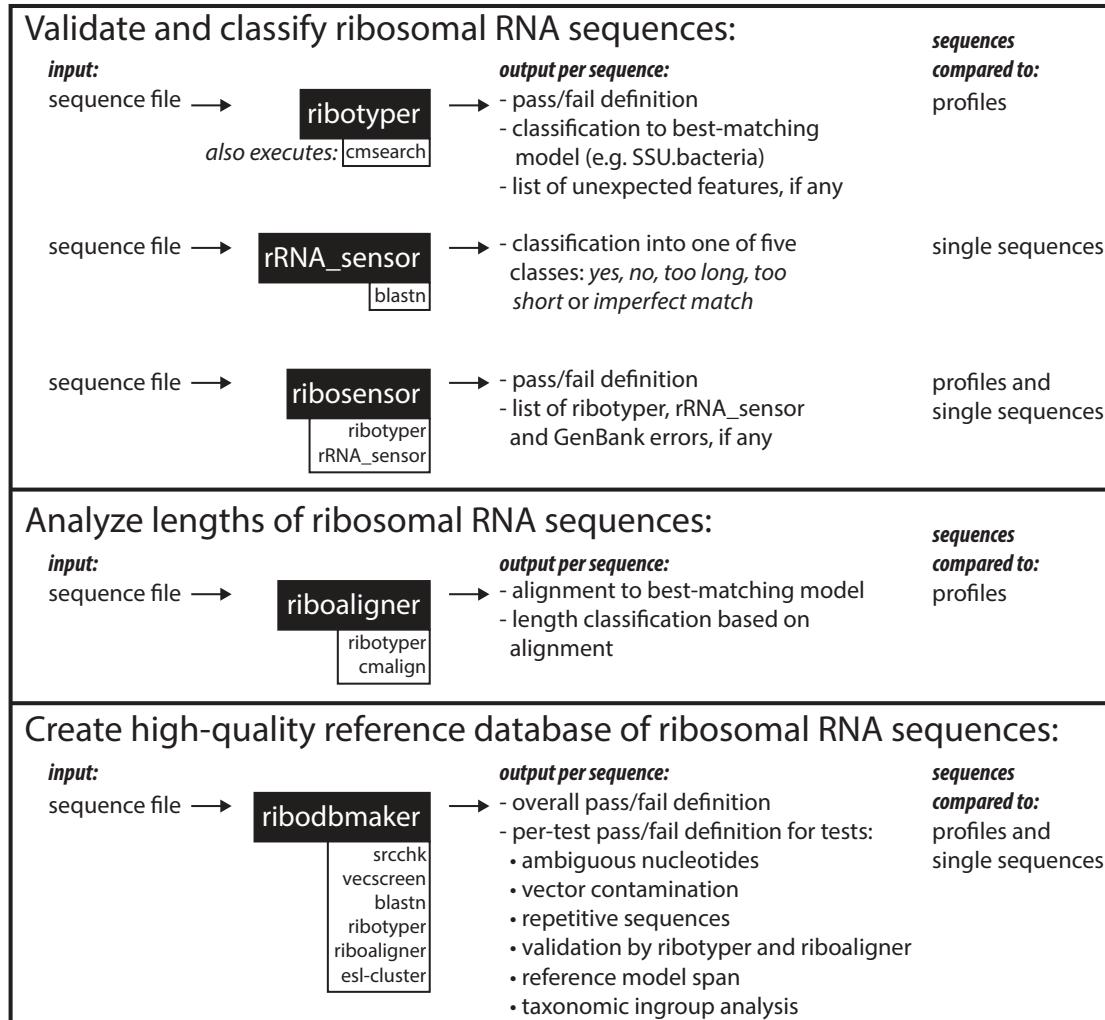
SOFTWARE

Open Access



Ribovore: ribosomal RNA sequence analysis for GenBank submissions and database curation

Alejandro A. Schäffer^{1,2}, Richard McVeigh², Barbara Robbertse², Conrad L. Schoch², Anjanette Johnston², Beverly A. Underwood², Ilene Karsch-Mizrachi² and Eric P. Nawrocki^{2*} 



Rfam 15: RNA families database in 2025

Nancy Ontiveros-Palacios  ¹, **Emma Cooke**  ², **Eric P. Nawrocki**  ³, **Sandra Triebel**  ^{4,5},
Manja Marz  ^{4,5}, **Elena Rivas**  ⁶, **Sam Griffiths-Jones**  ⁷, **Anton I. Petrov**  ⁸, **Alex Bateman**  ¹ and
Blake Sweeney  ^{1,*}

¹European Molecular Biology Laboratory, European Bioinformatics Institute (EMBL-EBI), Wellcome Genome Campus, Hinxton, Cambridge CB10 1SD, UK

²SciBite Limited, BioData Innovation Centre, Wellcome Genome Campus, Hinxton, Cambridge CB10 1DR, UK

³National Center for Biotechnology Information, U.S. National Library of Medicine, National Institutes of Health, Bethesda, MD 20894, USA

⁴RNA Bioinformatics and High-Throughput Analysis, Friedrich Schiller University Jena, 07743 Jena, Germany

⁵European Virus Bioinformatics Center, Friedrich Schiller University Jena, 07743 Jena, Germany

⁶Department of Molecular and Cellular Biology, Harvard University, Cambridge, MA 02138, USA

⁷School of Biological Sciences, Faculty of Medicine, Biology and Health, Michael Smith Building, The University of Manchester, Dover St, Manchester M13 9NT, UK

⁸Riboscope Ltd, Cambridge CB1 1AH, UK

*To whole correspondence should be addressed. Tel: +44 1223 494359; Email: bsweeney@ebi.ac.uk

RNAcentral: a hub of information for non-coding RNA sequences

The RNAcentral Consortium^{1–38,*}

D212–D220 Nucleic Acids Research, 2021, Vol. 49, Database issue
doi: 10.1093/nar/gkaa921

Published online 27 October 2020

RNAcentral 2021: secondary structure integration, improved sequence search and new member databases

RNAcentral Consortium^{1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,19,20,21,22,23,24,25,26,*}

ARTICLE



<https://doi.org/10.1038/s41467-021-23555-5>

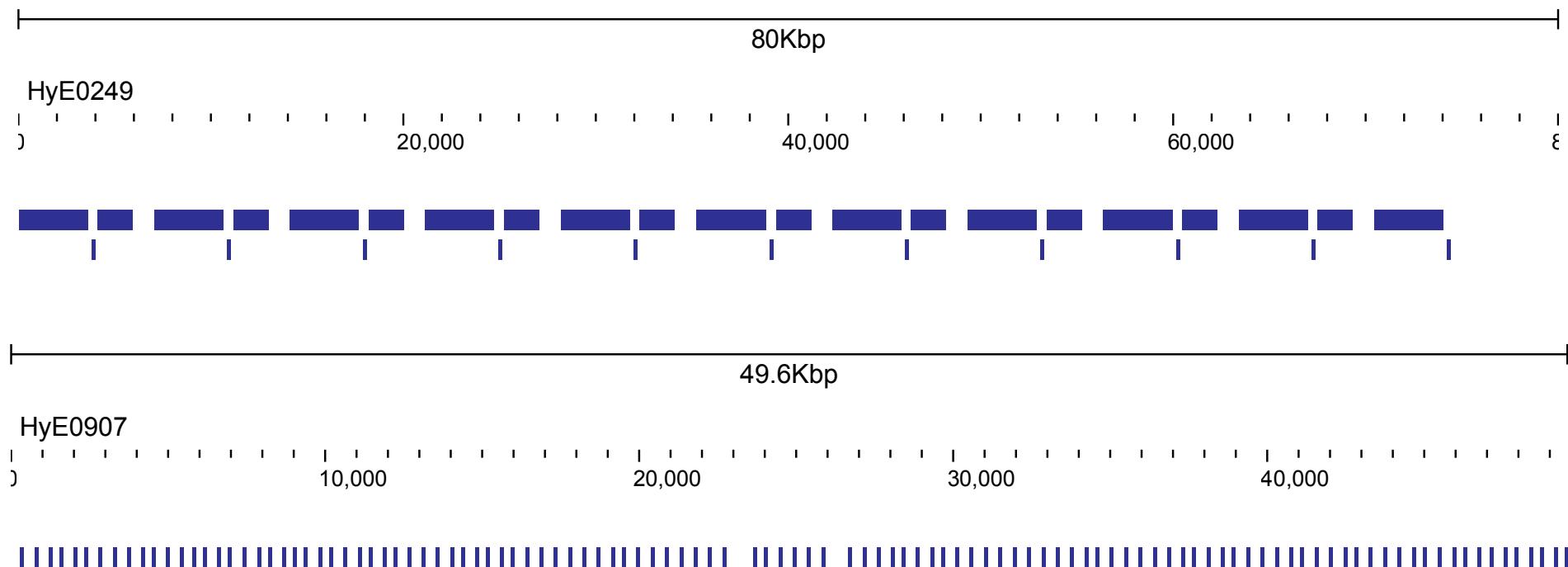
OPEN

R2DT is a framework for predicting and visualising RNA secondary structure using templates

Blake A. Sweeney^{1,7}, David Hoksza^{ID 2,7}, Eric P. Nawrocki³, Carlos Eduardo Ribas^{ID 1}, Fábio Madeira^{ID 1}, Jamie J. Cannone⁴, Robin Gutell⁴, Aparna Maddala⁵, Caeden D. Meade⁵, Loren Dean Williams^{ID 5}, Anton S. Petrov^{ID 5}, Patricia P. Chan^{ID 6}, Todd M. Lowe⁶, Robert D. Finn^{ID 1,8} & Anton I. Petrov^{ID 1,8✉}

The genome of the colonial hydroid *Hydractinia* reveals that their stem cells use a toolkit of evolutionarily shared genes with all animals

Christine E. Schnitzler,^{1,2} E. Sally Chang,^{3,4} Justin Waletich,^{1,2}
Gonzalo Quiroga-Artigas,^{1,2,5} Wai Yee Wong,⁶ Anh-Dao Nguyen,³ Sofia N. Barreira,³
Liam B. Doonan,⁷ Paul Gonzalez,³ Sergey Koren,³ James M. Gahan,^{7,8}
Steven M. Sanders,^{9,10} Brian Bradshaw,⁷ Timothy Q. DuBuc,^{7,11} Febrimarsa,^{7,12}
Danielle de Jong,^{1,2} Eric P. Nawrocki,⁴ Alexandra Larson,¹ Samantha Klasfeld,³
Sebastian G. Gornik,^{7,13} R. Travis Moreland,³ Tyra G. Wolfsberg,³ Adam M. Phillippy,³
James C. Mullikin,^{3,14} Oleg Simakov,⁶ Pauly Cartwright,¹⁵ Matthew Nicotra,^{9,10}
Uri Frank,⁷ and Andreas D. Baxevanis³



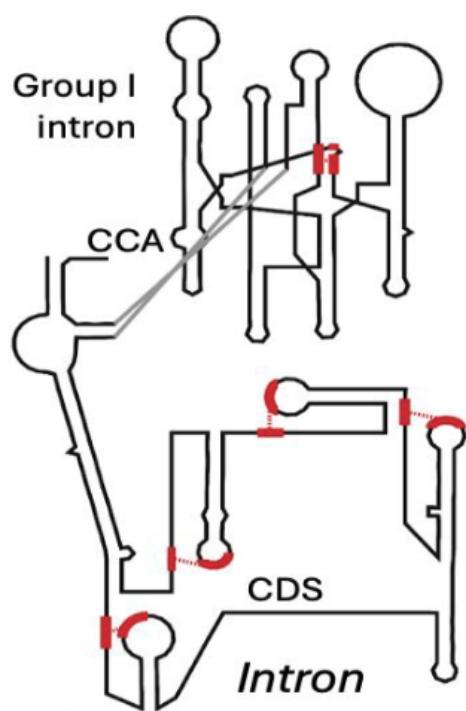
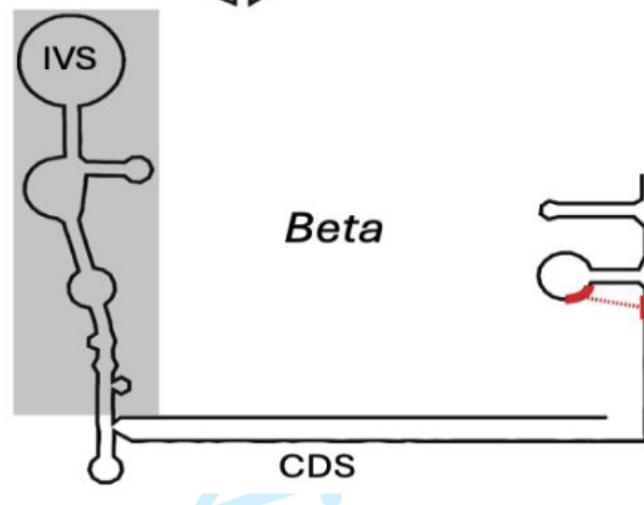
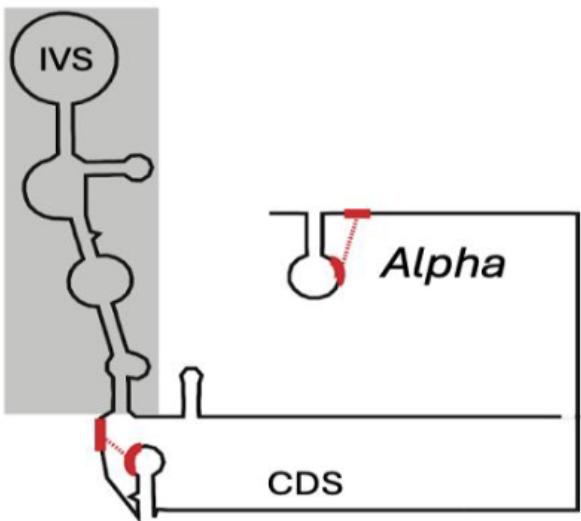
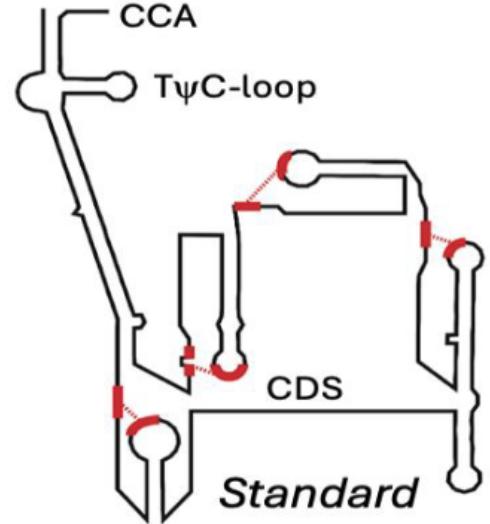
Expansion of the tmRNA sequence database and new tools for search and visualization

Authors

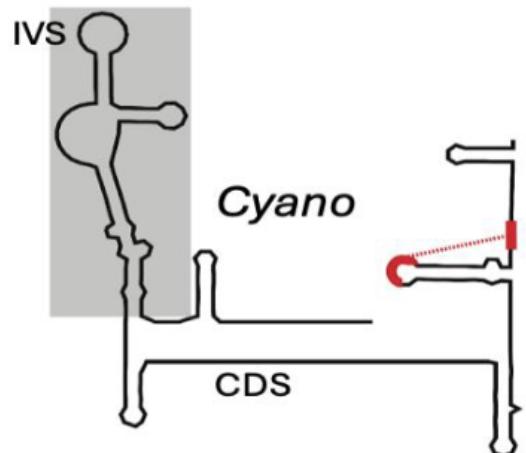
Eric P. Nawrocki¹, Anton I. Petrov², Kelly P. Williams³

Affiliations

1. Division of Intramural Research, U.S. National Library of Medicine, National Institutes of Health, Bethesda, MD, 20894, USA
0000-0002-2497-3427
eric.nawrocki@nih.gov
2. Riboscope Ltd, 23 King Street, Cambridge, CB1 1AH, UK
0000-0001-7279-2682
apetrov@riboscope.com
3. Sandia National Laboratories, Livermore CA 94550, USA
0000-0002-2606-9562
kwilli@sandia.gov



Model	Length	Bps	Stems	Pknots
<i>Standard</i>	358	117	12	4
<i>Intron</i>	616	180	22	5
<i>Alpha</i>	355	55	8	2
<i>Beta</i>	331	55	7	1
<i>Cyano</i>	288	72	9	1
<i>Mito</i>	77	19	3	0



Future directions for structural RNA research

- Further development of Infernal
 - iterative search
 - meta-models for clade-specific scoring
- Structural RNA annotation
 - group I introns: improved covariance models for Rfam
 - viral structural RNA annotations by VADR

Acknowledgements

NLM - VADR

Alejandro Schäffer
Rodney Brister
Ilene Mizrachi
Eneida Hatcher
Linda Yankie
Vince Calhoun
Susan Schafer
EB Dickinson

NLM - Ribovore

Alejandro Schäffer
Ilene Mizrachi
Rich McVeigh
Anji Johnston
Beverly Underwood
Alex Kotliarov
Barbara Robbertse
Conrad Schoch

NLM - RNA annotation

Françoise Thibaud-Nissen
Azat Badretdin
Terence Murphy
Michael DiCuccio
Tatiana Tatusova
Mark Borodovsky

Harvard/Janelia

Sean Eddy
Tom Jones
Diana Kolbe
Travis Wheeler
Elena Rivas
Michael Farrar

Rfam/RNACentral

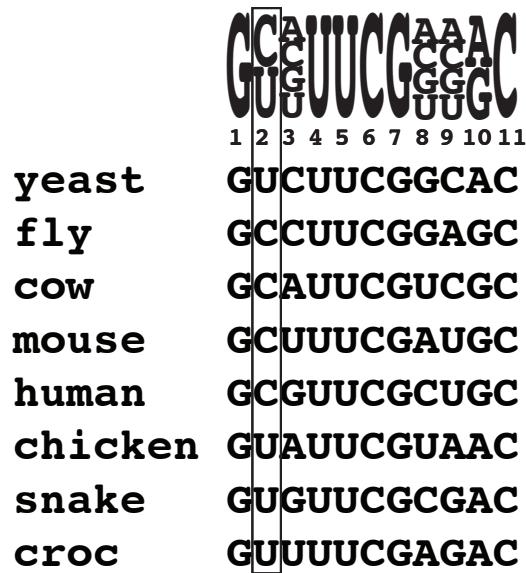
Anton Petrov
Blake Sweeney
Nancy Ontiveros
Kelly Williams
Sam Griffith-Jones
Paul Gardner

NLM - leadership

David Landsman
Richard Scheuermann
Steve Sherry
Kim Pruitt
Jim Ostell
David Lipman



Profile HMMs: sequence family models built from alignments



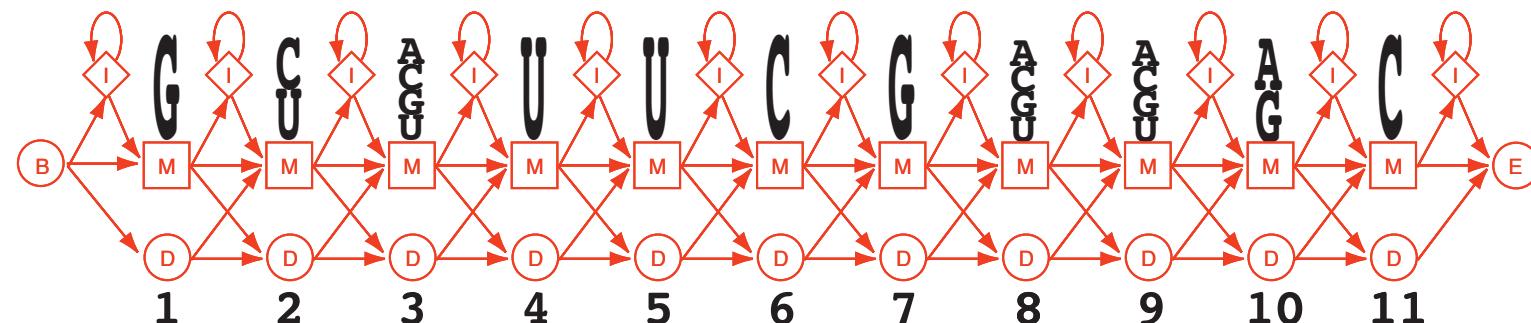
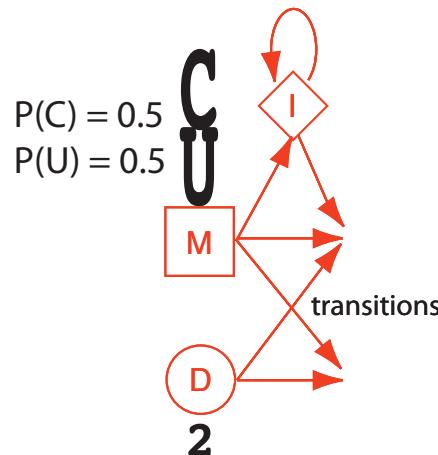
One HMM node per alignment column

3 states per node:

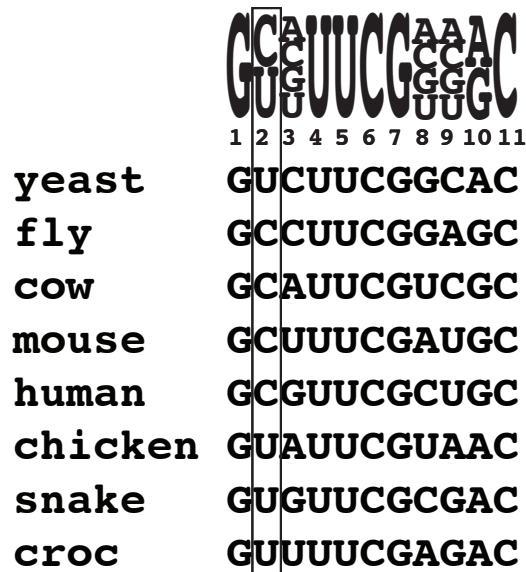
- (M) Match: emits residues
- (I) Insert: inserts extra residues
- (D) Delete: deletes residues

HMMs generate homologous sequences.

Node for column 2:



Profile HMMs: sequence family models built from alignments

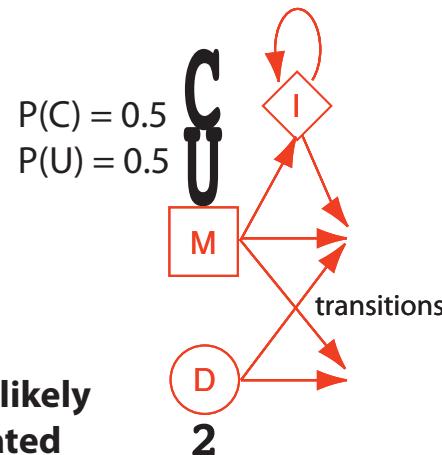


One HMM node per alignment column

3 states per node:

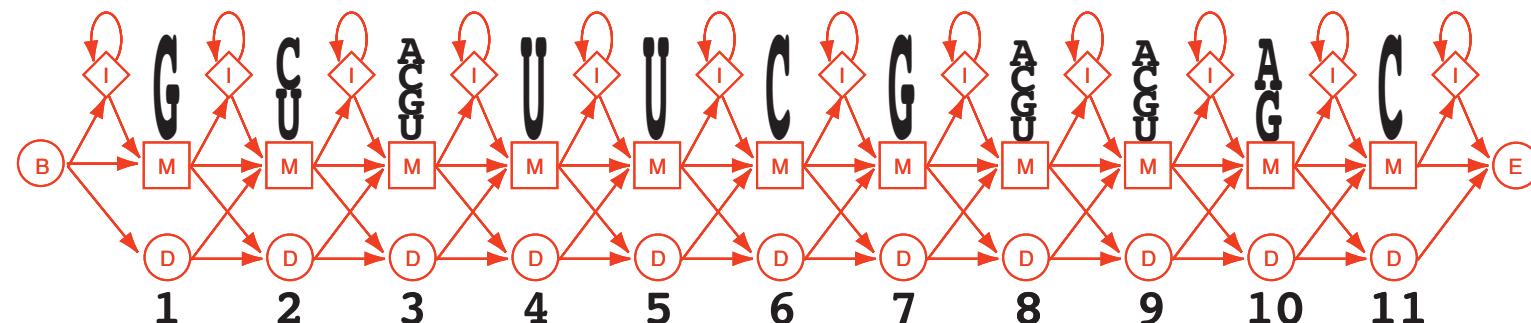
- (M) Match: emits residues
- (I) Insert: inserts extra residues
- (D) Delete: deletes residues

Node for column 2:



HMMs generate homologous sequences.

Given a sequence, the most likely path that could have generated that sequence can be computed.



Profile HMMs: sequence family models built from alignments

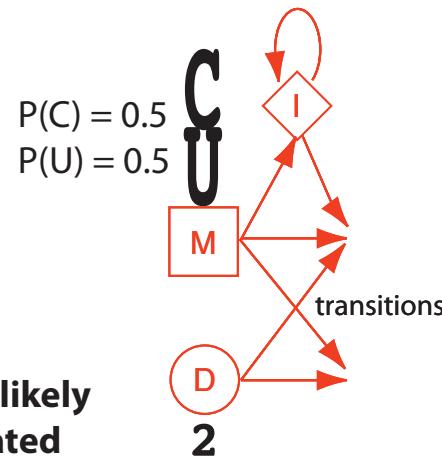
	yeast	GCA GUUUCGGAC 1 2 3 4 5 6 7 8 9 10 11
	fly	GCCUUUCGGAGC
	cow	GCAUUCGUCGC
	mouse	GCUUUCGAUGC
	human	GCGUUCGCUGC
	chicken	GUAUUCGUAAC
	snake	GUGUUCGCGAC
	croc	GUUUUCGAGAC
	worm	GCGUUCGCGGC

One HMM node per alignment column

3 states per node:

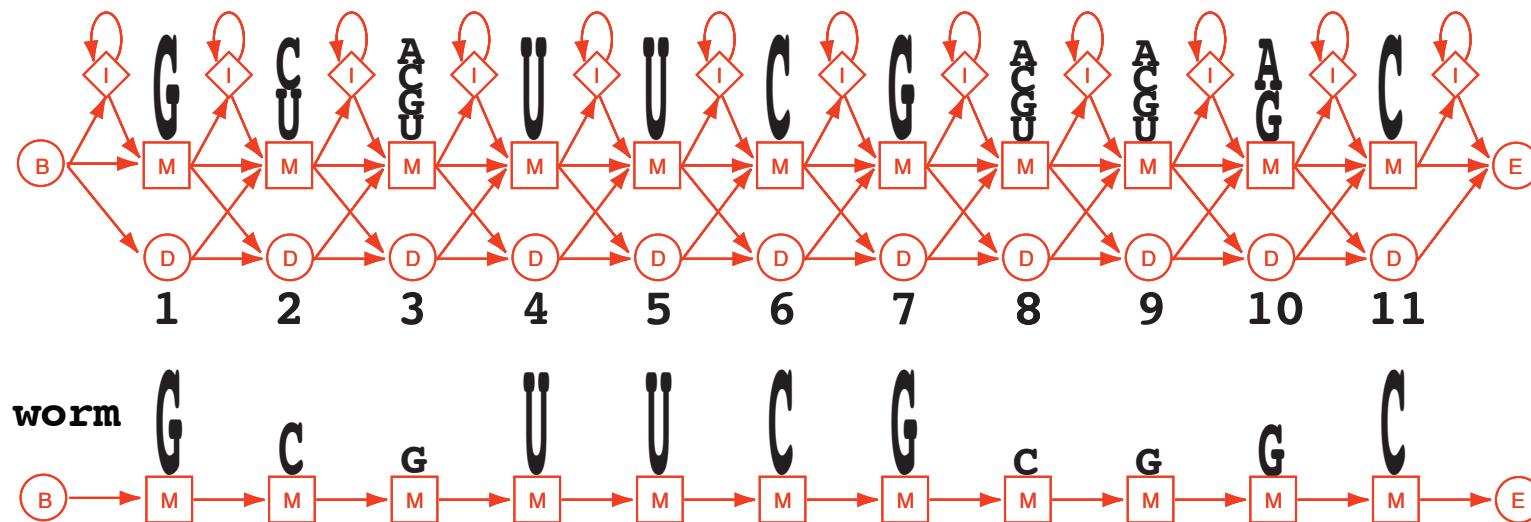
- (M) Match: emits residues
- (I) Insert: inserts extra residues
- (D) Delete: deletes residues

Node for column 2:

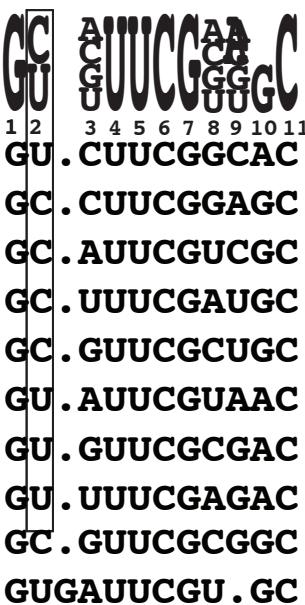


HMMs generate homologous sequences.

Given a sequence, the most likely path that could have generated that sequence can be computed.



Profile HMMs: sequence family models built from alignments

	
yeast	GU. CUUCGGCAC
fly	GC. CUUCGGAGC
cow	GC. AUUCGUCGC
mouse	GC. UUUCGAUGC
human	GC. GUUCGCUGC
chicken	GU. AUUCGUAAC
snake	GU. GUUCGCGAC
croc	GU. UUUCGAGAC
worm	GC. GUUCGCGGC
corn	GUGAUUCGU. GC

One HMM node per alignment column

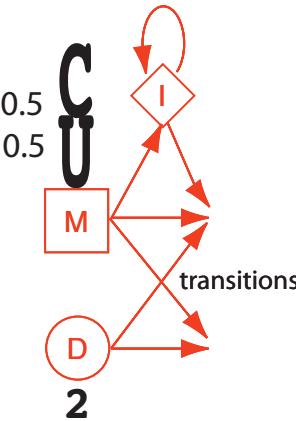
3 states per node:

- (M) Match: emits residues
- (I) Insert: inserts extra residues
- (D) Delete: deletes residues

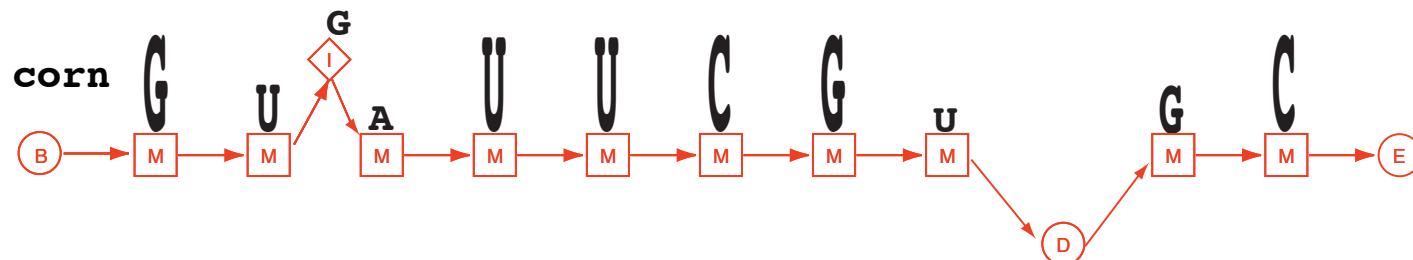
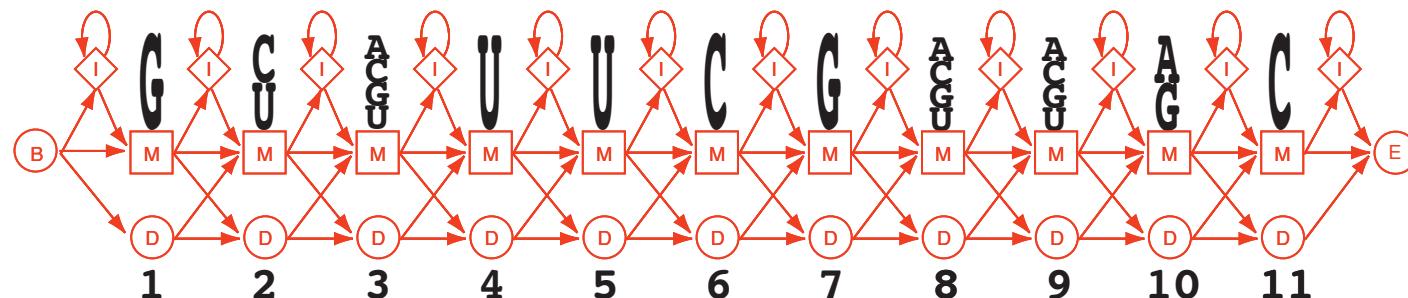
HMMs generate homologous sequences.

Given a sequence, the most likely path that could have generated that sequence can be computed.

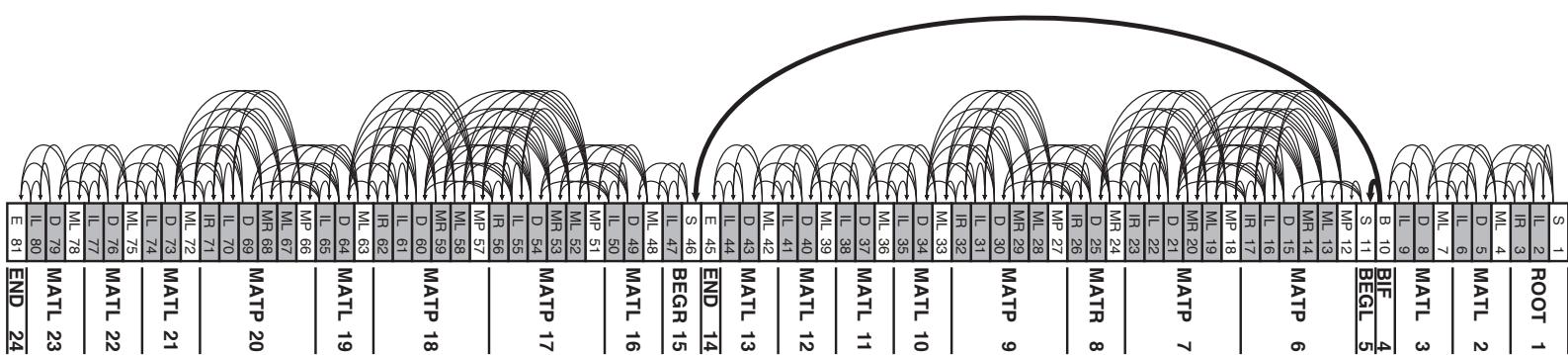
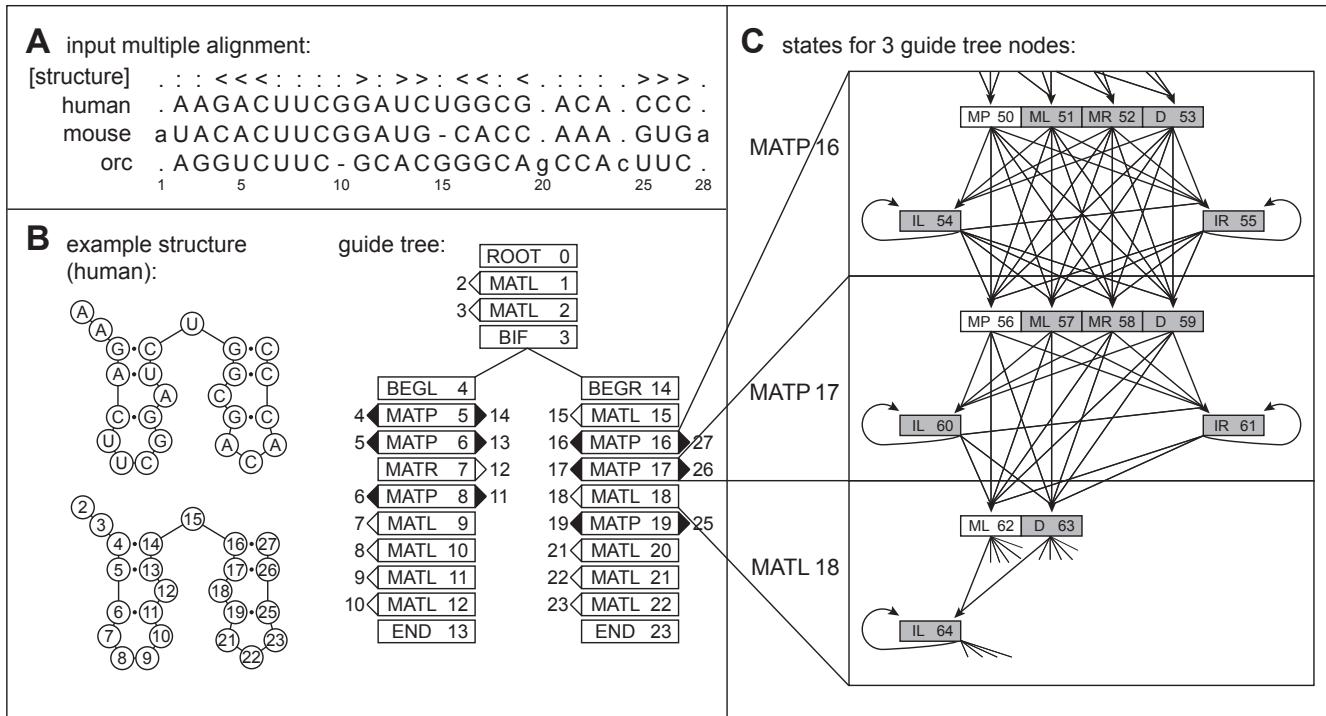
Node for column 2:



$$\begin{aligned} P(C) &= 0.5 \\ P(U) &= 0.5 \end{aligned}$$



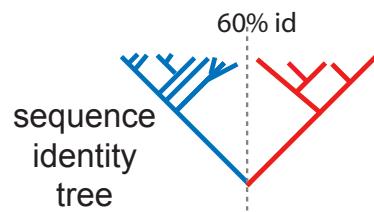
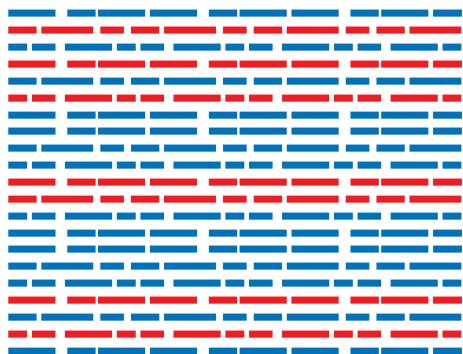
Covariance models (CMs) are built from structure-annotated alignments



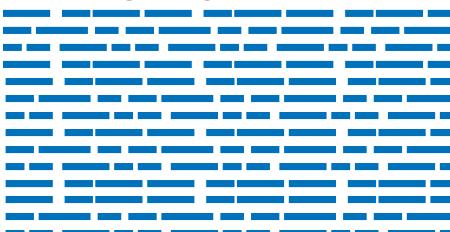
Is the added complexity worth it?

RMARK: a challenging internal RNA homology search benchmark

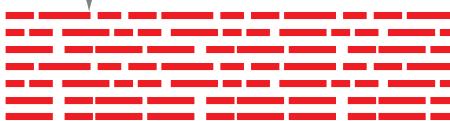
Rfam seed alignment:



training alignment

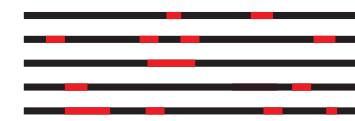


no train/test sequence pair is > 60% identical



test sequences

embed in
pseudo-genome

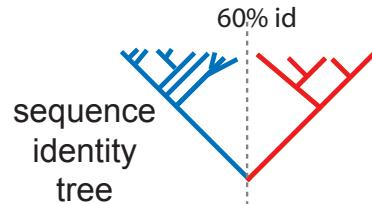
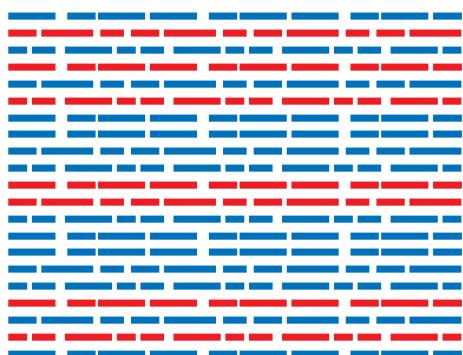


10 1Mb sequences
with 780 embedded
test seqs from 106 families

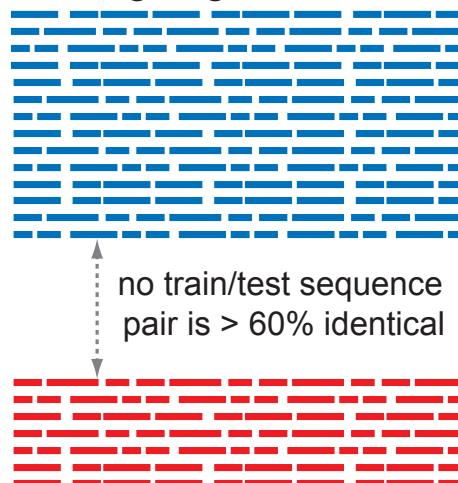
Is the added complexity worth it?

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



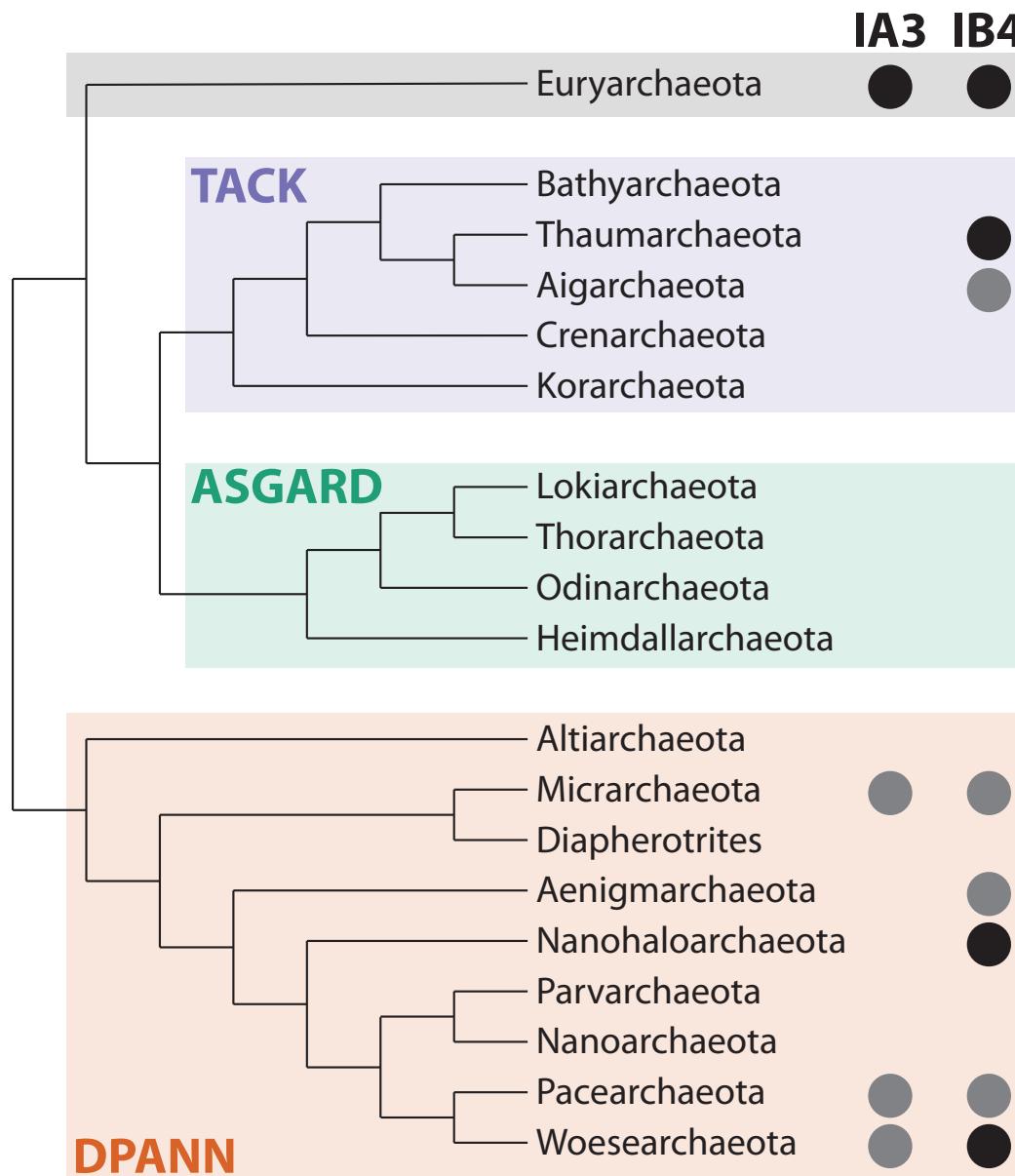
10 1Mb sequences
with 780 embedded
test seqs from 106 families

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 -

Group I introns are widespread in Archaea



Could archaeal group I introns have evolved into BHB introns?

Evolution of introns in the archaeal world

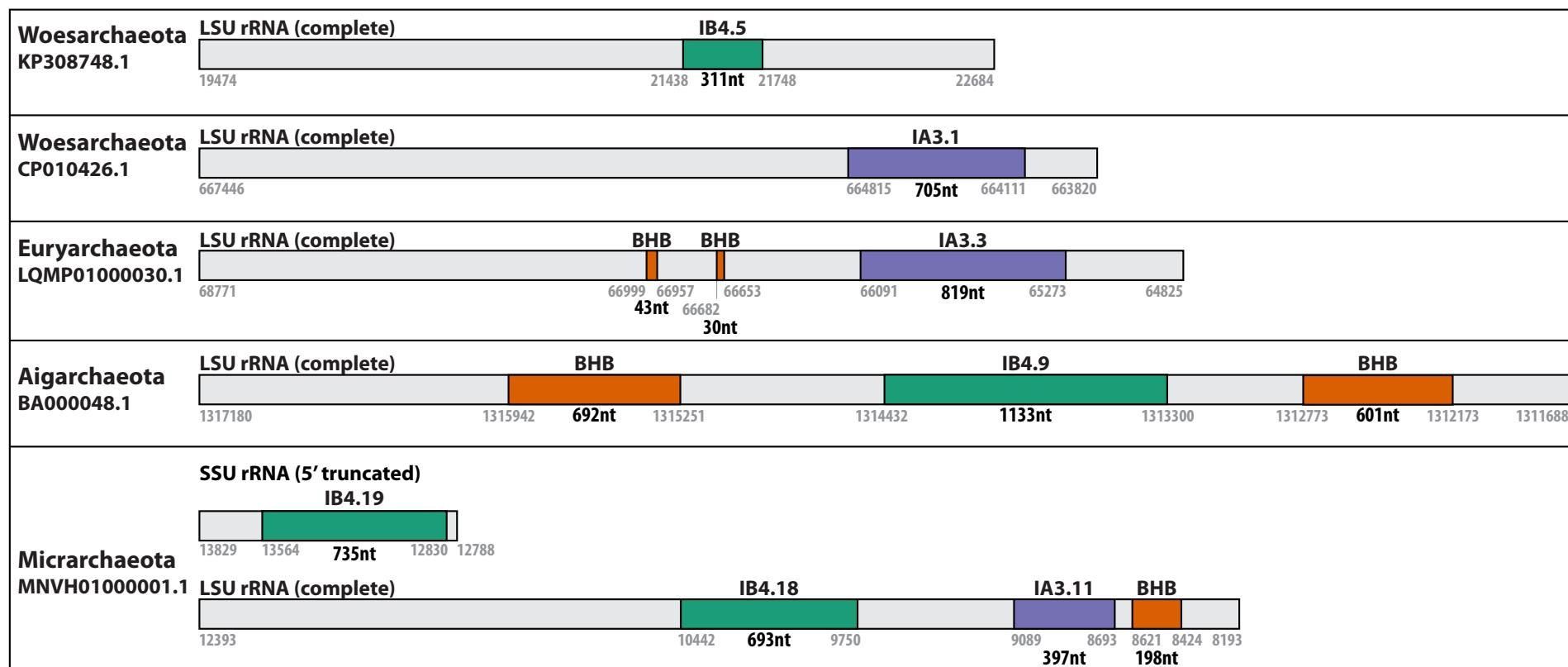
Giuseppe D. Tocchini-Valentini, Paolo Fruscoloni, and Glaucio P. Tocchini-Valentini¹

Istituto di Biologia Cellulare, Consiglio Nazionale delle Ricerche, Campus A, Buzzati-Traverso, Via Ramarini 32, Monterotondo Scalo, 00016 Rome, Italy

Contributed by Glaucio P. Tocchini-Valentini, January 24, 2011 (sent for review December 1, 2010)

*

Archaeal group I introns can occur in same host gene as BHB introns



Levels of sequence and structure conservation in RNA families

