

# RNAlien - Unsupervised RNA family model construction - Supplement

Florian Eggenhofer<sup>1,2</sup>, Ivo L. Hofacker<sup>1,3</sup> and Christian Höner zu  
Siederdisen<sup>4,1,5</sup>,

<sup>1</sup> Institute for Theoretical Chemistry, University of Vienna, Währingerstrasse  
17,A-1090 Vienna, Austria

<sup>2</sup> Bioinformatics Group, Department of Computer Science University of Freiburg,  
Georges-Köhler-Allee ,79110 Freiburg, Germany

<sup>3</sup> Bioinformatics and Computational Biology research group, University of Vienna,  
Währingerstrasse 17,A-1090 Vienna, Austria

<sup>4</sup> Bioinformatics Group, Department of Computer Science, University of Leipzig,  
D-04107 Leipzig

<sup>5</sup> Interdisciplinary Center for Bioinformatics, University of Leipzig, Härtelstraße  
16-18, D-04107 Leipzig, Germany

# Table of Contents

RNAlien - Unsupervised RNA family model construction - Supplement . . .	i
<i>Florian Eggenhofer, Ivo L. Hofacker and Christian Höner zu Siederdisen</i>	
A RNAlien detailed flowchart . . . . .	iii
B Implementation Details . . . . .	iv
B.1 Initial model construction . . . . .	iv
Search: . . . . .	iv
Filtering hits: . . . . .	iv
Model Construction: . . . . .	vi
Select Queries: . . . . .	vi
B.2 Model expansion . . . . .	vii
Search: . . . . .	vii
Filtering hits: . . . . .	vii
Model Construction: . . . . .	viii
Select Queries: . . . . .	ix
B.3 Model finalization: . . . . .	ix
Search, Filter: . . . . .	x
Reevaluation of potential candidates: . . . . .	x
Modelconstruction: . . . . .	x
B.4 Model evaluation . . . . .	x
B.5 Blast hit extension . . . . .	xi
C Rfam RNA families with known structure . . . . .	xii
D Diverse Rfam RNA families benchmark set . . . . .	xvi
E Negative control set . . . . .	xxiii
E.1 Random sequences . . . . .	xxiii
E.2 Ancestral repeats . . . . .	xxiii
E.3 Coding sequences . . . . .	xxiii
E.4 UTR regions . . . . .	xxiii

### A RNAlien detailed flowchart

Detailed flowchart representation of the RNAlien program flow.

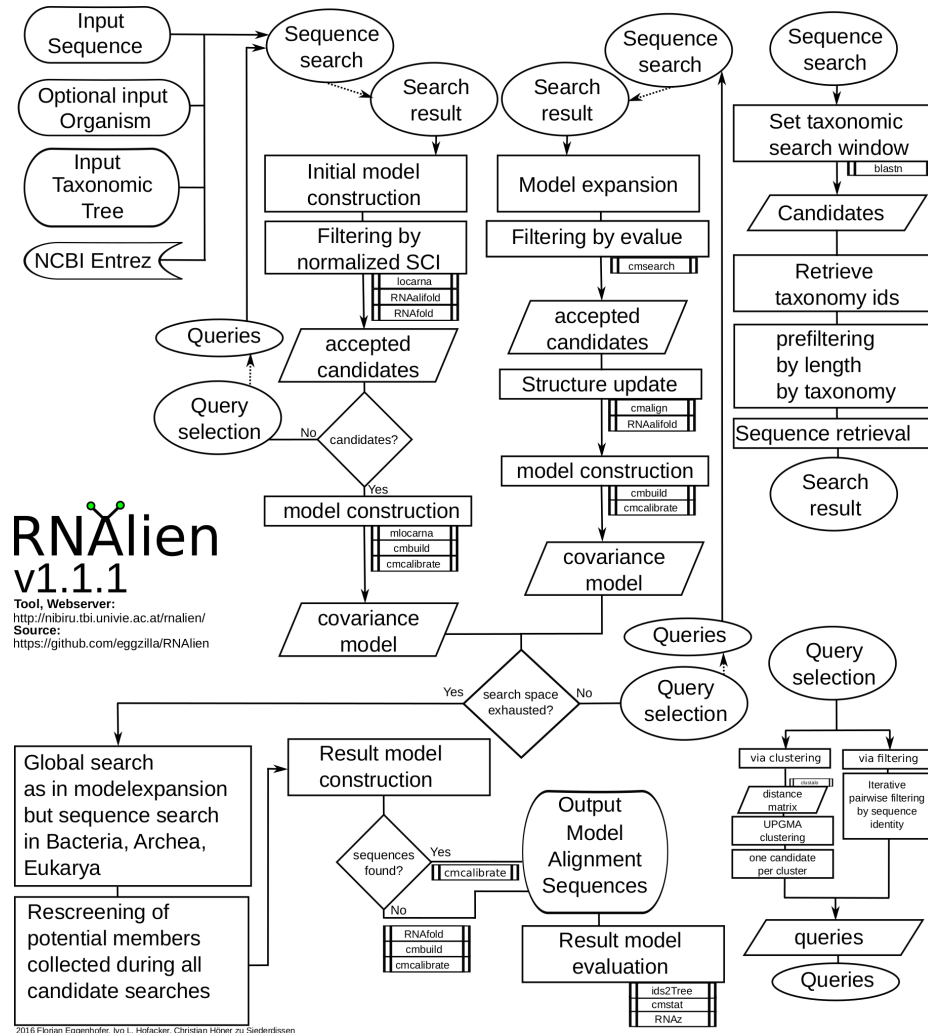


Fig. 1. Detailed program flow chart for RNAlien program.

## B Implementation Details

`RNAlieN` depends on several external tools and interfaces, which are listed in this section. System and function calls are included with their parameters and highlighted in *italic*. A starting point in the taxonomic tree is set, either specified by the input NCBI taxonomy id, or by running a nucleotide `BLAST` search via the NCBI REST interface and selecting the organism of the best hit. The model construction process starts at this organism and performs an initial model construction step. `RNAlieN` retrieves the taxonomic lineage of starting organism from the NCBI ENTREZ REST interface. After each of these steps `RNAlieN` proceeds to the taxonomic parent of the current taxonomic node. If a model was already constructed in a previous step then a model expansion step is performed, otherwise an initial model construction is attempted. Once the root of the taxonomic tree has been reached model expansion stops and the model finalization step is performed.

### B.1 Initial model construction

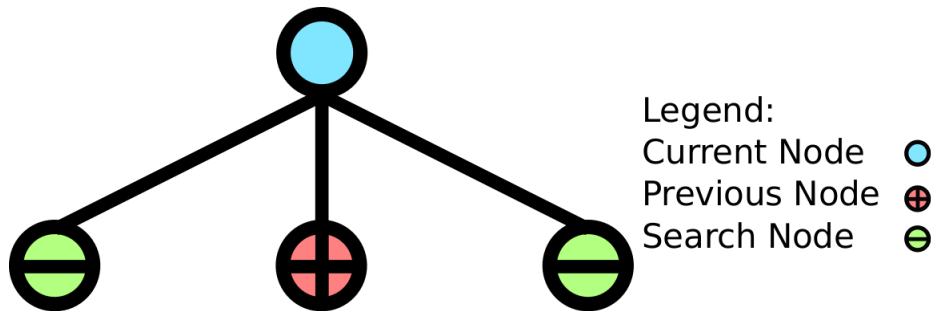
`RNAlieN` tries to establish an initial set of sequences related to the input sequence, that serve as seed for further expansion of the model.

**Search:** Candidate search is performed via the nucleotide `BLAST` REST interface which returns a list of hits. The organisms to be searched are restricted by the current taxonomic node of the step in two ways. To avoid overenrichment of sequences similar to included ones, already visited organisms are excluded. Only organisms that are associated with children of the current node are searched. For example if the current taxonomic position is *Enterobacteriaceae* and the previous node was *Enterobacter* all other organisms that belong to *Enterobacteriaceae* excluding *Enterobacter* are searched.

Summary of NCBI `BLAST` REST function call (one for each query, all other parameters default):

<b>blastHTTP</b>	
<i>parameter</i>	<i>value</i>
program	blastn
database	nt
querySequence	currentsequence
hitlistSize	5000
e-value	0.001
uppertaxonomylimit	currenttaxonomyid
lowertaxonomylimit	previoustaxonomyid

**Filtering hits:** `BLAST` hits are filtered by consecutively by following criteria: Hit has to achieve over 80% coverage of the query sequence.



**Fig. 2.** Organisms used for candidate search are determined as follows. All organisms and their corresponding genomes that are associated with the currently selected position in the taxonomic tree are used for searching. Excepted from this are organisms that have already been searched in previous rounds.

The hit must not exceed query length by factor of three.

Similarity of the hit to the query must be under 99%.

The remaining hits are expanded to query length as explained in subsection B.5. The gene id contained in the **BLAST** result and the expanded coordinates are used to retrieve nucleotide sequence from the Entrez REST interface. The sequences are filtered by normalized structure conservation index (nSCI). To compute the nSCI for each candidate sequence we need the minimum free secondary structure folding energy (MFE) of the candidate and the input sequence which is computed with RNAfold.

<b>RNAfold</b>	
<i>parameter</i>	<i>value</i>
--noPS	
inputfilePath	fastaFilePath
outputFilepath	foldFilePath

Furthermore the structure conservation index and the sequence identity of the candidate and input sequence are required. The candidate sequence is pair-wise aligned with free end-gap setting (semi-globally) to the input sequence. For each of these alignments the structure conservation index SCI is computed via RNAalifold.

<b>locarna</b>	
<i>parameter</i>	<i>value</i>
--write-structure	
--free-endgaps=+ + --	
--clustal	clustalFormatFilePath
inputFilepath1	inputFastaFilePath
inputFilepath2	inputFastaFilePath
outputFilepath	locarnaFilePath

<b>RNAalifold</b>	
<i>parameter</i>	<i>value</i>
inputFilepath	clustalFormatFilePath
outputFilepath	aliFoldFilePath

The sequence identity is computed via levenstein distance with following edit costs (delete,insert,substitution,transposition)=1. Candidate sequences are accepted for model construction if their nSCI exceeds one.

**Model Construction:** Candidate sequences that passed the nSCI filter are then used to build the initial model together with the input sequence. The sequences are structually aligned with mlocarna.

<b>mlocarna</b>	
<i>parameter</i>	<i>value</i>
inputFilepath	inputFastaFilePath
outputFilepath	mlocarnaFilePath

`cmbuild` is applied to the resulting structural stockholm alignment to construct a covariance model.

<b>cmbuild</b>	
<i>parameter</i>	<i>value</i>
--refine	
inputModelFilepath	cmFilePath
inputAlignmentFilepath	stockholmAlignmentFilePath
outputLogFilepath	logFilePath

The covariance model is used in the model expansion rounds to filter candidates and is therefore calibrated with `cmcalibrate`. This step is very time-consuming but sped up by using nonstandard (`--beta 10-4`) parameter. This affects the pre-filter steps of `cmsearch`, but not the final step where the sequence is aligned to the model via the CYK algorithm. Meaning that this increase in calibration speed reduces sensitivity but not specificity.

<b>cmcalibrate</b>	
<i>parameter</i>	<i>value</i>
--beta 1E-4	
inputModelFilepath	cmFilePath
outputFilepath	mlocarnaFilePath

**Select Queries:** At the end of the round queries for the candidate search of the next round are selected. `RNAlien` features a filtering and a clustering based method of query selection.

*Filtering based method* is the default method and iteratively removes all entries from the list of collected sequences, that do not have at most 95% pairwise sequence identity. This method has less specificity and sensitivity in the benchmarks (see 5, 6), but it is faster and removes the dependency on **clustalo**.

*Clustering based method* can alternatively be used by supplying the `-m` commandline switch with the value `clustering` to **RNAlien**. Clustal omega is used to compute a pairwise distance matrix of all collected sequences for clustering.

<b>clustalo</b>	
<i>parameter</i>	<i>value</i>
	<code>--full</code>
<code>--distmat-out</code>	<code>matrixFilePath</code>
<code>--infile</code>	<code>fastaFilePath</code>
<code>outputFilepath</code>	<code>clustaloFilePath</code>

**RNAlien** clusters the sequences via *unweighted pair group method with arithmetic mean* (UPGMA) and then incrementally increases the cutoff distance until 5 clusters can be formed. If less than 5 seqences have been collected, then each of them will be used as query.

## B.2 Model expansion

After a initial model has been constructed **RNAlien** enters into model expansion phase.

**Search:** Searching is performed as described in Initial model construction but with a relaxed e-value cutoff of 1 during the **BLAST** search.

<b>blastHTTP</b>	
<i>parameter</i>	<i>value</i>
<code>program</code>	<code>blastn</code>
<code>database</code>	<code>nt</code>
<code>querySequence</code>	<code>currentsequence</code>
<code>hitlistSize</code>	<code>5000</code>
<code>e-value</code>	<code>1</code>
<code>uppertaxonomylimit</code>	<code>currenttaxonomyid</code>
<code>lowertaxonomylimit</code>	<code>previoustaxonomyid</code>

**Filtering hits:** Filtering of **BLAST** hits and hit expansion is performed as described in Initial model construction.

Sequences are also retrieved via the NCBI Entrez REST interface but then filtered with a different approach. We use the calibrated covariance model of the previous round and apply it with `cmsearch` to the candidate sequences. Candidates are accepted into the growing model if their e-value is below 0.001 or as specified by the `inputValueCutoff` commandline argument.

To ensure a meaningful e-value cutoff we need to consider the size of the database. We reuse the size of the blast database the hit originates from.

The value is not by itself contained in the blast XML output, but all the parameters needed to compute it. The relationship of E-value and bitscore (Equation 3 adopted from [1]):

$$e = d * q * 2^{-b} \quad (1)$$

where  $d$  = databasesize

$e$  = e-value

$b$  = bitscore

$q$  = querylength

We compute the database size in Mbases that was used for the blast search as follows, by rearranging the equation above:

$$d = (e * 2^b) / q \quad (2)$$

where  $d$  = databasesize

$e$  = e-value

$b$  = bitscore

$q$  = querylength

Candidates are accepted into the growing model if their `cmsearch` E-value is below 0.001 or as specified by the `inputValueCutoff` commandline argument

<b>cmsearch</b>	
<i>parameter</i>	<i>value</i>
--notrunc	
-Z	databaseSize
-g	covarianceModelPath
inputFilepath	sequenceFilePath
outputFilepath	cmsearchFilePath

**Model Construction:** Candidates that were accepted by `cmsearch` and already collected sequences are structurally aligned with the covariance model of the previous round.

<b>cmalign</b>	
<i>parameter</i>	<i>value</i>
inputModelFilepath	cmFilePath
inputSequenceFilepath	fastaFilePath
outputAlignmentFilepath	stockholmAlignmentFilePath



As the secondary structure of the resulting stockholm alignment is not updated in this process, a consensus secondary structure of the new alignment is computed via RNAalifold, with settings specifically optimized to consider covariance contributions. The old consensus secondary structure is replaced with the new one in the alignment.

<b>RNAalifold</b>	
<i>parameter</i>	<i>value</i>
	-r
	--cfactor
-Z	databaseSize
-g	covarianceModelPath
inputFilepath	sequenceFilePath
outputFilepath	cmsearchFilePath

cmbuild is used to construct a updated covariance model.

<b>cmbuild</b>	
<i>parameter</i>	<i>value</i>
	--refine
inputModelFilepath	cmFilePath
inputAlignmentFilepath	stockholmAlignmentFilePath
outputLogFilepath	logFilePath

The model is calibrated with cmcalibrate for the following candidate search.

<b>cmcalibrate</b>	
<i>parameter</i>	<i>value</i>
	--beta 1E-4
inputModelFilepath	cmFilePath
outputFilepath	mlocarnaFilePath

**Select Queries:** Search candidates for the next round are selected as described in Initial model construction.

### B.3 Model finalization:

Model finalization serves to collect family members that could not be included in earlier rounds, because the model was too specific at that point and make the results available for the user. First individual candidate searches are performed in Archea, Bacteria, and Eukaria or as specified by the taxonomyRestriction commandline argument. The results are pooled and then processed as described in model expansion. The resulting model is then used to reevaluate collected potential candidates. These sequences are filtered as described in modelexpansion and if accepted included into the model. This final model is then calibrated with

default options to make it immediately useable for further homology search by the user.

**Search, Filter:** as in modelexpansion for 3 kingdoms (Archea - taxid 2157, Bacteria - taxid 2, Eukaria - taxid 2759)  
**Modelconstruction** as in Modelexpansion

**Reevaluation of potential candidates:** Filter like in Modelexpansion

**Modelconstruction:** as described above

<b>cmbuild</b>	
<i>parameter</i>	<i>value</i>
--refine	
inputModelFilepath	cmFilePath
inputAlignmentFilepath	stockholmAlignmentFilePath
outputLogFilepath	logFilePath

Calibration is done without speedup by --beta 1E-4 for the final model

<b>cmcalibrate</b>	
<i>parameter</i>	<i>value</i>
inputModelFilepath	cmFilePath
outputFilepath	mlocarnaFilePath

#### B.4 Model evaluation

In this step descriptors for the result files are computed. The covariance model is used as input for **cmstat**, which computes among other features the cm and hmm content of the model. **cmalign** is used to generate a **clustalw** format result alignment which is prefiltered by **rnazSelectSeqs.pl** (auxiliary script packaged with **RNAz**). This filtered alignment is used as input for **RNAz** set to use the decision model for structural alignments. The most relevant output of **RNAz** in this case is if it predicts the input to be structured RNA, which is a indicator for successful model constructions.

<b>cmalign</b>	
<i>parameter</i>	<i>value</i>
--outformat=Clustal	
inputModelFilepath	cmFilePath
outputFilepath	mlocarnaFilePath

<b>rnazSelectSeqs.pl</b>	
<i>parameter</i>	<i>value</i>
inputFilePath	clustalFilePath
outputFilePath	selectedClustalFilePath

<b>RNAz</b>	
<i>parameter</i>	<i>value</i>
	-1
inputFilePath	selectedClustalFilePath
outputFilePath	rnazFilePath

<b>cmstat</b>	
<i>parameter</i>	<i>value</i>
	-1
inputFilePath	covarianceModelPath
outputFilePath	cmstatFilePath

## B.5 Blast hit extension

RNALien expands found BLAST hits to the query length if possible.

*Same strand* BLAST hit are extended as follows,

$$\begin{aligned}
 t &= h - q \\
 T &= H + (L - Q) \\
 s(t) &= \begin{cases} t, & \text{if } t \geq 0 \\ 0, & \text{otherwise} \end{cases} \\
 E(T) &= \begin{cases} b, & \text{if } T \geq b \\ T, & \text{otherwise} \end{cases}
 \end{aligned}$$

where  $h$  is the start coordinate of the hit,  $t$  is the extended start coordinate,  $q$  is the start coordinate of the hit on the query,  $H$  is the end coordinate of the hit,  $T$  is the extended endcoordinate,  $Q$  is the end coordinate of the hit on the query,  $L$  is the length of the query sequence  $b$  is the length of the sequence the hit maps to  $s$  is the start coordinate of the extended sequence checked for being within the available coordinates of the hit sequence,  $E$  is the end coordinate of the extended sequence checked for being within the available coordinates of the hit sequence



**Fig. 3.** Extension of BLAST hit and query on the same strand to query length, where  $h$  is the start coordinate of the hit,  $t$  is the extended start coordinate,  $q$  is the start coordinate of the hit on the query,  $H$  is the end coordinate of the hit,  $T$  is the extended endcoordinate,  $Q$  is the end coordinate of the hit on the query,  $L$  is the length of the query sequence  $b$  is the length of the sequence the hit maps to  $s$  is the start coordinate of the extended sequence checked for being within the available coordinates of the hit sequence,  $E$  is the end coordinate of the extended sequence checked for being within the available coordinates of the hit sequence

*Different Strand* BLAST hit are extended as follows,

$$t = h + q$$

$$T = H - (L - Q)$$

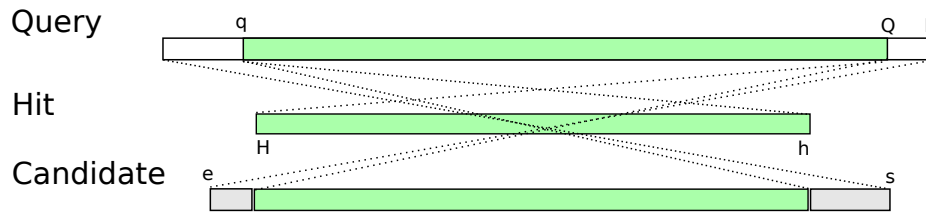
$$s(t) = \begin{cases} b, & \text{if } t \geq b \\ t, & \text{otherwise} \end{cases}$$

$$e(T) = \begin{cases} T, & \text{if } T \geq 0 \\ 0, & \text{otherwise} \end{cases}$$

where  $h$  is the start coordinate of the hit,  
 $t$  is the extended start coordinate,  
 $q$  is the start coordinate of the hit on the query,  
 $H$  is the end coordinate of the hit,  
 $T$  is the extended endcoordinate,  
 $Q$  is the end coordinate of the hit on the query,  
 $L$  is the length of the query sequence  
 $b$  is the length of the sequence the hit maps to  
 $s$  is the start coordinate of the extended sequence checked for being within the available coordinates of the hit sequence,  
 $E$  is the end coordinate of the extended sequence checked for being within the available coordinates of the hit sequence

## C Rfam RNA families with known structure

This section contains additional plots for the RNA families with known structure featured in the paper. The first 2 plots show the changes of specificity and



**Fig. 4.** Extension of BLAST hit and query on different strands to query length

sensitivity after subsequently applying the suggestions of the reviewers. The original version before the review was RNAlien 1.0.0, the one including all changes listed here has version 1.1.1.

Inclusion of paralogs and toggling of the refine switch for cmbuild were included first, this has improved both specificity, as well as recall. Additionally to this, we changed the method for selecting queries for searching candidates from clustering all collected sequences and picking one sequence per cluster to filtering all sequence that do not have a pairwise sequence identity of less than 95%.

While the specificity is slightly only lower, there is a decrease in specificity. Nevertheless we have selected the new query selection method as default, because it is substantially faster and it drops the dependency on clustal-omega.

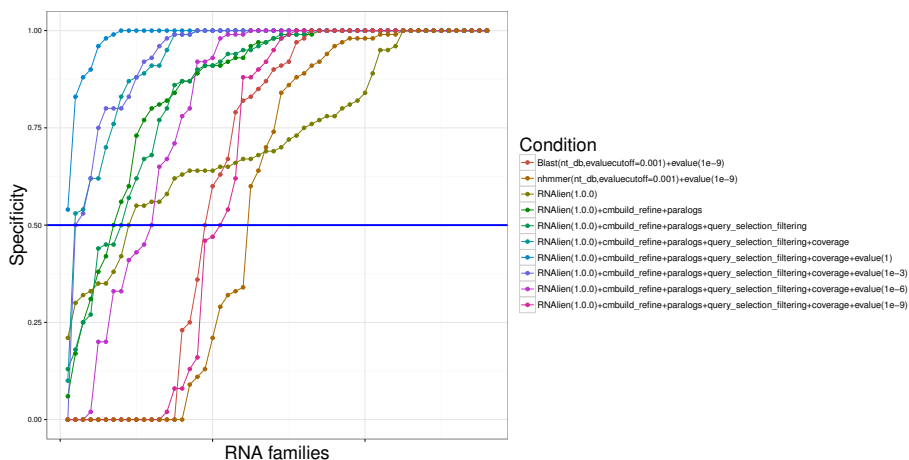
Blast hits are now also checked for the hit to have at least 80% coverage of the query. This feature should have been included in RNAlien 1.0.0, but was faulty.

Query sequences submitted to blast can be softmasked with conservation information from /cmalign. This feature is not considered in the shown benchmarks, but can be activated via commandline switch.

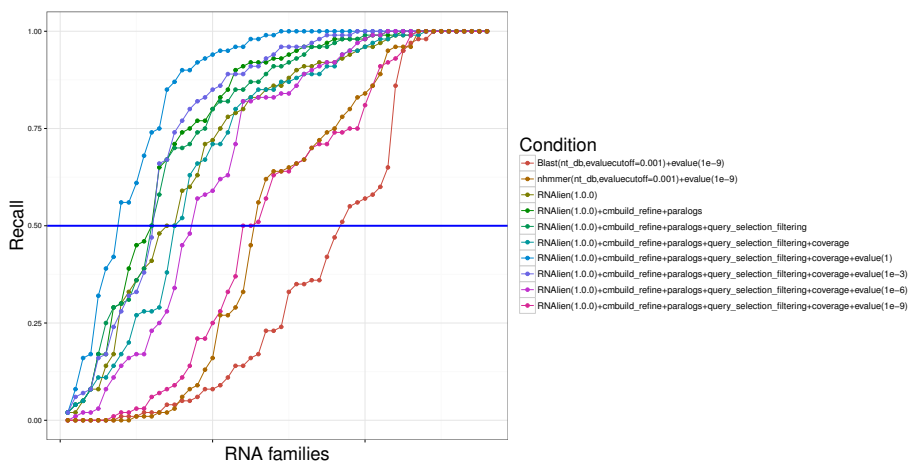
All of the newly introduced features can be controlled via commandline switches, with exception of the cmbuild refinement.

The runtime of RNAlien for the structured RNA test set

Following is the table of sequences from the Rfam 12.0 seed alignments of families with known structure, that was used in the result section. The first sequence of the family was picked with the exception of sequences that are associated with metagenomic tax ids that could not be processed by the NCBI REST BLAST interface.



**Fig. 5.** Specificity for 56 RNA families with known 3D structure.



**Fig. 6.** Recall for 56 RNA families with known 3D structure.

Table 1: RNA families with known structure benchmark table. Column names A to N are placeholders for the following names: Specificity\_Alien (=A) Sens\_Alien (=B) Spec+paralogs+refine (=C) Sens+paralogs+refine (=D) Spec+filterings (=E) Sens+filtering (=F) Spec+coverage (=G) Sens+coverage (=H) Spec\_value (=I) Sens\_value (=J) Spec\_nhmmer\_value (=K) Sens\_nhmmer\_value (=L) Spec\_blast\_value (=M) Sens\_blast\_value (=N). The column names annotated with *value* were computed with a *value* cutoff of  $1^{-3}$  and a databasesize of  $10^9$  bases per default, with the exception of families that can be found exclusively in prokaryotes and viruses.

Rfam id	Rfam name	A	B	C	D	E	F	G	H	I	J	K	L	M	N
5S_rRNA	RF00001	0.64	0.91	0.99	0.92	0.77	0.82	1	0.87	1	0.83	0.94	0.53	0.72	0.55

Continued on next page

Table 1 – continued from previous page

Rfam id	Rfam name	A	B	C	D	E	F	G	H	I	J	K	L	M	N
5.8S_rRNA	RF00002	0.95	0.9	1	0.9	1	0.85	1	0.85	1	0.89	0.96	0.95	1	0.74
U1	RF00003	0.58	0.88	0.97	1	0.99	1	1	1	1	0.99	0.86	0.99	1	0.75
U2	RF00004	0.62	0.99	0.89	0.98	0.99	0.99	0.99	0.95	0.99	0.96	0.84	0.99	0.96	0.89
tRNA	RF00005	0.77	0.48	1	0.75	1	0.7	1	0.63	0.75	0.47	0.9	0.15	1	0.04
Hammerhead_3	RF00008	1	0.63	1	0.74	1	0.74	1	0.74	1	0.74	1	0.74	1	0.74
RNaseP_bact_a	RF00010	0.56	1	0.93	0.98	0.94	0.98	1	1	1	1	0.98	1	1	1
RNaseP_bact_b	RF00011	0.55	1	0.99	1	1	1	1	1	1	1	0.99	1	0.91	1
Metazoa_SRP	RF00017	0.35	0.92	0.06	0.96	0.13	0.96	1	0.99	1	0.99	0.95	0.96	1	0.99
tmRNA	RF00023	0.65	0.92	0.98	0.93	0.98	0.92	0.99	0.88	0.99	0.91	0.98	0.95	0.98	0.61
U6	RF00026	0.64	0.83	0.98	0.83	0.99	0.82	1	0.82	1	0.8	0.88	0.89	0.93	0.71
Intron_gpI	RF00028	0.32	0.08	0.17	0.17	0.27	0.25	1	0.08	1	0.08	0.28	0.08	1	0.08
Intron_gpII	RF00029	0.75	0.41	0.31	0.65	0.92	0.58	0.89	0.2	0.92	0.16	0.84	0.2	1	0.09
Histone3	RF00032	1	0.02	1	0.02	1	0.02	1	0.02	0.5	0.02	1	0.02	1	0.02
IRE_I	RF00037	0.78	0.92	0.99	0.92	0.95	0.87	0.95	0.89	0.96	0.89	1	0.05	1	0.05
Phage_pRNA	RF00044	0.8	1	0.8	1	0.8	1	1	1	1	1	0.8	1	1	1
FMN	RF00050	0.21	0.79	0.6	1	0.18	1	0.1	0.66	1	1	1	1	0.85	1
TPP	RF00059	0.63	0.83	0.77	0.92	0.57	0.91	0.87	0.83	0.8	0.82	1	0.5	1	0.3
S15	RF00114	0.56	0.85	1	0.85	1	0.85	1	0.85	1	0.85	1	0.83	1	0.82
SAM	RF00162	0.55	0.72	0.81	0.99	0.67	0.98	0.76	0.91	0.98	1	0.99	0.99	1	0.79
s2m	RF00164	1	1	1	1	1	1	1	0.97	1	0.97	1	0.18	1	0.87
Purine	RF00167	0.66	0.71	1	1	1	1	1	0.99	1	1	1	0.99	1	0.94
Lysine	RF00168	0.82	0.94	1	0.98	1	0.89	1	0.85	1	0.94	1	0.98	1	0.45
Bacteria_small_SRP	RF00169	0.64	0.75	0.99	0.77	0.99	0.75	1	0.27	1	0.67	1	0.41	1	0.39
Cobalamin	RF00174	0.73	0.6	0.87	0.99	0.91	1	1	0.94	1	0.96	1	0.97	0.9	0.9
HIV-1_DIS	RF00175	1	0.99	1	0.99	1	0.99	1	0.99	1	1	1	0.91	1	0.91
SSU_rRNA_bacteria	RF00177	1	1	1	1	1	1	1	1	1	1	1	1	1	1
K10_TLS	RF00207	1	1	1	1	1	1	1	1	0.8	1	1	1	1	1
IRES_Pesti	RF00209	0.95	1	0.91	1	0.91	1	1	1	1	1	0.96	1	1	1
glmS	RF00234	0.69	1	0.97	1	1	1	1	0.89	1	0.89	1	1	1	0.89
Gammaretro_CES	RF00374	0.64	1	0.87	1	0.87	1	1	1	1	1	0.99	1	1	1
ykoK	RF00380	0.7	0.86	1	0.94	1	0.96	1	0.89	1	0.99	1	0.99	0.95	0.99
IRES_Cripavirus	RF00458	0.33	0.14	0.25	0.29	0.25	0.29	1	0.29	1	1	1	0.86	1	1
HIV_FE	RF00480	1	0.98	1	0.98	1	0.98	1	0.98	1	0.99	1	0.99	1	0.97
TCV_H5	RF00500	1	0.8	1	0.8	1	0.8	1	0.8	1	1	1	1	1	1
Glycine	RF00504	0.69	0.59	0.91	0.77	0.87	0.7	0.91	0.52	0.99	0.66	1	0.09	1	0.09
mir-228	RF00843	1	1	1	1	1	1	1	1	1	1	1	1	1	1
mir-689	RF00871	0.5	0.08	0.5	0.08	0.5	0.08	0.83	0.38	0.83	0.38	0.92	0.38	1	0.46
c-di-GMP-I	RF01051	0.81	0.97	1	0.97	0.94	0.94	1	0.96	1	0.98	1	0.76	1	0.69
preQ1-II	RF01054	0.67	0.93	1	0.93	1	0.93	1	0.71	1	0.93	1	0.86	1	0.71
GP_knot1	RF01073	0.96	0.86	0.96	0.71	0.96	0.71	0.91	0.71	0.93	0.86	1	0.43	0.88	0.71
PK-G12rRNA	RF01118	0.65	0.99	0.73	0.99	0.68	0.97	0.99	0.99	1	1	1	1	1	1
HIV-1_SD	RF01380	1	0.05	1	0.05	1	0.05	1	0.05	1	0.77	1	0	1	0

Continued on next page

Table 1 – continued from previous page

Rfam id	Rfam name	A	B	C	D	E	F	G	H	I	J	K	L	M	N
MFR	RF01510	1	0.33	1	0.67	1	0.67	1	0.67	1	1	1	1	0.67	1
AdoCbl-variant	RF01689	1	0.91	1	0.91	1	0.91	1	0.91	1	0.91	1	0.86	1	0.91
crcB	RF01734	0.84	0.36	0.93	0.45	1	0.36	0.88	0.28	0.88	0.32	1	0.03	1	0.03
c-di-GMP-II	RF01786	0.67	0.02	1	0.04	1	0.04	1	0.04	1	0.07	1	0.04	1	0.02
THF	RF01831	0.76	0.96	1	0.96	0.86	0.87	0.7	0.14	0.8	0.24	1	0.24	1	0.16
tRNA-Sec	RF01852	0.89	0.3	0.92	0.3	0.9	0.3	0.53	0.28	0.53	0.28	0.88	0.29	0.92	0.28
Protozoa_SRP	RF01856	0.72	0.39	0.91	0.39	0.95	0.39	1	0.11	1	0.33	1	0.33	1	0.33
Archaea_SRP	RF01857	0.35	0.96	0.82	0.96	0.45	0.96	1	0.87	1	0.96	1	1	1	0.15
group-II-D1D4-1	RF01998	0.38	0.5	0.38	0.46	0.45	0.31	0.62	0.11	0.62	0.06	0.83	0.02	0.5	0.02
group-II-D1D4-3	RF02001	0.3	0.78	1	0.98	1	0.99	1	0.98	1	0.96	0.94	0.94	0.97	0.49
mir-2985-2	RF02095	0.68	0.95	0.84	0.95	0.97	1	1	0.95	1	1	1	1	1	1
IRE_II	RF02253	0.42	0.17	0.42	0.17	0.44	0.17	0.54	0.17	0	0.17	0	0.03	0	0.17
ToxI	RF02519	0.78	0.5	0.56	0.5	0.62	0.5	0.62	0.5	1	1	1	1	1	1

## D Diverse Rfam RNA families benchmark set

The Rfam database features following tags to group families: Cis-reg, frameshift element, IRES, leader, riboswitch, thermoregulator, antisense, antitoxin, CRISPR, lncRNA, miRNA, ribozyme, rRNA, snRNA, snoRNA, CD-box, HACA-box, scaRNA, splicing, Gene, sRNA, tRNA, Intron.

To obtain a representative sample of Rfam families, for each of these tags the alphanumerically first 10 families (if available for that tag) were selected. As some families have multiple tags, the list was filtered to contain each family only once.

The benchmark was conducted in the same manner as for the families with known 3D structure. The plots show different combinations of e-value cutoffs and databasesizes. Without explicitly setting the database size cmsearch uses twice the sequence length (forward/backward strand).

The setting comparable to the one used for the structured dataset is diverse(ev-1e-3,db-1e-9), meaning a cmsearch e-value cutoff of 1e-3 and a databasesize of  $10^9$  bases in general and  $10^6$  bases for bacterial and viral RNA families.

The result with comparable settings to the structured dataset has 191 of 192 cases (99%) with at least half of the sequences collected by RNAlie are recognized as belonging to the Rfam model. In 170 (89%) families all sequences included by RNAlie are recognized as belonging to the Rfam model

In of 163 cases (85%) at least half the sequences in the Rfam seed alignment are correctly recognized by the RNAlie model. In 123 of 191 cases (64%) all sequences in the Rfam seed alignment are correctly recognized by the RNAlie model.



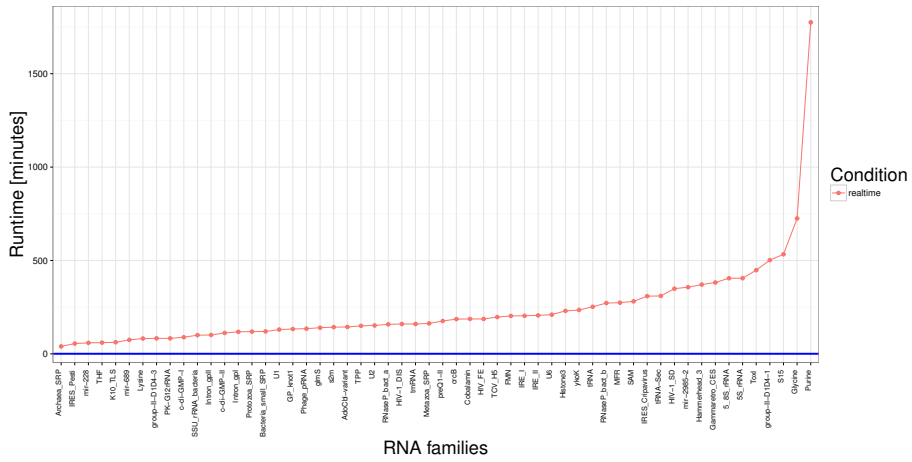


Fig. 7. Alien program runtime in minutes for structured families

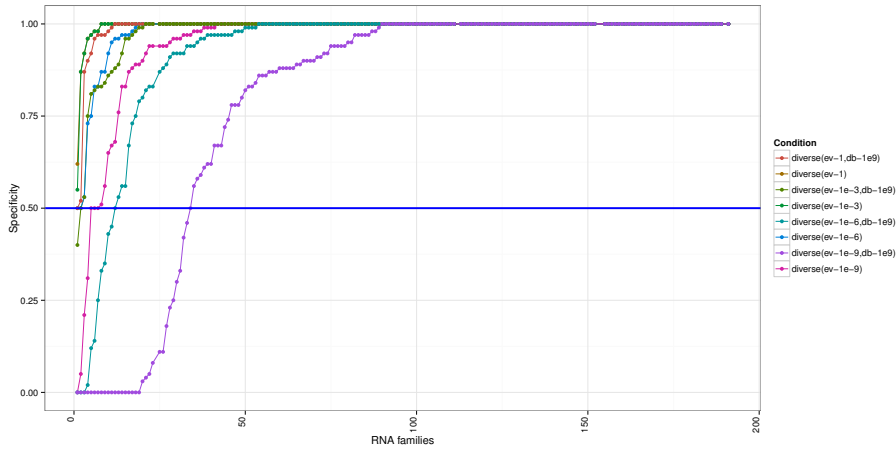
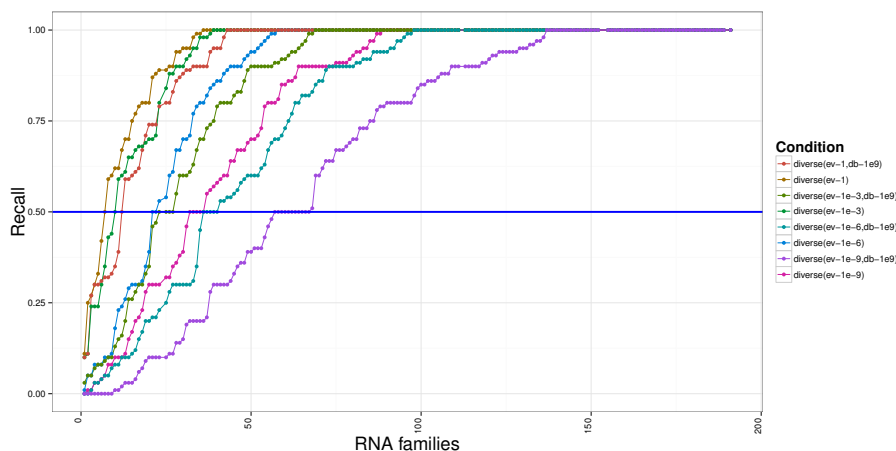


Fig. 8. Specificity of RNAlien homology search. The plot shows the fraction of homologs predicted by RNAlien that are recognized by the original Rfam model. The legend indicates the e-evalue cutoff (ev-) and the database size used. The e-evalue cutoffs start at 1 and are made stricter in 1e-3 steps up to 1e-9. The result with comparable settings to the structured dataset has 191 of 192 cases (99%) with at least half of the sequences collected by RNAlien are recognized as belonging to the Rfam model. In 170 (89%) families all sequences included by RNAlien are recognized as belonging to the Rfam model

Following is the table of families from the Rfam 12.0 used in the as a second benchmark set.



**Fig. 9.** Recall for 191 RNA families, selected up to 10 for each family tag. To test our method, sRNA Rfam family models were reconstructed by **RNALien** from a random sequence picked from the family seed sequences. This plot shows how many Rfam seed sequences are recognized by the reconstructed **RNALien** model using the model gathering score (used by Rfam to establish full models). In of 163 cases (85%) at least half the sequences in the Rfam seed alignment are correctly recognized by the **RNALien** model. In 123 of 191 cases (64%) all sequences in the Rfam seed alignment are correctly recognized by the **RNALien** model.

Table 2: Diverse RNA families benchmark set. Column names A to D are placeholders for following names: Specificity\_value\_1 (=A) Sensitivity\_value\_1 (=B) Specificity\_value\_1e-3 (=C) Sensitivity\_value\_1e-3 (=D) Specificity\_value\_1e-6 (=E) Sensitivity\_value\_1e-6 (=F) Specificity\_value\_1e-9 (=G) Sensitivity\_value\_1e-9 (=H)

Rfam name	Rfam id	A	B	C	D	E	F	G	H
5S_rRNA	RF00001	1	0.88	1	0.75	0.97	0.54	0.58	0.48
5_8S_rRNA	RF00002	1	0.89	1	0.82	1	0.75	1	0.69
U1	RF00003	1	1	1	1	0.98	0.98	0.94	0.94
U2	RF00004	1	1	0.99	0.96	0.92	0.83	0.88	0.75
tRNA	RF00005	1	0.61	0.75	0.47	0.02	0.25	0	0.03
RNaseP_nuc	RF00009	1	0.11	1	0.09	1	0.08	1	0.07
RNaseP_bact_b	RF00011	1	1	1	1	1	1	1	1
U4	RF00015	0.97	0.74	0.96	0.55	0.94	0.12	0.9	0.06
Y_RNA	RF00019	0.52	0.59	0.4	0.33	0.12	0.31	0.05	0.31
U5	RF00020	1	0.9	1	0.61	0.97	0.21	0.78	0.09
U6	RF00026	1	0.94	1	0.8	0.99	0.63	0.88	0.51
PrfA	RF00038	1	1	1	1	1	1	1	1
CopA	RF00042	1	1	1	1	1	1	1	0.97
FMN	RF00050	1	1	1	1	0.99	0.99	0.8	0.78

Continued on next page

Table 2 – continued from previous page

Rfam name	Rfam id	A	B	C	D	E	F	G	H
TPP	RF00059	1	0.95	1	0.83	1	0.53	0.67	0.39
U7	RF00066	1	0.71	0.99	0.63	0.79	0.53	0.18	0.39
SNORD29	RF00070	1	0.3	1	0.3	1	0.3	0	0.3
mir-29	RF00074	0.96	1	0.82	0.6	0.56	0.6	0.11	0.4
RNAI	RF00106	1	1	1	1	0.97	1	0.91	0.9
SIB_RNA	RF00113	1	1	1	1	1	1	1	1
snoZ159	RF00160	1	0.3	1	0.1	0.33	0.1	0	0.1
Hammerhead_1	RF00163	1	0.31	1	0.1	0	0.03	0	0.03
Purine	RF00167	1	0.89	1	0.83	1	0.55	1	0.14
SSU_rRNA_bacteria	RF00177	1	1	1	1	1	1	1	1
IRES_Bag1	RF00222	1	1	1	1	1	1	1	1
glmS	RF00234	1	1	1	0.89	1	0.89	1	0.33
ctRNA_pGA1	RF00236	1	1	1	1	1	0.6	1	0.2
RNA-OUT	RF00240	1	1	1	1	1	1	1	1
ctRNA_pT181	RF00242	1	1	1	0.94	1	0.69	1	0.62
IRES_L-myc	RF00261	1	1	1	1	0.82	1	0.23	1
SCARNA18	RF00283	1	1	1	1	1	1	0.86	0.95
SCARNA8	RF00286	1	1	1	1	1	1	1	1
snoR86	RF00303	1	1	1	1	1	1	0.88	1
snoZ157	RF00333	1	1	1	1	0.94	0.9	0.83	0.8
snoR60	RF00339	1	1	0.96	1	0.96	1	0.72	0.9
ydaO-yuaA	RF00379	1	1	1	0.99	1	0.94	0.9	0.84
Antizyme_FSE	RF00381	1	1	1	0.92	0.99	0.62	0.91	0.46
Pox_AX_element	RF00384	1	1	1	1	1	1	1	1
IBV_D-RNA	RF00385	1	1	1	1	1	1	1	0.9
SNORA30	RF00415	1	1	1	0.73	1	0.73	1	0.64
SCARNA24	RF00422	1	1	1	1	1	1	1	0.87
SCARNA15	RF00426	1	1	1	1	1	0.77	0.88	0.64
SCARNA23	RF00427	1	1	1	1	1	1	1	0.94
Hsp90_CRE	RF00433	1	1	1	1	1	1	1	1
ROSE	RF00435	1	1	1	0.46	1	0.15	0	0.15
IRES_HIF1	RF00449	1	1	1	1	1	0.94	0.87	0.88
IRES_mnt	RF00457	1	1	1	1	1	1	1	0.95
HCV_SLVII	RF00468	1	1	1	1	1	1	1	0.91
HCV_SLIV	RF00469	1	1	1	1	1	1	1	0.94
SCARNA6	RF00478	1	1	1	0.94	1	0.94	1	0.94
HIV_FE	RF00480	1	1	1	0.99	1	0.99	0.92	0.86
IRES_Cx43	RF00487	1	1	1	1	0.98	1	0.83	0.93
U1_yeast	RF00488	1	1	1	1	1	1	1	0.8
ctRNA_p42d	RF00489	0.97	1	0.81	0.9	0.53	0.6	0.03	0.2
IRES_Hsp70	RF00495	1	0.86	1	0.86	1	0.86	1	0.86
RNAIII	RF00503	1	1	1	1	0.56	1	0.56	1
Thr_leader	RF00506	1	1	1	1	1	1	0.95	0.96
snosnR64	RF00509	1	1	1	1	0.91	1	0.82	0.9
Leu_leader	RF00512	1	1	1	1	1	1	1	1
Trp_leader	RF00513	1	0.95	1	0.91	0.98	0.59	0.78	0.36
His_leader	RF00514	1	1	1	1	1	0.97	0.99	0.79

Continued on next page

Table 2 – continued from previous page

Rfam name	Rfam id	A	B	C	D	E	F	G	H
PreQ1	RF00522	1	0.74	0.84	0.26	0.14	0.17	0	0.14
Flavivirus_DB	RF00525	1	1	1	1	1	0.92	1	0.68
snoMe28S-G3255	RF00527	1	1	1	0.5	1	0.3	0	0.1
IRES_TrkB	RF00547	1	1	1	1	1	0.94	1	0.94
IRES_c-sis	RF00549	1	1	1	1	1	1	1	1
L13_leader	RF00555	1	0.35	1	0.35	1	0.35	0.95	0.35
L19_leader	RF00556	1	1	1	0.6	0.88	0.2	0	0.2
L20_leader	RF00558	1	0.79	1	0.3	1	0.21	0.94	0.21
L21_leader	RF00559	1	0.87	1	0.74	0.96	0.45	0.46	0.11
SCARNA3	RF00565	1	1	1	1	1	1	1	0.96
SCARNA14	RF00582	1	1	1	1	1	1	1	0.86
CoTC_ribozyme	RF00621	1	1	1	1	1	1	1	0.9
CPEB3_ribozyme	RF00622	1	1	1	1	1	0.92	1	0.67
P1	RF00623	1	1	1	1	1	0.86	1	0.5
P24	RF00629	1	1	1	1	1	1	1	1
MIR169.2	RF00645	1	0.32	1	0.07	1	0.03	1	0.03
MIR168	RF00677	1	1	1	1	0.97	0.9	0.87	0.6
MIR162.2	RF00742	1	1	1	0.9	0.75	0.1	0.25	0.1
mir-342	RF00760	1	1	1	1	1	1	1	1
mir-541	RF00777	1	1	1	0.9	0.92	0.9	0.62	0.9
mir-1255	RF00994	0.99	0.9	0.86	0.9	0.45	0.9	0.11	0.1
WLE3	RF01046	1	1	1	0.7	1	0.7	1	0.6
Sacc_telomerase	RF01050	1	1	1	1	1	1	1	1
preQ1-II	RF01054	1	1	1	0.93	1	0.71	1	0.64
MOCO_RNA_motif	RF01055	1	0.33	1	0.13	1	0.07	1	0.02
RF_site2	RF01076	1	1	0.83	1	0.67	1	0.67	1
RF_site3	RF01079	1	1	1	1	1	1	1	0.5
RF_site5	RF01093	1	1	1	1	1	0.58	0.9	0.5
RF_site9	RF01098	1	1	1	1	1	1	1	1
PK-G12rRNA	RF01118	1	1	1	1	1	1	1	1
snoZ30a	RF01196	1	1	1	1	1	1	1	1
snoR103	RF01213	0.87	1	0.87	1	0.87	0.82	0.87	0.73
snoR442	RF01232	1	1	1	1	0.25	0.7	0	0.1
snR161	RF01237	1	1	1	0.9	1	0.9	1	0.5
snR36	RF01242	1	1	1	1	1	1	1	1
snR8	RF01248	1	1	1	1	1	0.91	1	0.91
snR190	RF01249	1	1	1	1	1	0.8	1	0.8
snR5	RF01252	1	1	1	1	1	0.82	1	0.82
snR35	RF01255	1	1	1	1	1	1	1	1
snR191	RF01263	1	1	1	1	1	1	1	1
SCARNA2	RF01268	1	0.95	1	0.95	1	0.95	1	0.95
snoR2	RF01292	1	1	1	1	1	1	1	1
SCARNA7	RF01295	1	1	1	1	1	0.94	1	0.94
AHBV_epsilon	RF01313	1	1	1	1	1	1	0.88	1
CRISPR-DR2	RF01315	1	0.74	1	0.05	0	0.05	0	0
CRISPR-DR3	RF01316	0.5	0.1	0.5	0.05	0	0	0	0
CRISPR-DR5	RF01318	1	1	1	0.08	1	0.08	0	0

Continued on next page

Table 2 – continued from previous page

Rfam name	Rfam id	A	B	C	D	E	F	G	H
CRISPR-DR7	RF01320	1	0.9	1	0.2	1	0.1	0	0
CRISPR-DR35	RF01345	1	1	1	1	1	1	1	0
CRISPR-DR53	RF01366	1	1	1	1	1	1	1	0
CRISPR-DR60	RF01373	1	1	1	1	1	0.5	0	0.5
CRISPR-DR61	RF01374	1	1	0.83	1	0.83	1	0	0.5
CRISPR-DR65	RF01378	1	1	1	1	1	1	0	0
isrA	RF01385	1	1	1	1	0.97	1	0.97	1
istR	RF01400	1	1	1	1	1	1	0.97	1
NrrF	RF01416	1	1	1	1	1	1	1	1
IsrR	RF01419	1	0.98	1	0.97	1	0.91	1	0.88
VrrA	RF01456	1	1	1	1	0.95	1	0.84	1
Afu_300	RF01509	1	1	1	1	1	1	0.61	0.5
MFR	RF01510	1	1	1	1	1	1	1	0.67
Afu_309	RF01512	1	1	1	1	1	1	1	1
Dictyostelium_SRP	RF01570	1	1	1	1	1	1	1	1
RNase_P	RF01577	1	1	1	1	1	1	1	1
AdoCbl-variant	RF01689	1	0.62	1	0.03	1	0.01	1	0.01
Lnt	RF01711	1	0.9	1	0.8	1	0.3	0	0.3
cspA	RF01766	1	1	1	1	1	1	1	1
SMK_box_riboswitch	RF01767	1	0.6	1	0.08	1	0.04	1	0.04
rnk_leader	RF01771	0.97	1	0.97	1	0.97	0.85	0.97	0.85
RatA	RF01776	1	1	0.88	1	0.35	0.56	0.04	0.5
blv_FSE	RF01785	1	1	1	1	1	0	0	0
FourU	RF01795	1	1	1	1	1	1	0.94	1
fstAT	RF01797	1	1	1	1	0.94	1	0.94	0.73
HSUR	RF01802	1	0.5	1	0.5	1	0.5	1	0.5
Lambda_thermo	RF01804	1	1	1	1	1	1	1	1
GIR1	RF01807	1	1	0.89	0.92	0.89	0.92	0.89	0.92
MicX	RF01808	1	1	1	1	1	1	1	1
symR	RF01809	1	1	1	1	1	1	1	1
PtaRNA1	RF01811	1	1	1	1	1	1	1	0.75
rdlD	RF01813	1	1	1	1	1	1	1	0.98
ROSE_2	RF01832	1	1	1	1	0.99	1	0.94	1
HIV_FS2	RF01835	1	1	1	1	1	1	1	0.79
ovine_lenti_FSE	RF01840	1	1	1	1	1	1	1	0.93
veev_FSE	RF01841	1	1	1	1	0.5	0.9	0.5	0.7
alpha_tmRNA	RF01849	1	1	1	1	1	1	0.98	1
tRNA-Sec	RF01852	0.9	0.32	0.53	0.28	0.43	0.28	0.08	0.28
MIAT_exon1	RF01874	1	1	1	1	1	1	0.98	0.9
MIAT_exon5_2	RF01876	1	1	1	1	1	1	1	1
HSR-omega_2	RF01886	1	1	1	1	1	1	0.86	1
mir-2241	RF01899	1	1	1	0.5	1	0.5	1	0.5
mir-284	RF01901	1	1	1	1	1	1	1	1
HEARO	RF02033	1	0.27	1	0.16	1	0.05	1	0.01
STnc630	RF02052	1	1	1	1	1	1	1	1
STnc370	RF02064	1	1	1	1	1	0.8	1	0.8
STnc180	RF02079	1	0.8	1	0.5	1	0.3	1	0.3

Continued on next page

Table 2 – continued from previous page

Rfam name	Rfam id	A	B	C	D	E	F	G	H
OrzO-P	RF02083	1	1	1	1	1	1	1	0.43
Yar_1	RF02085	1	0.9	1	0.8	1	0.5	1	0.3
tfoR	RF02100	1	1	1	1	1	1	1	1
IS009	RF02111	1	1	1	1	0.97	1	0.74	0.73
FAM13A-AS1.1	RF02114	0.92	1	0.92	0.8	0.92	0.3	0.92	0.2
FAM13A-AS1.2	RF02115	1	0.8	1	0.6	1	0.3	1	0.3
MEG8.3	RF02147	1	1	1	1	0.92	0.6	0.42	0.4
PVT1.4	RF02167	1	1	1	0.9	1	0.5	1	0.4
HPnc0260	RF02194	1	0.39	1	0.26	1	0.23	1	0.19
WT1-AS.1	RF02203	1	1	1	1	1	1	0.9	0.8
sX5	RF02224	1	1	1	1	1	1	0.59	0.8
sX11	RF02230	1	1	1	1	1	1	1	1
Six3os1.3	RF02248	1	1	1	1	1	0.9	0.89	0.8
Hammerhead_II	RF02276	1	0.83	1	0.79	0.73	0.54	0	0
Hammerhead_HH10	RF02277	1	1	1	1	1	1	1	1
hsp17	RF02358	1	0.67	1	0.67	0.83	0.67	0.33	0.67
PyrG_leader	RF02371	1	1	1	0.7	0.8	0.2	0.3	0.2
PyrD_leader	RF02373	1	0.59	1	0.15	1	0.11	1	0.11
Ms_AS-8	RF02466	1	1	1	1	1	1	0.78	0.8
GLRNase_MRP	RF02472	1	1	1	1	1	1	1	1
GLU1	RF02491	1	1	1	1	1	1	1	1
GLU2	RF02492	1	1	1	1	1	1	1	1
GLU4	RF02493	1	1	1	1	1	1	1	1
GLU6	RF02494	1	1	1	1	1	1	1	1
ohsC.RNA	RF02495	1	1	1	1	1	0.97	1	0.97
mir-2494	RF02518	1	1	1	0.9	1	0.9	1	0.7
ToxI	RF02519	1	1	1	1	1	1	0.62	0.5
ROSE.3	RF02523	1	1	1	1	1	1	1	0.88
NRF2_IRES	RF02531	0.98	1	0.98	1	0.97	0.95	0.86	0.9
MNV_3UTR	RF02532	1	1	1	1	1	1	1	1
ODC_IRES	RF02535	1	1	1	1	1	1	0.97	0.85
mt-tmRNA	RF02544	1	1	1	0.91	1	0.82	0.67	0.36

## E Negative control set

We used coding sequences, ancestral repeats, untranslated regions (UTRs) and random sequences to perform a negative control. According to the procedure for structured and diverse RNA families the sequences of the negative control set were used as a input sequence for **RNAlien**. Taxonomic start points for the construction were set as below using taxids from NCBI taxonomy [2]. The results were summarized for each subset individually.

### E.1 Random sequences

A test with 300 different 100 nucleotides long random sequences was performed. 100 Sequences each were used in *Escherichia coli*, *Homo sapiens* and *Sulfolobus solfataricus*. The sequences were created with a inhouse *randseq* program, source code will be provided on request by Ivo L. Hofacker (ivo@tbi.univie.ac.at).

### E.2 Ancestral repeats

All 62 entries tagged with ancestral repeat from the **Dfam** [3] database were used with *Homo sapiens* as starting point for RNAlien, if the repeat was present there. The exceptions are the following list of pairs, with the first element containing the family name and the second the taxonomic start point: (Charlie12\_Rodent,*Mus musculus*), (DNA9TA1\_DR,*Danio rerio*), (L2-1\_DR,*Danio rerio*), (Jockey2,*Drosophila melanogaster*), (DIVER2\_I,*Drosophila melanogaster*)

### E.3 Coding sequences

50 Protein coding sequences were checked for *Escherichia coli*, *Sulfolobus solfataricus* and *Homo sapiens*. *Escherichia coli* sequences are the first 50 annotated CDS sequences from regulonDB 9.0 [4] ([http://regulondb.ccg.unam.mx/menu/download/datasets/files/Gene\\_sequence.txt](http://regulondb.ccg.unam.mx/menu/download/datasets/files/Gene_sequence.txt)) . *Sulfolobus solfataricus* sequences are retrieved from the reference genbank [5] assembly for *Sulfolobus solfataricus* *GCF\_000007005.1\_ASM700v1*. *Homo sapiens* sequences are from Ensemble [6] (Release 84, GRCh38.p5), chromosome2.

### E.4 UTR regions

50 3-prime and 5-prime untranslated regions from *E.coli* and *Homo sapiens* were checked. *Escherichia coli* sequences are from regulonDB version 9.0 [4] ([http://regulondb.ccg.unam.mx/menu/download/datasets/files/UTR\\_5\\_3\\_sequence.txt](http://regulondb.ccg.unam.mx/menu/download/datasets/files/UTR_5_3_sequence.txt)), *Homo sapiens* sequences are from Ensemble [6] (Release 84, GRCh38.p5), chromosome2. For *Sulfolobus* we could not find a UTR dataset.

Table 3: Table for negative control set construction results. Shown are selected result fields of RNAs, RNACode and cmdstat. Column names A to Y are placeholders for following names: Name(=A), alienFastaNumber(=B), meanPairwiseIdentity(=C), shannonEntropy(=D), gcContent(=E), meanSingleSequenceMFE(=F), consensusMFE(=G), energyContribution(=H), covarianceContribution(=I), combinationsPair(=J), meanZScore(=K), SCI(=L), svmDecisionValue(=M), svmRNAClassProbability(=N), prediction(=O), RNACodeLowestP-value(=P), RNACodeClassification(=Q), statSequenceNumber(=R), statEffectiveSequences(=S), statConsensusLength(=T), statW(=U), statBasepairs(=V), statBifurcations(=W), relativeEntropyCM(=X), relativeEntropyHMM(=Y)

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y
hs_random1	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.7	100	118	28	1	0.591	0.318
hs_random2	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.61	100	118	34	1	0.59	0.266
hs_random3	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.75	100	118	24	2	0.589	0.369
hs_random4	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.78	100	118	26	3	0.59	0.34
hs_random5	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.75	100	118	27	2	0.591	0.335
hs_random6	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.75	100	118	27	4	0.592	0.335
hs_random7	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.69	100	117	27	2	0.589	0.335
hs_random8	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.69	100	119	28	2	0.59	0.319
hs_random9	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.71	100	118	29	0	0.589	0.322
hs_random10	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.72	100	132	26	2	0.59	0.348
hs_random11	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.65	100	117	33	1	0.589	0.276
hs_random12	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.78	100	118	28	1	0.59	0.321
hs_random13	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.8	100	118	20	1	0.589	0.401
hs_random14	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.62	100	118	33	2	0.589	0.278
hs_random15	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.85	100	117	19	2	0.591	0.414
hs_random16	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.79	100	117	23	0	0.588	0.371
hs_random17	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.75	100	118	27	1	0.59	0.336
hs_random18	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.72	100	118	26	1	0.589	0.343
hs_random19	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.79	100	118	26	2	0.591	0.34
hs_random20	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.76	100	118	29	2	0.591	0.328
hs_random21	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.65	100	118	33	2	0.59	0.282
hs_random22	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.69	100	117	31	2	0.589	0.299
hs_random23	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.69	100	118	27	1	0.59	0.331
hs_random24	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.69	100	118	31	1	0.591	0.291
hs_random25	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.69	100	118	30	3	0.589	0.303
hs_random26	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.71	100	118	31	1	0.59	0.287
hs_random27	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.66	100	117	30	1	0.59	0.306
hs_random28	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.73	100	119	28	2	0.589	0.326
hs_random29	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.69	100	117	30	1	0.591	0.308
hs_random30	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.62	100	117	35	1	0.59	0.249
hs_random31	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.76	100	118	27	1	0.589	0.334
hs_random32	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.94	100	118	18	1	0.592	0.432
hs_random33	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.87	100	118	22	0	0.59	0.39
hs_random34	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.69	100	119	32	1	0.59	0.288
hs_random35	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.68	100	118	28	2	0.59	0.318
hs_random36	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.73	100	118	26	2	0.591	0.341
hs_random37	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.75	100	118	25	2	0.589	0.35
hs_random38	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.67	100	119	32	2	0.59	0.281
hs_random39	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.77	100	118	24	2	0.59	0.365
hs_random40	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.82	100	117	24	1	0.591	0.369
hs_random41	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.75	100	118	26	2	0.589	0.345
hs_random42	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.82	100	118	25	1	0.591	0.356
hs_random43	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.76	100	117	25	2	0.59	0.352
hs_random44	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.73	100	118	27	2	0.59	0.33
hs_random45	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.65	100	118	30	1	0.59	0.305
hs_random46	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.78	100	118	23	2	0.591	0.371
hs_random47	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.82	100	118	22	1	0.588	0.376
hs_random48	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.59	100	118	34	1	0.59	0.259
hs_random49	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.65	100	118	30	1	0.589	0.299
hs_random50	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.76	100	117	26	1	0.59	0.347
hs_random51	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.71	100	118	29	0	0.59	0.313
hs_random52	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.69	100	118	31	2	0.589	0.289
hs_random53	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.73	100	118	26	1	0.589	0.351
hs_random54	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.71	100	118	28	2	0.589	0.328
hs_random55	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.64	100	118	31	1	0.591	0.289
hs_random56	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.77	100	118	25	0	0.591	0.354
hs_random57	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.71	100	117	31	1	0.591	0.298
hs_random58	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.74	100	118	25	1	0.592	0.358
hs_random59	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.78	100	118	26	3	0.591	0.346
hs_random60	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.67	100	118	29	1	0.59	0.319
hs_random61	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.67	100	118	31	0	0.591	0.287
hs_random62	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.72	100	118	28	0	0.59	0.322
hs_random63	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.75	100	119	26	1	0.59	0.342
hs_random64	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.83	100	118	25	1	0.591	0.356
hs_random65	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.91	100	117	21	0	0.589	0.395
hs_random66	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.8	100	118	25	1	0.589	0.363
hs_random67	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.73	100	118	29	2	0.59	0.317
hs_random68	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.73	100	118	30	2	0.589	0.31
hs_random69	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.68	100	118	32	1	0.589	0.278
hs_random70	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.74	100	118	28	1	0.59	0.322
hs_random71	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.69	100	118	30	1	0.59	0.31
hs_random72	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.75	100	118	26	1	0.591	0.344
hs_random73	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.72	100	118	29	3	0.588	0.312
hs_random74	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.91	100	118	18	2	0.591	0.424
hs_random75	1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1	0.79	100	118	29	1	0.59	0.314
hs_random76	1																							

















Table 3 – continued from previous page

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y
ec_3putr_696	5	91.77	0.13538	0.5147	-35.72	-35.7	-35.2	-0.5	1.16	-5.73	1	5.44	0.999999	RNA	0.992	OTHER	5	1.12	77	93	26	1	0.733	0.435
ec_3putr_697	70	83.08	0.33097	0.42827	-21.95	-18.78	-20.67	1.89	1.07	-5.52	0.86	5.69	1	RNA	0.083	OTHER	70	3.7	46	68	17	1	1.198	0.995
ec_3putr_698	209	86.65	0.22301	0.57768	-19.67	-16.6	-16.6	0	1	-4.67	0.84	4.24	0.999987	RNA	0.118	OTHER	209	209	28	85	11	0	1.53	1.225
ec_3putr_699	14	93.44	0.12758	0.45678	-25.38	-25.33	-25.22	-0.11	1.11	-4.97	1	4.7	0.999996	RNA	0.384	OTHER	14	1.7	60	76	19	1	0.931	0.678
ec_3putr_700	18	83.8	0.28603	0.4945	-27.3	-24.13	-25.43	1.31	1.1	-4.75	0.88	4.82	0.999997	RNA	0.831	OTHER	18	1.31	76	93	21	2	0.743	0.501
ec_3putr_701	13	91.58	0.15036	0.4783	-18.98	-19.12	-18.98	-0.14	1.15	-3.78	1.01	3.63	0.999935	RNA	0.311	OTHER	13	4.28	39	52	14	0	1.4	1.252
ec_3putr_702	8	97.22	0.04507	0.42794	-16.38	-16.75	-16.38	-0.38	1.14	-2.76	1.02	1.94	0.994105	RNA	0.375	OTHER	8	1.9	54	69	16	1	1.027	0.805
ec_3putr_703	24	79.63	0.3742	0.51288	-14.57	-10.68	-11.02	0.33	1	-2.6	0.73	2.01	0.995079	RNA	0.204	OTHER	24	4.72	46	77	11	1	1.197	1.09
ec_3putr_704	42	80.14	0.38538	0.56406	-23.02	-20.48	-20.37	-0.11	1.08	-4.79	0.89	5.37	0.999999	RNA	0.797	OTHER	42	8.04	47	62	12	0	1.173	1.062
ec_3putr_705	5	62.53	0.60602	0.54072	-111.92	-82.64	-81.45	-1.19	1.19	-6.08	0.74	6.75	1	RNA	0.268	OTHER	5	1.13	243	279	76	5	0.59	0.315
ec_3putr_706	21	89.47	0.18397	0.37186	-21.18	-20.71	-21.1	0.39	1.07	-7.12	0.98	6.66	1	RNA	0.858	OTHER	21	2.79	47	61	14	0	1.173	0.99
ec_3putr_707	10	92.97	0.1238	0.50556	-20.56	-20.14	-20.22	0.08	1.08	-6.22	0.98	5.71	1	RNA	0.76	OTHER	10	4.71	37	51	12	0	1.472	1.349
ec_3putr_708	24	83.68	0.3139	0.55579	-13.98	-13.98	-13.98	0	1	-2.67	1	3.14	0.999759	RNA	0.906	OTHER	24	24	32	44	7	0	1.617	1.591
ec_3putr_709	191	81.28	0.3428	0.49671	-59.19	-42.89	-42.95	0.06	1.16	-2.89	0.72	2.17	0.996782	RNA	9.228E-07	PROTEIN	191	2.06	165	193	45	2	0.59	0.392
ec_3putr_710	46	76.97	0.41504	0.45547	-22.67	-18.6	-18.55	-0.05	1.08	-4.94	0.82	5.35	0.999999	RNA	0.191	OTHER	46	2.78	57	75	16	0	0.977	0.796
ec_3putr_711	5	87.62	0.17647	0.4123	-17.63	-18.32	-18.1	-0.22	1.09	-4.58	1.04	4.72	0.999996	RNA	0.525	OTHER	5	5	35	48	11	0	1.456	1.344
ec_3putr_712	15	84.81	0.26892	0.70876	-20.07	-18.3	-18.3	0	1	-3.61	0.91	3.62	0.999933	RNA	0.809	OTHER	15	9.39	34	47	12	0	1.595	1.483

## References

1. Statistics of of sequence similarity scores. <http://www.ncbi.nlm.nih.gov/BLAST/tutorial/>. Accessed: 2016-03-19.
2. David L Wheeler, Tanya Barrett, Dennis A Benson, Stephen H Bryant, Kathi Canese, Vyacheslav Chetvernin, Deanna M Church, Michael DiCuccio, Ron Edgar, Scott Federhen, et al. Database resources of the national center for biotechnology information. *Nucleic acids research*, 35(suppl 1):D5–D12, 2007.
3. Travis J Wheeler, Jody Clements, Sean R Eddy, Robert Hubley, Thomas A Jones, Jerzy Jurka, Arian FA Smit, and Robert D Finn. Dfam: a database of repetitive dna based on profile hidden markov models. *Nucleic acids research*, 41(D1):D70–D82, 2013.
4. Heladia Salgado, Martin Peralta-Gil, Socorro Gama-Castro, Alberto Santos-Zavaleta, Luis Muñiz-Rascado, Jair S García-Sotelo, Verena Weiss, Hilda Solano-Lira, Irma Martínez-Flores, Alejandra Medina-Rivera, et al. Regulondb v8. 0: omics data sets, evolutionary conservation, regulatory phrases, cross-validated gold standards and more. *Nucleic acids research*, 41(D1):D203–D213, 2013.
5. Dennis A Benson, Mark Cavanaugh, Karen Clark, Ilene Karsch-Mizrachi, David J Lipman, James Ostell, and Eric W Sayers. Genbank. *Nucleic acids research*, page gks1195, 2012.
6. Andrew Yates, Wasiu Akanni, M Ridwan Amode, Daniel Barrell, Konstantinos Billis, Denise Carvalho-Silva, Carla Cummins, Peter Clapham, Stephen Fitzgerald, Laurent Gil, et al. Ensembl 2016. *Nucleic acids research*, 44(D1):D710–D716, 2016.