# RNAcentral

## an International Database of ncRNA Sequences

Anton Petrov

apetrov@ebi.ac.uk

*Benasque*

*July 27th, 2015*

EMBL-EBI

# RNAcentral at Benasque 2012

**Monday, July 30**

11:00h **RNA-proteins**

N. Rajewski and J. Bujnicki

**Incorporating RNA-Protein Interactions into RNA Secondary Structure Prediction**

R. Bundschuh

18:00h **Databases**

S. Griffiths-Jones, P. Gardner and R. Knight

**RNAcentral**

S. Griffiths-Jones

**RNASTAR, greengenes and an environmental seq database**

R. Knight

**miRBase**

A. Kozomara

**Modomics and RNA processing**

J. Bujnicki

**Rfam**

S. Burge

# Why do we need RNAcentral?

Before RNAcentral:

- lots of specialized databases

- no single entry point

  for ncRNA sequence analysis

- lack of standard identifiers

*http://www.officesignspro.com/Funny-Road-Signs-2/*

EMBL-EBI

# What is RNAcentral?

- RNAcentral is a **comprehensive** and **up-to-date** database of **accessioned** ncRNA sequences that collates and integrates information from an international consortium of established RNA sequence databases.

- RNAcentral provides **broad coverage** of ncRNA types and the taxonomic space.

- Four releases since **June 2014**.

# What does RNAcentral provide?

- **unified access** to data from multiple sources

- **stable identifiers** for distinct RNA sequences

- **cross-references** to other databases

- sequence and metadata **search**

- **API** for programmatic data access

- **FTP** archive

EMBL-EBI

# Where does the data come from?



**Expert Databases**

(such as miRBase or Vega)

supply data to

RNAcentral.

*Bateman et al., 2011*

# What data is in RNAcentral?

**15 Expert Databases** imported so far:

# > 20 more Expert Databases to import

# Data from **INSDC** is imported automatically

International Nucleotide Sequence Database Collaboration =

**ENA** European Nucleotide Archive, EMBL-EBI +

**GenBank** NCBI +

**DDBJ** DNA Data Bank of Japan



http://www.insdc.org/

# Demo

http://rnacentral.org

RNAcentral homepage: http://rnacentral.org/

EMBL-EBI

Example Expert Database page: http://rnacentral.org/expert-database/mirbase

# Exploring sequences by length



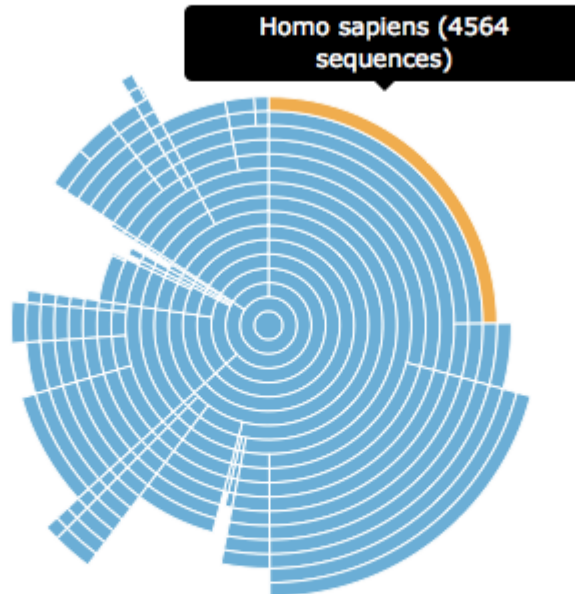You can explore sequence length distribution and **launch searches** using the interactive graph.

# Exploring species distribution



You can view what species the data comes from using the interactive **sunburst diagram**.

EMBL-EBI

# Unique RNA Sequence identifiers

Each distinct sequence gets a unique RNAcentral identifier regardless of what species it is coming from.

- Format: **URS + 10-digit hexadecimal number**

- Example: URS00000B15DA

- Sequences must be **at least 10 nucleotides** long
  - ~ 8.6 million ids assigned so far
  - > 1 trillion possible ids

EMBL-EBI

# Species-specific identifiers

RNAcentral also provides **species-specific identifiers**.

- Format: **URS identifier / NCBI taxid**

- Examples:

  - URS00003B7674**/9606**

    can also use **underscore** instead of **slash**

Gene Ontology uses RNAcentral IDs for annotating miRNAs:

http://amigo.geneontology.org/amigo/gene_product/RNAcentral:URS00004C9052_9606

# Demo

Example RNAcentral sequence page:
http://rnacentral.org/rna/URS0000086A4D

# Switching between species using taxonomic tree

# Genomic mapping

Many entries in RNAcentral come from **reference genomes**. These entries can be viewed in their genomic context using an **embedded genome browser** and their coordinates can be downloaded in GFF/GFF3/BED formats.

# Viewing entries in their genomic context

# Metadata search

- **faceting** helps to explore and filter the data

- **advanced search**

  logical operators, field-specific search and <u>more</u>

- results can be **exported** in several formats

# Sequence search

- powered by *nhmmer*

- results are stored for **7 days** and can be accessed using **unique URLs**

- results can be **sorted** by sequence identity, coverage etc

http://rnacentral.org/sequence-search

EMBL-EBI

# Instant retrieval of exact sequence matches



- The **exact sequence match** is retrieved instantly (*if it exists*)

- You can **cancel** the search if you only need the exact sequence match

Example sequence search result:
http://rnacentral.org/sequence-search/?id=07a325aa-c909-4c8f-a7a5-aee1e553ec1b

EMBL-EBI

# Demo

# Example search result

- **RNAcentral search is fast and intuitive**

- **by default all metadata associated with all entries is searched**

- **one can construct specific searches using the query syntax**



http://rnacentral.org/search?q=RNA

EMBL-EBI

# Search facets

- **Facets** allow to quickly filter search results while also exposing the kinds of data that are available

- For example, the **RNA types** facet shows how many sequences of each type are present in RNAcentral.

**Keyboard shortcut**: hitting "/" puts the cursor in the search box

**RNA types**

- [ ] rRNA (6,013,717)
- [ ] misc RNA (1,142,722)
- [ ] tRNA (864,846)
- [ ] piRNA (208,933)
- [ ] other (134,086)
- [ ] miRNA (95,673)
- [ ] snRNA (90,635)
- [ ] snoRNA (81,115)
- [ ] lncRNA (47,491)
- [ ] siRNA (45,060)
- [ ] hammerhead ribozyme (40,236)
- [ ] antisense RNA (23,920)
- [ ] precursor RNA (20,552)
- [ ] SRP RNA (14,462)
- [ ] RNase P RNA (9,524)
- [ ] tmRNA (4,716)
- [ ] scRNA (969)
- [ ] ribozyme (927)
- [ ] RNase MRP RNA (622)
- [ ] autocatalytically spliced intron (599)
- [ ] vault RNA (456)
- [ ] rasiRNA (325)
- [ ] telomerase RNA (311)
- [ ] guide RNA (133)
- [ ] ncRNA (37)
- [ ] Y RNA (18)

EMBL-EBI

# Example search: HOTAIR lncRNA

1.  Search for HOTAIR  http://rnacentral.org/search?q=HOTAIR

2.  Exclude HOTAIRM1 entries (type "not hotairm1")

    http://rnacentral.org/search?q=HOTAIR%20not%20hotairm1

3.  Focus on sequences from Vega

    using the Expert Databases facet:

    http://rnacentral.org/search?q=HOTAIR%20not%20hotairm1%20AND%20expert_db:%22VEGA%22

4.  Get just the human sequences

    using the Organisms facet:

    http://rnacentral.org/search?q=HOTAIR%20not%20hotairm1%20AND%20expert_db:%22VEGA%22%20AND%20TAXONOMY:%229606%22

EMBL-EBI

# Programmatic access

- **Learn just one API** instead instead of dozens

- Browse the API to get a sense of how it works

  http://rnacentral.org/api/v1/

- Documentation: http://rnacentral.org/api

- Example Python script:

  http://rnacentral.org/api#v1-example-script

EMBL-EBI

# FTP archive

For each release the archive contains:

- **sequences** in FASTA format

- **mapping** between RNAcentral and Expert Database identifiers

- **MD5** values for all sequences that can be used to match large numbers of sequences to RNAcentral IDs

- **genomic coordinates** in multiple formats

ftp://ftp.ebi.ac.uk/pub/databases/RNAcentral

EMBL-EBI

# RNA modifications (coming soon)

RNAcentral will import modification information from PDBe and Modomics.

For example, position **628** in URS000080DFCD is **1MA**

(6-HYDRO-1-METHYLADENOSINE)

Positional IDs:

**URS000080DFCD.628** = 1MA

Example: http://test.rnacentral.org/rna/URS000080DFCD

EMBL-EBI

# Future plans for RNAcentral

- import **more data**

- do more **genomic mapping**

- organise **related sequences** from the same species

- import **2D and 3D** structure information

- **integrate** closer with Rfam and other resources

  (Ensembl, EuropePMC, PDBe, others)

EMBL-EBI

# New RNA positions at EBI

RNAcentral and Rfam:

- RNA project leader (EBI_00582)

Rfam:

- Software developer (EBI_00551)

- Database biocurator (EBI_00552)

EMBL-EBI

# Getting in touch

- by email: apetrov@ebi.ac.uk or helpdesk@rnacentral.org

- via the website: http://rnacentral.org/contact

- on Twitter: @rnacentral

- on GitHub: https://github.com/RNAcentral/rnacentral-webcode/issues

EMBL-EBI

# Acknowledgements

- Alex Bateman
- Paul Kersey
- Guy Cochrane
- Simon Kay
- Richard Gibson
- Dan Staines
- Rob Finn
- Elspeth Bruford
- Mathew Wright
- Sameer Velankar

- DBAs
- Systems
- Web Production and other teams

# Acknowledgements - Expert Databases

Sam Griffiths-Jones

Jennifer Harrow

Christian Zwieb

Todd Lowe
Patricia Chan

Kelly Williams
Corey Hudson

Michael Clark
Camelia Quek

Mike Cherry

Paul Sternberg

James Cole
Benli Chai

Kim Pruitt

Naoya Kenmochi

Tanya Berardini

All RNAcentral Consortium members:
http://rnacentral.org/expert-databases

EMBL-EBI

# Thank you!

@RNAcentral

http://blog.rnacentral.org/

EMBL-EBI