# What's new in RNAcentral and Rfam

Anton Petrov
apetrov@ebi.ac.uk

*Benasque - July 20, 2018*

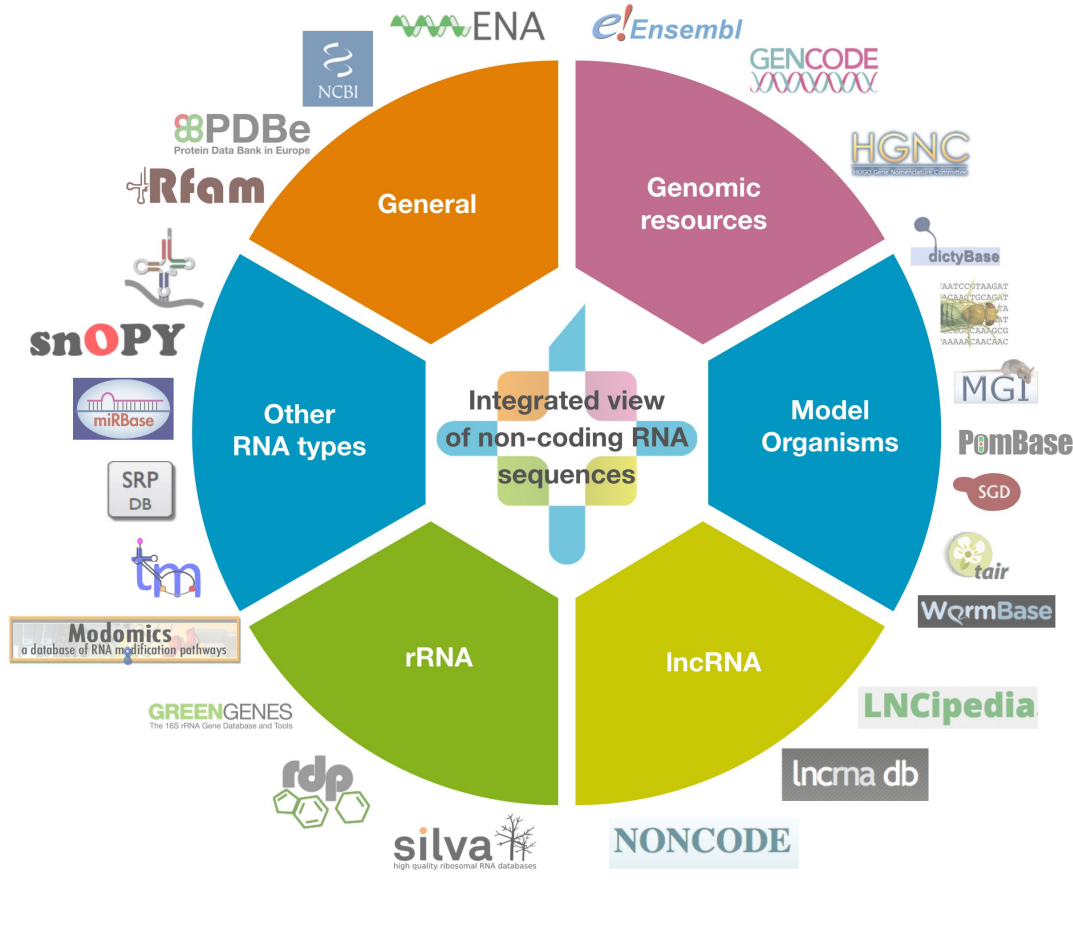EMBL-EBI

# What is RNAcentral

EMBL-EBI

# RNAcentral

## The non-coding RNA sequence database

[rnacentral.org](rnacentral.org)

- >10 million sequences
- 27 databases
- 800,000 species

EMBL-EBI

# RNAcentral has lots of useful data

- Sequence
- Description
- RNA type
- Links to other databases
- Genome locations
- Publications
- RNA modifications from MODOMICS and PDB

## Sequence URS00005A4DCF

Homo sapiens (human) microRNA hsa-miR-125a-5p

**24** nucleotides   **3** databases (ENA, miRBase, RefSeq)   Found in **6** species   miRNA

| Overview | Taxonomy | | ⬇ Download ⌄ |

### Annotations 4 total

| Database | Description |
| --- | --- |
| RefSeq | **Homo sapiens (human) hsa-miR-125a-5p.**<br>› RefSeq: NR_029693.1 ☑ - mature miRNA (precursor URS000075D168)<br>› NCBI GeneID: 406910 ☑ |
| miRBase | **Homo sapiens (human) microRNA hsa-miR-125a-5p**<br>› miRBase: MI0000469 ☑ - mature miRNA (precursor URS000075D168) |
| miRBase | **Homo sapiens (human) microRNA hsa-miR-125a-5p**<br>› miRBase: MIMAT0000443 ☑ |
| ENA | **Homo sapiens (human) miscellaneous RNA**<br>› ENA: HB859642.1:1..24:misc_RNA ☑ |

http://rnacentral.org/rna/URS00005A4DCF/9606

EMBL-EBI

# RNAcentral

Search RNAcentral    🔍 Search

*Examples:* human HOTAIR, Homo sapiens, tRNA, miRBase, 4V4Q

## RNAcentral: The non-coding RNA sequence database

More about RNAcentral →

## Getting started

### 🔍 Text search

Search by *gene, species, ncRNA type* or any other keyword

Browse sequences

### 🧩 Sequence search

Search for similar sequences or look up your sequence in RNAcentral

Search by sequence

### 📍 Genome browser

Explore RNAcentral sequences in your favorite genome locations

Browse genomes

## ncRNA data provided by 26 databases:

**GENCODE**

*47,753* sequences
Example
`Updated`

*12,011* sequences
Example

**GG**

*1.0 million* sequences
Example

*956* sequences
Example

**MGI**

*16,719* sequences
Example

Until recently

just data aggregation,

now additional analysis

# Two important new features

1. **Quality control** using Rfam models

2. Comprehensive **genome mapping**

EMBL-EBI

# 1. Rfam models are used to annotate RNAcentral



**With Rfam annotation**
50.2%

**New Rfam matches**
39.4%

RNAcentral
release 5

9 386 112
sequences

**Not suitable for Rfam**
5.2%

**Potential new Rfam families**
1.8%

**Undetected tRNA, rRNA and tmRNA**
3.4%

- **~90%** of RNAcentral sequences match Rfam models

- about **2%** of RNAcentral sequences can be used to build new Rfam models

**Natalia Quiñones Olvera**

EMBL-EBI

# Rfam annotations help detect:

- **truncated** sequences

- potential **contamination**

- **missing** annotations

# 2. Comprehensive genome mapping for **>250 species**

- genome mapping Ensembl genomes and **blat**

- **>95% of sequences mapped**
  for human, mouse,
  and other key species

- one of the **largest collections**
  of ncRNA genome annotations



EMBL-EBI

# Here is an Ensembl miRNA

**Transcript: n-TSaga9-201** ENSMUST00000197675.1

| | |
|---|---|
| **Description** | nuclear encoded tRNA serine 9 (anticodon AGA) [Source:MGI Symbol;Acc:MGI:4414029 ] |
| **Location** | Chromosome 4: 10,874,064-10,874,170 forward strand. |
| **About this transcript** | This transcript has 1 exon, is associated with 308 variations and maps to 64 oligo probes. |
| **Gene** | This transcript is a product of gene ENSMUSG00000106355 [Hide transcript table] |

[Show/hide columns (1 hidden)]  Filter

| Name | Transcript ID | bp | Protein | Biotype | CCDS | Flags |
|---|---|---|---|---|---|---|
| n-TSaga9-201 | ENSMUST00000197675.1 | 107 | No protein | miRNA | - | TSL:NA   GENCODE basic |

http://www.ensembl.org/Mus_musculus/Transcript/Summary?g=ENSMUSG00000106355;r=4:10874064-10874170;t=ENSMUST00000197675

EMBL-EBI

# But is it a miRNA or a tRNA?

EMBL-EBI

# RNAcentral shows a match to a tRNA Rfam model



http://rnacentral.org/rna/URS0000A85A32/10090

# … and other annotation in this genomic location



The other sequence is a well-annotated tRNA from GtRNAdb:
http://rnacentral.org/rna/URS000038D8D3/10090

EMBL-EBI

# RNAcentral makes data consistent across databases

- **Automatically reconcile** annotations for all sequences

- **Report** problems to member databases

- **Prioritise** sequences without inconsistencies

EMBL-EBI

# Overcoming important limitation of RNAcentral

- Can we **group** related sequences into "**genes**"
  using genomic location, Rfam annotations and sequence metadata?

# Manually assigned **Gene Ontology terms**

hsa-mir-126 involved in heart development:



## Gene Ontology annotations ❓

| Qualifier | | GO Term | Evidence Code |
|---|---|---|---|
| involved_in | 🔍 | positive regulation of blood vessel endothelial cell migration | direct assay evidence used in |
| involved_in | 🔍 | positive regulation of MAPK cascade | mutant phenotype evidence us |
| involved_in | 🔍 | negative regulation of vascular endothelial cell proliferation | mutant phenotype evidence us |
| involved_in | 🔍 | negative regulation of endothelial cell apoptotic process | mutant phenotype evidence us |

View in QuickGO ↗

Huntley et al., 2018    http://rnajournal.cshlp.org/content/24/8/1005.long

EMBL-EBI

# **Secondary structures** from GtRNAdb

Benasque 2018
top-secret project:

display 2Ds for all rRNAs
in RNAcentral
using **standard layouts**

This is how RNA biologists want to see rRNA 2D

This is how it is shown in RNA databases

Secondary Structure: small subunit ribosomal RNA

*Thermus thermophilus*
(X07998)
1.cellular organisms 2.Bacteria
3.Thermus/Deinococcus group
4.Thermus group 5.Thermus
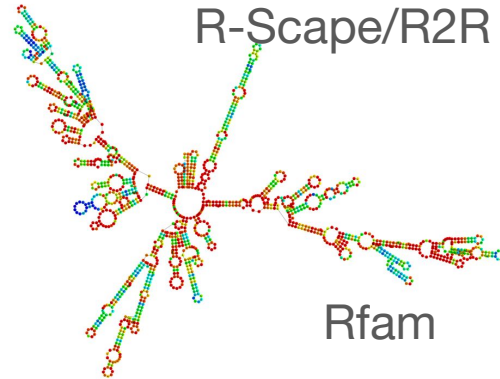September 2001

Citation and related information available at http://www.rna.icmb.utexas.edu

NDB

R-Scape/R2R

Rfam

# 1. **Comparative RNA Website** (Robin Gutell's lab)

- Contains ~1,000 of RNA secondary structures in standard layouts
- Authoritative source of rRNA data

# 2. **Traveler** software

**BMC Bioinformatics**

**SOFTWARE**                                              **Open Access**

CrossMark

# TRAVeLer: a tool for template-based RNA secondary structure visualization

Richard Elias and David Hoksza[*]

**Abstract**

**Background:** Visualization of RNA secondary structures is a complex task, and, especially in the case of large RNA structures where the expected layout is largely habitual, the existing visualization tools often fail to produce suitable visualizations. This led us to the idea to use existing layouts as templates for the visualization of new RNAs similarly to how templates are used in homology-based structure prediction.

**Results:** This article introduces Traveler, a software tool enabling visualization of a target RNA secondary structure using an existing layout of a sufficiently similar RNA structure as a template. Traveler is based on an algorithm which converts the target and template structures into corresponding tree representations and utilizes tree edit distance coupled with layout modification operations to transform the template layout into the target one. Traveler thus accepts a pair of secondary structures and a template layout and outputs a layout for the target structure.

**Conclusions:** Traveler is a command-line open source tool able to quickly generate layouts for even the largest RNA structures in the presence of a sufficiently similar layout. It is available at http://github.com/davidhoksza/traveler.
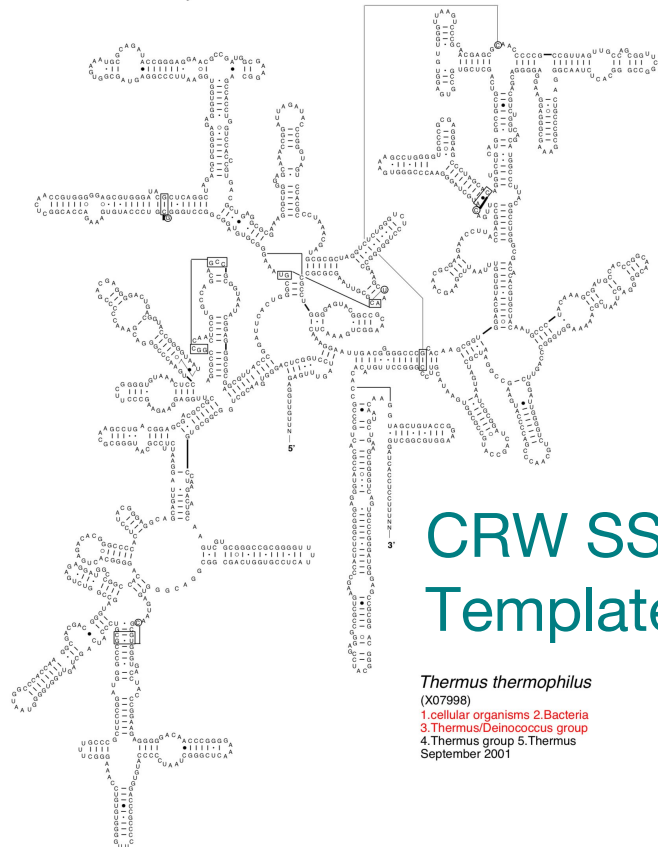
**Keywords:** Visualization, RNA secondary structure, Template-based modeling, Software tool

David Hoksza

Assistant Professor
Charles University, Prague

EMBL-EBI

Secondary Structure: small subunit ribosomal RNA
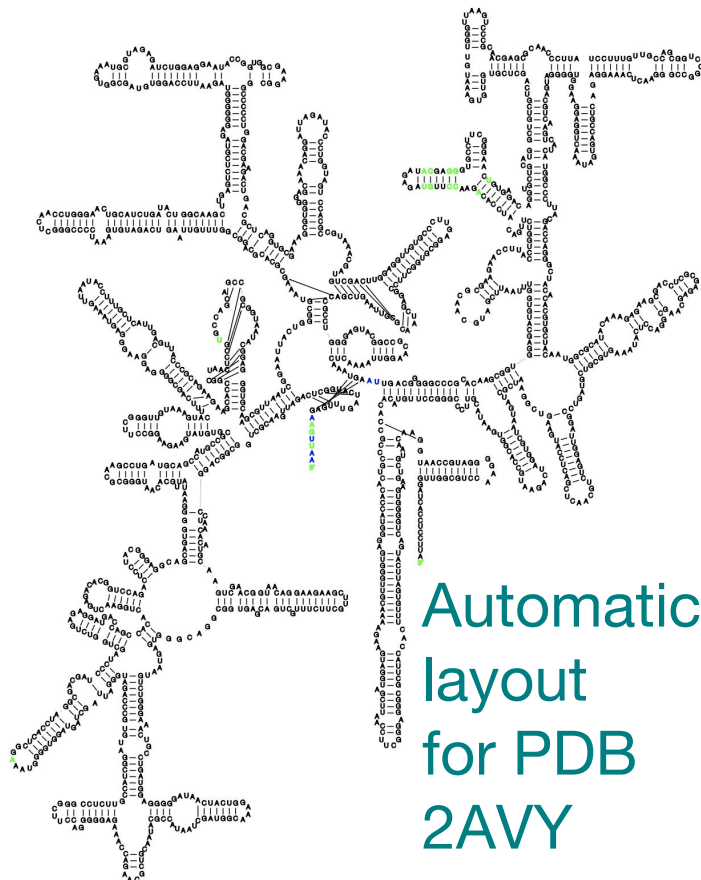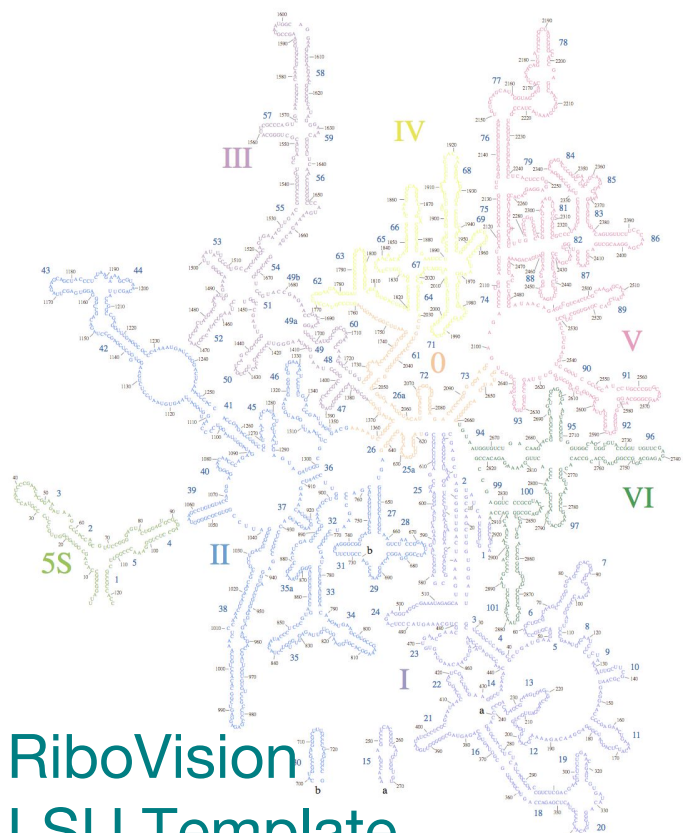
CRW SSU Template

*Thermus thermophilus*
(X07998)
1.cellular organisms 2.Bacteria
3.Thermus/Deinococcus group
4.Thermus group 5.Thermus
September 2001

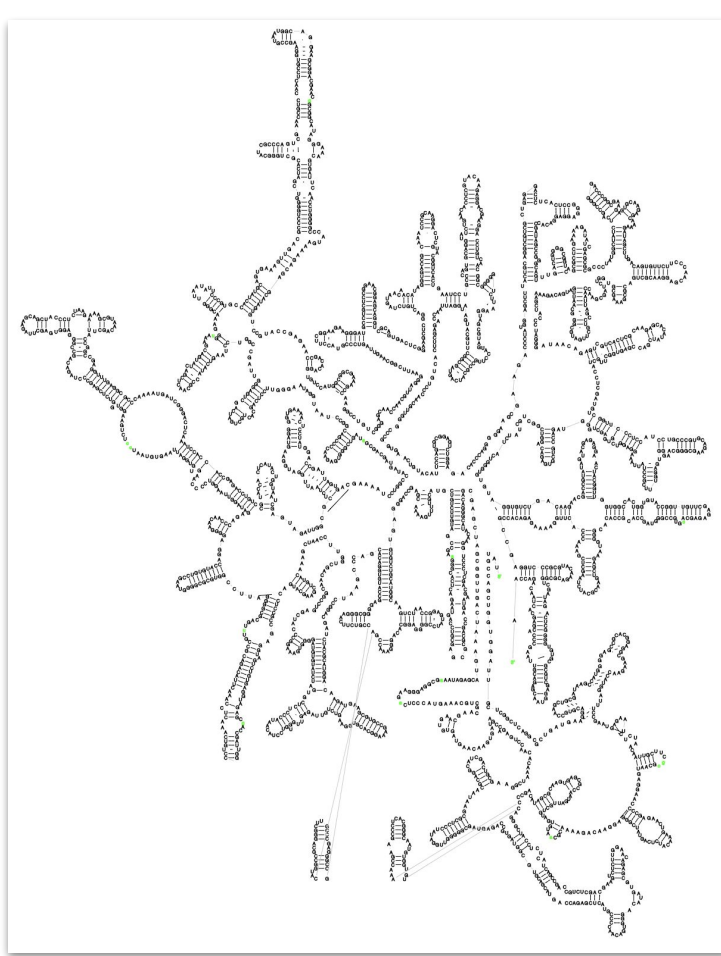Citation and related information available at http://www.rna.icmb.utexas.edu

Automatic layout for PDB 2AVY

RiboVision
LSU Template

Haloarcula marismortui
large subunit ribosomal RNA

A 3D-based secondary structure, generated by RiboVision.
Saved on 6/25/2018, 1:51:48

Automatic
layout
for PDB
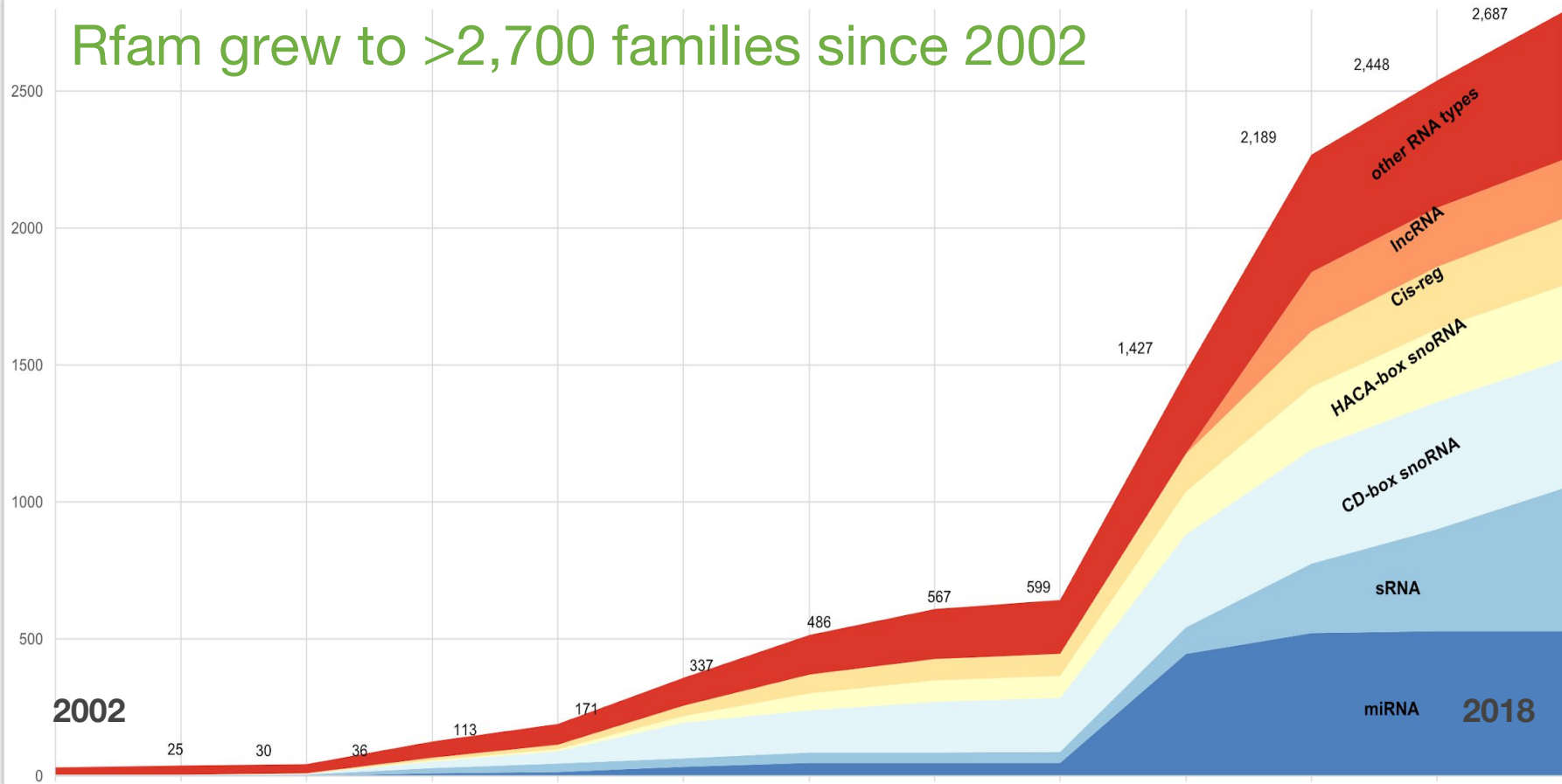1S72

EMBL-EBI

# RNAcentral keeps evolving
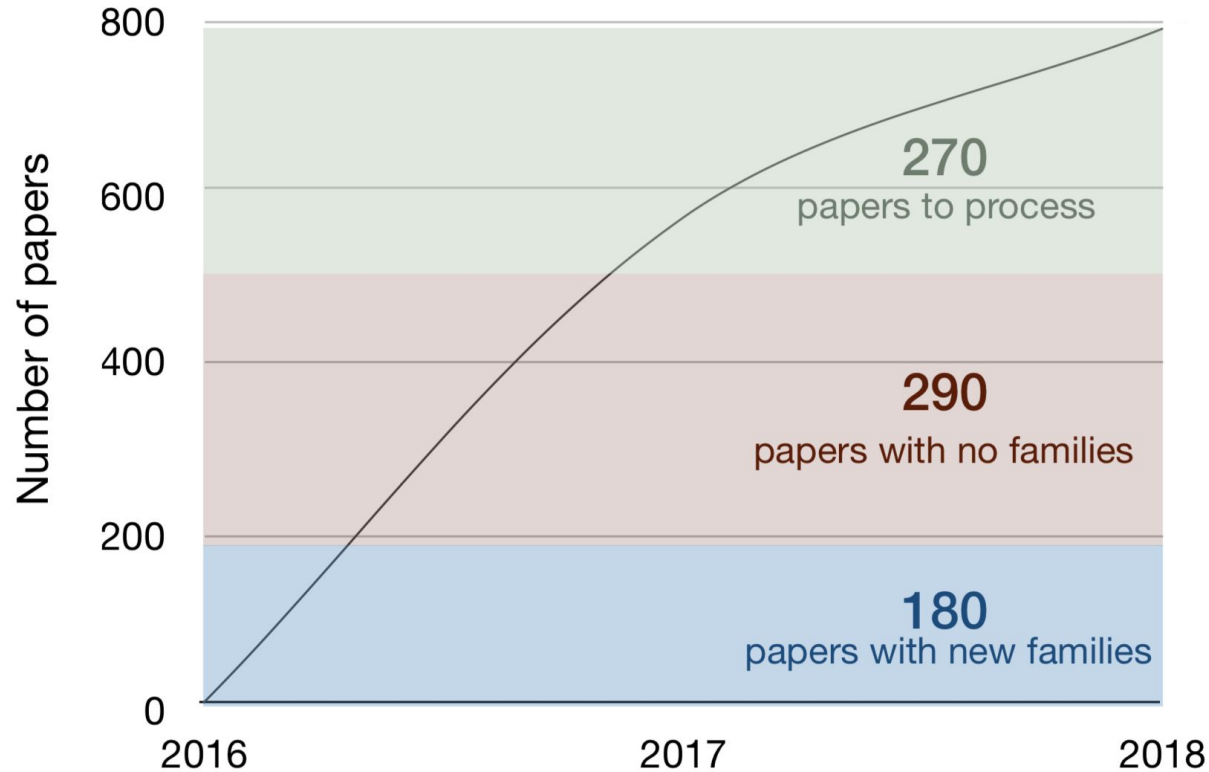
- Try it and send us your **feedback**

- Help us **improve**

EMBL-EBI

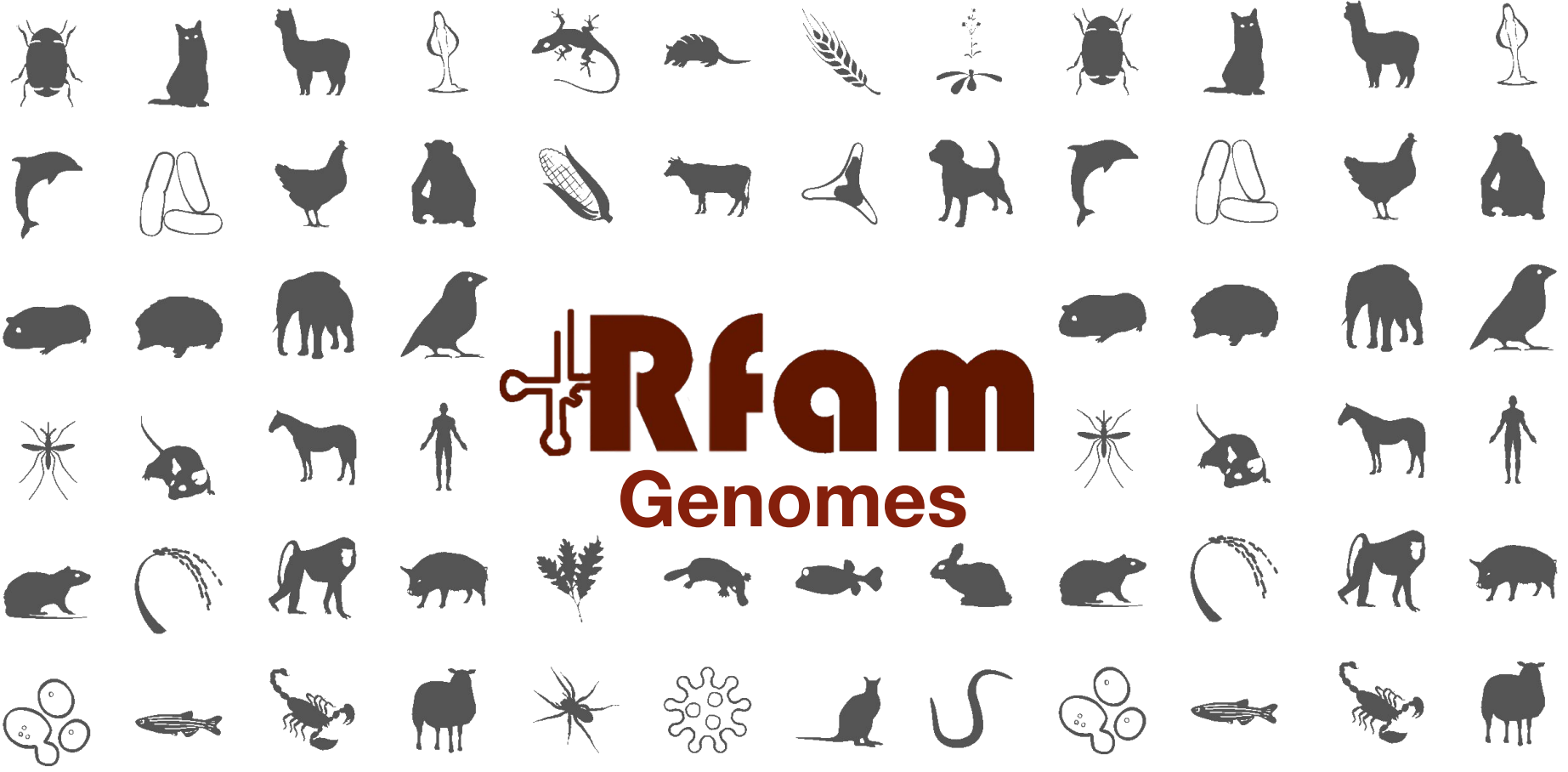# There are still more RNA families to be added



- **350** new families created since 2016
- **30%** of papers still waiting to be curated

EMBL-EBI

# What's new in Rfam

EMBL-EBI

Can you guess

**the most common user request**
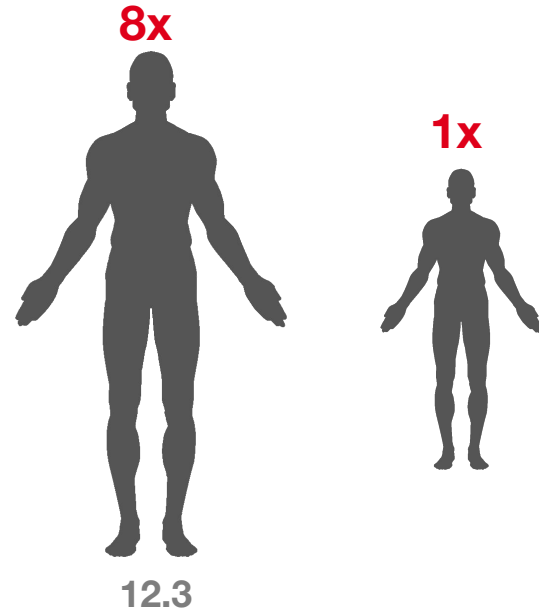
in Rfam 12.*?

Rfam
Genomes

# Previously Rfam analysed WGS and STD sequences from ENA and GenBank

- The data were **redundant**

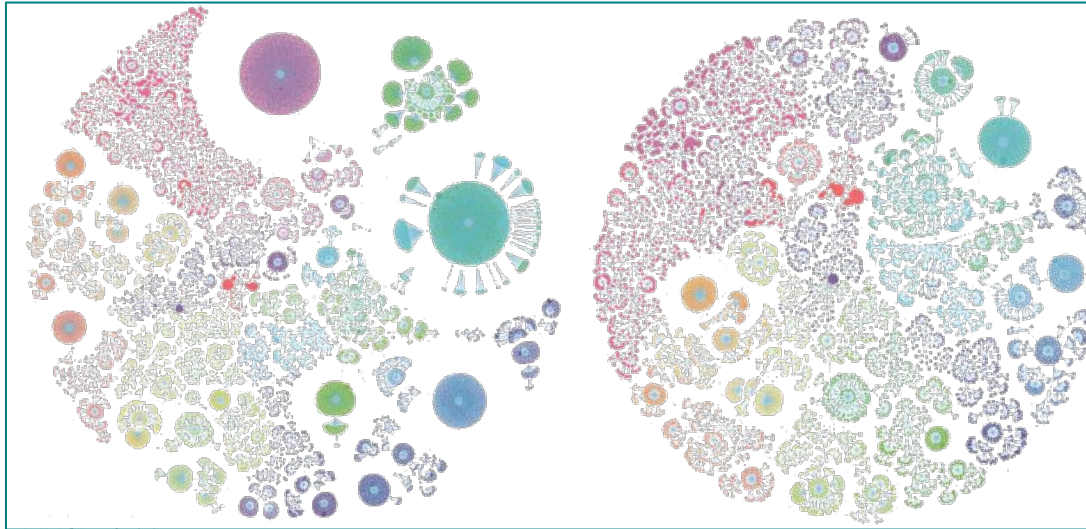- Taxonomic comparisons were **difficult**

## Redundant strains

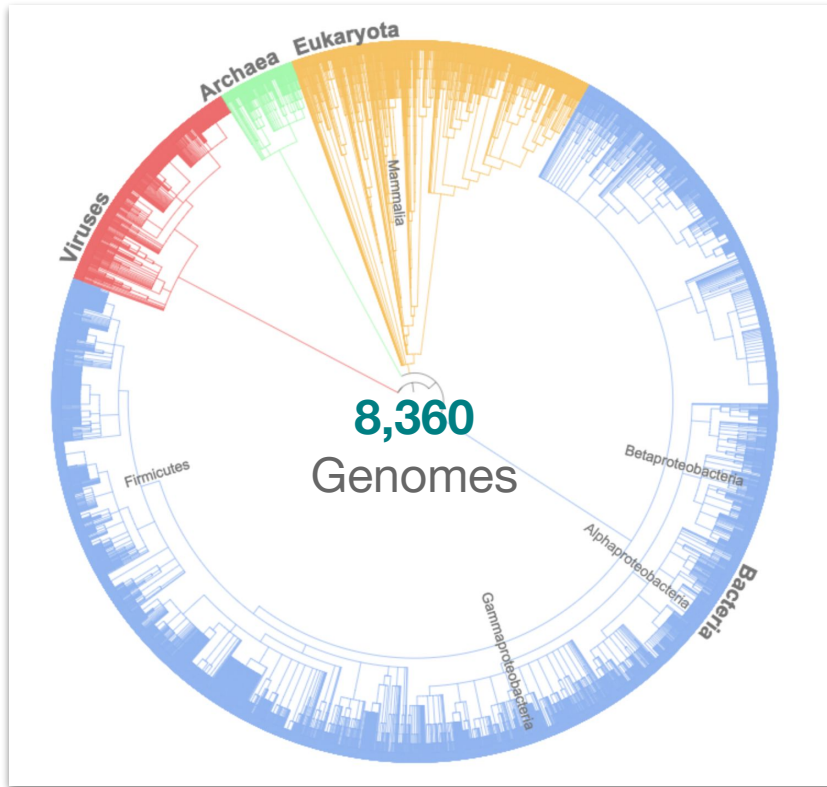|  | Strains |
|---|---|
| *Mycobacterium tuberculosis* | 77 |
| *Escherichia coli* | 348 |

## Over-represented genomes



**8x**

**1x**

12.3

EMBL-EBI

# Rfam now annotates **complete**, **representative**, and **non-redundant** genomes

Based on the **UniProt** Reference Proteome collection



https://doi.org/10.1093/database/baw139

# Rfam 13.0 is based on **8,360 genomes**



EMBL-EBI

# Need to further expand Rfam sequence database

- Viruses

- Metagenomes

- 3D structures

- RNAcentral sequences that do not match Rfam families

# Find out more about genome-centric Rfam

## Rfam 13.0: shifting to a genome-centric resource for non-coding RNA families

Ioanna Kalvari[1], Joanna Argasinska[1], Natalia Quinones-Olvera[2], Eric P. Nawrocki[3], Elena Rivas[4], Sean R. Eddy[5], Alex Bateman[1], Robert D. Finn[1] and Anton I. Petrov[1,*]

[1]European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Genome Campus, Hinxton, Cambridge CB10 1SD, UK, [2]Systems Biology Graduate Program, Harvard University, Cambridge, MA 02138, USA, [3]National Center for Biotechnology Information; National Institutes of Health; Department of Health and Human Services; Bethesda, MD 20894, USA, [4]Department of Molecular and Cellular Biology, Harvard University, Cambridge, MA 02138, USA and [5]Howard Hughes Medical Institute, Harvard University, 16 Divinity Avenue, Cambridge, MA 02138, USA

https://academic.oup.com/nar/article/4588106

EMBL-EBI

# New
# website functionality

# New **faceted** text search and **search API**

# R-Scape visualisations

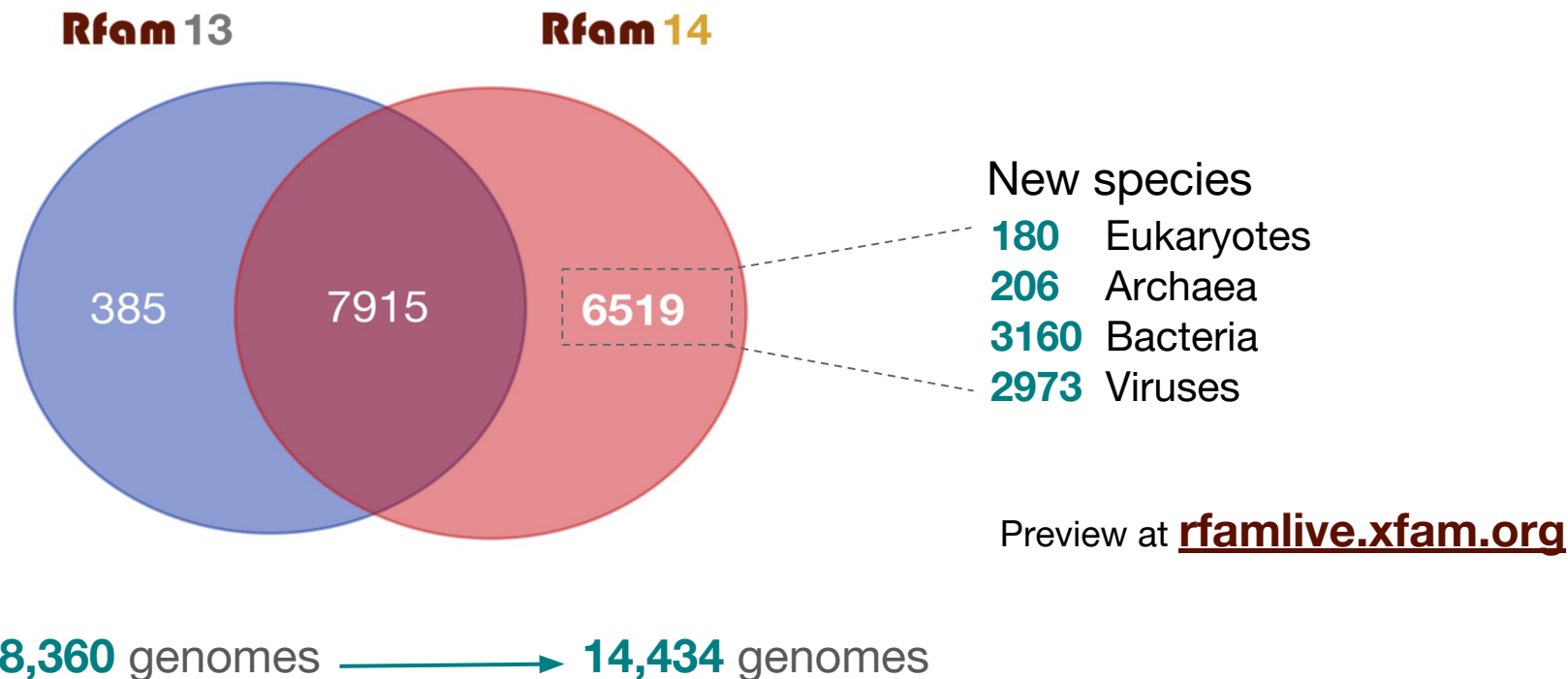# Find how to search Rfam, query public MySQL database, and more



UNIT

## Non-Coding RNA Analysis Using the Rfam Database

Ioanna Kalvari, Eric P. Nawrocki, Joanna Argasinska, Natalia Quinones-Olvera, Robert D. Finn, Alex Bateman, Anton I. Petrov
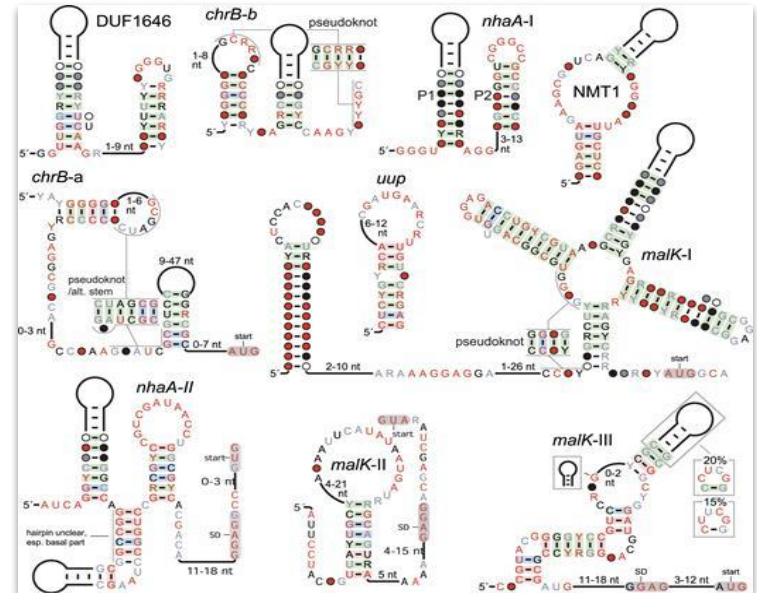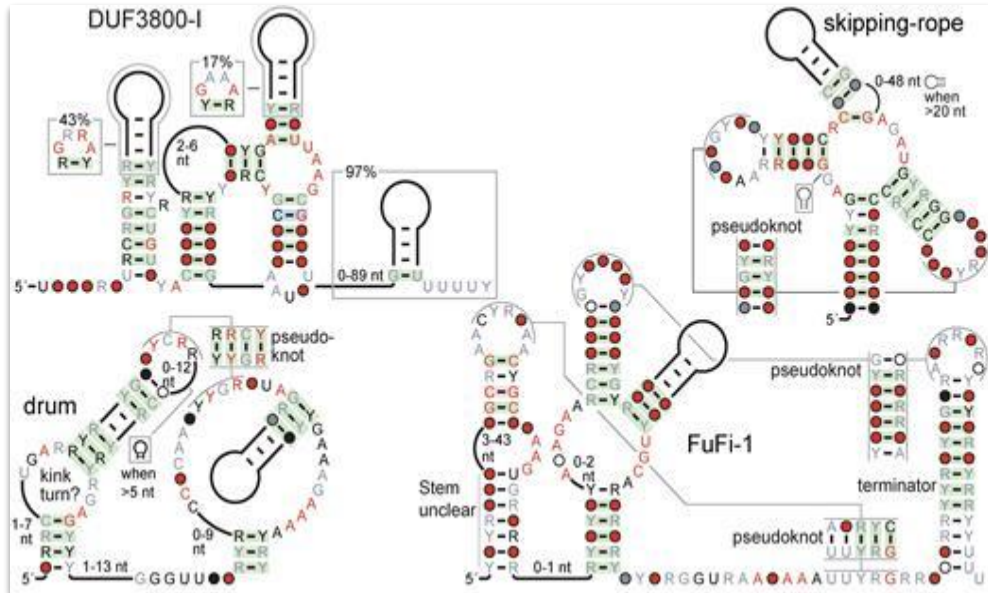
First published: 05 June 2018 | **https://doi.org/10.1002/cpbi.51**

**https://doi.org/10.1002/cpbi.51**

EMBL-EBI

# Expect Rfam 14.1 later this year

No new genomes but lots of new families from Zasha Weinberg

# Do you want to build Rfam families?

- Family curation by approved **experts**

- **Cloud-based** Rfam pipeline

- Command line or **Galaxy** access

# Special session

Wednesday 3pm

# Acknowledgements

## EMBL-EBI

Alex Bateman
Robert Finn
Ioanna Kalvari
Joanna Argasinska
Blake Sweeney
Boris Burkov

## NCBI

Eric Nawrocki

## Harvard University

Sean Eddy
Elena Rivas
Natalia Quinones-Olvera