

RNAcentral in 2026: genes and literature integration

The 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,*}

¹European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Genome Campus, Hinxton, Cambridge CB10 1SD, United Kingdom

²Independent Researcher, Santos, 11000-000, Brazil

³Dept. of Physiology, Development and Neuroscience, University of Cambridge, Cambridge CB2 3DY, United Kingdom

⁴Department of Molecular Genetics, Weizmann Institute of Science, Herzl 234, Rehovot 7610010, Israel

⁵Department of Biomolecular Engineering, University of California Santa Cruz, Santa Cruz, CA 95064, United States

⁶Department of Haematology, University of Cambridge, School of Clinical Medicine, Cambridge CB2 0PT, United Kingdom

⁷Faculty of Biology, Medicine and Health, University of Manchester, Manchester M13 9PT, United Kingdom

⁸The Arctic University of Norway, UiT - The Arctic University of Norway, 9006, Tromsø, Norway

⁹Department of Biological Sciences, Dartmouth College, Hanover, NH 03755, United States

¹⁰Laboratory of Bioinformatics and Protein Engineering, International Institute of Molecular and Cell Biology in Warsaw, Warsaw 02-109, Poland

¹¹State Key Laboratory of Epigenetic Regulation and Intervention, Institute of Biophysics, Chinese Academy of Sciences, Beijing 100101, China

¹²Department of Biochemistry, University of Cambridge, Cambridge CB2 1GA, United Kingdom

¹³Functional Gene Annotation, Preclinical and Fundamental Science, UCL Institute of Cardiovascular Science, University College London, London WC1E 6BT, United Kingdom

¹⁴Department of Biosciences, Biotechnologies and Environment, University of Bari Aldo Moro, via Orabona 4, 70125 Bari, Italy

¹⁵Guangdong Provincial Key Laboratory of Malignant Tumor Epigenetics and Gene Regulation, Guangdong-Hong Kong Joint Laboratory for RNA Medicine, Medical Research Center, Sun Yat-Sen Memorial Hospital, Sun Yat-Sen University, Guangzhou, 510120, China

¹⁶School of Chemistry and Biochemistry, Georgia Institute of Technology, Atlanta, GA, 30332, United States

¹⁷Department of Biomedical Data Science, Stanford University, Stanford, CA, 94305, United States

¹⁸Department of Computer Science, University of Milan, Milan, 20133, Italy

¹⁹Department of Molecular and Cellular Biology, Harvard University, Cambridge, MA 02138, United States

²⁰Department of Computer Science and Software Engineering, The University of Western Australia, Crawley, WA 6009, Australia

²¹Leibniz Institute DSMZ – German Collection of Microorganisms and Cell Cultures, Braunschweig, 38124, Germany

²²Department of Biochemistry and Functional Genomics, Université de Sherbrooke, Sherbrooke, Québec, J1E 4K8, Canada

²³Department of Integrative Biology, University of Texas, Austin, TX, 78712, United States

²⁴Sandia National Laboratories, Livermore, CA, 87185, United States

²⁵Australian Centre for Ecogenomics, School of Chemistry and Molecular Biosciences, The University of Queensland, Brisbane, QLD 4072, Australia

²⁶Riboscope Ltd, 23 King St, Cambridge CB1 1AH, United Kingdom

*To whom correspondence should be addressed. Email: agb@ebi.ac.uk

Abstract

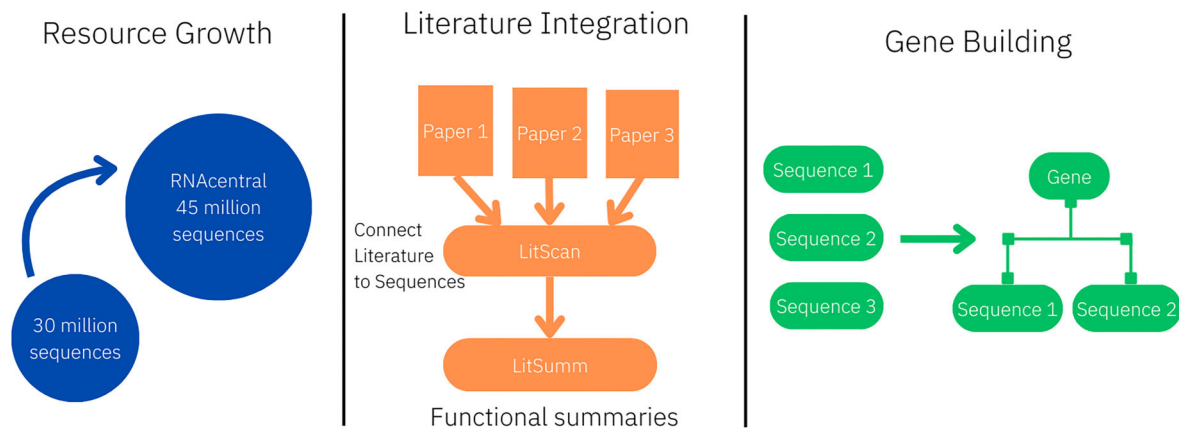
RNAcentral was founded in 2014 to serve as a comprehensive database of non-coding RNA sequences. It began by providing a single unified interface to more specialized resources and now contains 45 million sequences. It has grown beyond providing a single interface to many specialized resources and now provides several services and analyses. These include secondary structure prediction with R2DT, sequence search, and analysis with Rfam. Since its last publication in 2021, RNAcentral has developed two major features. First, literature integration with the development of LitScan and LitSumm. LitScan automatically identifies and links relevant publications to RNA entries, while LitSumm uses natural language processing to generate functional summaries from the literature. Together, these tools address the critical challenge of connecting sequence data with scattered functional knowledge across thousands of publications. Second, RNAcentral has created gene-level entries. Gene-level entries represent a large structural change to RNAcentral. While RNAcentral previously organized data exclusively at the sequence level, we now group related transcripts into gene-centric views. This allows researchers to explore all isoforms, splice variants, and related sequences for a gene in a unified interface, better reflecting biological organization and facilitating comparative analyses. RNAcentral is freely available at <https://rnacentral.org>.

Received: September 15, 2025. Revised: October 24, 2025. Accepted: October 27, 2025

© The Author(s) 2025. Published by Oxford University Press.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

Graphical abstract



Introduction

RNAcentral was established in 2014 as a comprehensive database unifying all non-coding RNA (ncRNA) sequences into a single, searchable resource [1]. This was aimed at solving the issue of fragmentation in the ncRNA field, where there were many high-quality databases that focused on specific subtypes of ncRNA. There was no easy-to-use resource that compiled all ncRNA data and RNAcentral sought to fill that gap. RNAcentral is designed around an “expert database” model, where each contributing database provides its expertise and data in a single area, and then RNAcentral combines, standardizes, and integrates these data into a single resource. Since its inception it has grown to contain over 45 million sequences and has moved from providing sequences and metadata to more complex data types, such as Gene Ontology (GO) annotations [2], Sequence Ontology (SO) terms [3], interaction data from resources such as IntAct [4] and RNA-KG [5, 6], disease associations from GeneCards [7], and analysis, such as Rfam and R2DT predictions [8]. RNAcentral remains the primary source of sequence data for ncRNA science, with tens of thousands of users and hundreds of citations per year.

Since its last publication in 2021 [8], RNAcentral has had 9 releases, grown from 30 to 45 million sequences, and now includes 52 expert databases, including 10 new databases and major updates to existing sources. We also expanded secondary structure predictions from 13 to 30 million sequences. Beyond these quantitative changes, we have also updated and improved how researchers can interact with our ncRNA data. We extended our search to allow searching by taxonomic descendants (e.g. all primates or all fungi) and transitioned to CC0 licensing for unrestricted data reuse enabling both academic and commercial reuse.

This article details the advances in RNAcentral release 26: new database integrations and major updates to existing resources; the implementation of our literature integration system (LitScan and LitSumm); the creation of gene-level entries; technical infrastructure improvements; and future development priorities.

Database content and growth

Since our 2021 publication, RNAcentral has grown from 30 to 45 million sequences (50% increase) across 52 expert databases. This expansion includes 10 new database integrations, major updates to three core resources, and the retire-

Table 1. A summary of the key features added since the last publication (releases 17–26)

Release	Sequences (millions)	Key features
17	30	Added piRBase, added Ribovore
18	30	Finalized improved rRNA RNA types
19	31	Added PSICQUIC, CPAT based ORF prediction
20	31	Migrated to CC0, Added RiboVision, Created LitScan
21	33	Added PlncDB, Added Expression Atlas, Added SwissBioPics
22	34	Added EVlncRNAs, Added Ribocentre, Integrated Expression Atlas Viewer
23	36	Added MGnify, Added REDiportal, Improved genome browser, Added LitSumm
24	35	Added tmRNA Website 2.0, Improved Phylogenetic search
25	45	Added Rfam 15.0, Improved taxonomic search
26	45	Created gene level entries, Added RNA-KG

ment of one database, VEGA. Table 1 summarizes the major changes across releases 17–26.

New database integrations

We integrated ten new expert databases to expand both taxonomic coverage and functional annotations. piRBase [9] contributes 200 000 piRNA sequences from 21 organisms for studying transposon silencing. PSICQUIC [10] provides 4420 manually curated RNA–protein interactions to expand our coverage of RNA interactions. PLncDB [11] adds plant lncRNAs from 80 species, while linking to EVlncRNAs [12] provides a connection to their experimentally validated lncRNAs from 19 000 publications with disease associations. Ribocentre and Ribocentre-switch [13, 14] provide structural data for ribozymes and riboswitches, respectively, enabling structure–function studies. MGnify [15] expands environmental coverage with sequences from over 28 000 metagenome-assembled genomes spanning 1132 isolates. RiboVision [16] provides templates to R2DT, allows secondary structure editing, and connects sequences to 3D structures and detailed annotations. REDiportal [17] adds RNA editing events for post-transcriptional modification analysis. Expression Atlas [18] provides tissue-specific expression profiles with interactive vi-

sualizations, expanding sequence data with functional information. RNA-KG [5, 6] integrates around 100 M interactions involving RNAs from 91 linked open data repositories and ontologies. Relationships are characterized by standardized properties that capture the specific context (e.g. cell line, tissue, pathological state) in which they have been identified.

A key database RNAcentral now incorporates is Expression Atlas [18]. This resource hosts gene and protein expression at bulk and single-cell levels for a range of species. The data provided is often at the tissue level. RNAcentral has imported the bulk expression data and mapped their gene-level entries to our sequence pages. We have also integrated the visualization provided by Expression Atlas to allow users to explore the expression of their sequence of interest.

Overall, RNAcentral increased the range of ncRNAs included, adding piRNAs with piRBase, and expanded to new organisms, e.g. plants and metagenomes from PLncDB and MGnify, but also enriched the data included by adding expression and modification data via Expression Atlas and REDIPortal. Users interested in exploring data from any of these resources can use our text search to filter by database.

Major database updates

In addition to adding new resources, we have also had several resources undergo major changes. First, Rfam updated to version 15.0 [19]. This was a major update to Rfam, which included updating the matches to all families. This led to increasing the size of Rfam from 1.88 M to 5.9 M sequences, with refreshed genomes and updated alignments for all RNA families, substantially improving family coverage and annotation quality. Second, TarBase is now at version 9.0 [20]. TarBase now contains experimentally validated miRNA–gene interactions from a range of species, including virally encoded miRNAs. Finally, we are now the primary site hosting Transfer-messenger (tmRNA) Website data [21]. The site has been discontinued, but the data has been updated in RNAcentral. This now expands RNAcentral to include an archival role for expert databases, where the data can continue to be publicly accessible after resources become unavailable. In release 24, we updated the data and included detailed annotations of sequence features within tmRNAs [22].

Enhanced analysis pipelines

Since release 17, RNAcentral has added two new analyses. First, we now use Ribovore to analyze all sequences from ENA [1] prior to import [23]. Ribovore is a tool that can quickly and accurately detect complete and fragmentary rRNAs. We found that the rate of fragmentary rRNAs was doubling the size of RNAcentral but provided little value to users as these are generally poorly annotated sequences. Thus, we sought to exclude them from future import and used Ribovore to detect such sequences. This does not prevent the import of complete rRNAs or non-metagenomic rRNAs.

In addition, we now analyze sequences from four organisms—human, fly, zebrafish, and mouse—with CPAT, which detects open reading frames [24]. Results indicate that ~6% of sequences contain detectable ORFs, providing quality metrics that help distinguish coding transcripts from legitimate ncRNAs. These sequences are given a QC warning as containing an open reading frame. Finally, we used the existing Rfam analysis pipeline to reanalyze all sequences with all Rfam families after each major release of Rfam. This will keep

the annotations up-to-date with the latest families and meta-data.

Literature integration: connecting RNA sequences with scientific knowledge

A longstanding challenge for RNAcentral has been providing users with comprehensive, up-to-date connections between RNA sequences and their functional descriptions in the scientific literature. To address this need, RNAcentral has developed two complementary tools to automatically extract and synthesize information from the vast corpus of RNA research: LitScan for literature mining and LitSumm for automated summarization [25] that work with any ncRNA in RNAcentral.

LitScan: mining the literature for RNA mentions

LitScan systematically identifies and extracts literature mentions of RNA identifiers in the open access literature. The system queries the EuropePMC [2] API using RNAcentral's comprehensive collection of identifiers and their aliases.

For each matching article, LitScan extracts the complete sentences with the mentioned RNA. It uses RNAcentral's cross-references to capture alternative names, e.g. the lncRNA THRIL is also known as 'Linc1992' and 'NR_110375.' Results are stored in RNAcentral's database and displayed through an interactive widget allowing filtering by journal, publication year, article section, and publication type (Fig. 1). Additionally, some expert databases provide annotations for key papers that discuss an RNA. RNAcentral captures this information and provides it as a filter in the viewer.


As of release 26, RNAcentral has processed 40 million papers with 13.6 million identifiers, resulting in 1.5 million ncRNA sequences connected to at least one paper. Most RNAs have very few papers with well-studied RNAs such as human XIST and MALAT1 appearing in over 5000 publications.

LitSumm: LLM-powered literature summarization

LitSumm addresses the challenge of synthesizing scattered functional knowledge by generating structured summaries from LitScan-extracted sentences. The system employs a multi-stage pipeline, where sentences from LitScan are selected via topic modeling, an initial summary is generated using ChatGPT4, and then checked for self-consistency and sufficient citations [25]. We then utilize GPT-4-turbo (gpt-4-1106-preview) to summarize and validate the summaries. Expert review of summaries showed that 94% of the summaries were rated good or excellent. Where the summary was rated poorly, the main reason was a failure to properly synthesize facts from multiple sources. This issue, which automated checking could not catch, has also been observed in other work [26]. The other most common issue was references being either irrelevant or assigned to the incorrect sentence. A detailed analysis of the limitations of LitSumm is given in [25].

Currently, LitSumm has generated summaries for ~4600 human ncRNAs prioritized by community interest (HGNC [3], miRBase [4], mirGeneDB [5], and snoDB [6] entries). For miRNAs, we restrict to species-specific identifiers (e.g. "hsa-mir-21" rather than "mir-21") to prevent cross-species confusion. The LitSumm summaries are displayed on each sequence page accompanied by a warning about AI-generated

Open Access Publications 429 total 

Sort by score - default  Text search within results

Keyword

- MT-RNR2 (362)
- MTRNR2 (45)
- ENSG00000210082 (14)
- ENST00000387347 (8)
- URS000047A7F4_9606 (0)

Article type

- Research article (361)
- Review article (57)
- Brief report (5)
- Case report (2)
- Other (2)
- Systematic review (1)
- Letter (1)

Paper section

- Title (1)
- Abstract (16)
- Main text (428)

Mentioned Organism

- Metazoa (143)

The miR-27a/FOXJ3 Axis Dysregulates Mitochondrial Homeostasis in Colorectal Cancer Cells
Barisciano, Giovannina et al., *Cancers*, 2021. Cited by 2 articles. [\[Europe PMC\]](#)

*Mitochondrial DNA content was evaluated with a commonly used method based on the mitochondrial to nuclear DNA (mtDNA/nDNA) ratio, in which we quantified the mitochondrial encoded genes tRNA^{Leu} and **MT-RNR2** versus nuclear encoded ones, Claudin1 and SOX9, by q-PCR [29].*

[View abstract](#)

Cell Type-Specific Predictive Models Perform Prioritization of Genes and Gene Sets Associated With Autism
Guan, Jinting et al., *Frontiers in Genetics*, 2021. [\[Europe PMC\]](#)

MT-RNR2 found in an image, table or supplementary material

[View abstract](#)

Isoform-resolved mRNA profiling of ribosome load defines interplay of HIF and mTOR dysregulation in kidney cancer
Sugimoto, Yoichiro et al., *Nature Structural & Molecular Biology*, 2022. [\[PubMed\]](#) [\[Europe PMC\]](#)

*The pre-processed reads were first mapped to cytoplasmic rRNAs (NR_023363.1 and NR_046235.1), mitochondrial ribosomal rRNAs (ENSG00000211459 and **ENSG00000210082**), and ERCC external control RNAs (https://www-s.nist.gov/srmors/certificates/documents/SRM2374_putative_T7_products_NoPolyA_v1.fasta) using Bowtie2 software (2.4.1)57 with the following parameters: -N 1-un-conc-gz.*

Figure 1. Image of an RNAcentral sequence page for human MT-RNR2 (URS000047A7F4_9606) with the LitScan widget displayed.

content and a link to provide feedback (Fig. 2). Entries with summaries are available on RNAcentral by searching for ‘has_litsumm:“True”’ (https://rnacentral.org/search?q=has_litsumm:%22True%22) or can be fetched at <https://huggingface.co/datasets/RNAcentral/litsumm-v1.5>, and the code is available at <https://github.com/RNAcentral/litscan-summazation>.

Development of LitSumm is ongoing, aiming to alleviate the limitations of the original. Among these is the sensitivity of the original system to mentions “in passing,” where a paper merely mentions an ncRNA but is not the main focus of the study. We are developing an agentic approach, including literature triaging and prioritization to filter out such mentions, and provide much better initial context for the summary. The second version of LitSumm will continue to prioritize information provenance above all else and is expected to enter production in mid-2026.

RNAcentral gene building

In release 26 of RNAcentral, we have created the first set of genes for all sequences that could be found in a genome, excluding piRNAs, in 204 organisms. Until now RNAcentral has been a sequence-based resource, which means each unique sequence is given a Unique RNA Sequence identifier (URS id) [1] and treated as a separate entry. There are many examples of sequences that may differ in only a few nucleotides. When these sequences are present in the genome, they often represent variants in sequence of the same transcript. For example, they include cases of one transcript a single nucleotide longer than another being treated as separate, unrelated entries. For genes that are commonly used for fingerprinting, such as rRNAs, there will be thousands of examples of sequences that dif-

fer slightly in overall length and have a few changes between them. Conversely, for some ncRNAs, such as miRNAs, there may be several identical copies in a genome, which are merged into a single URS entry. Both of these situations can also be confusing to users. Additionally, many biological experiments and data are at the gene, not transcript level. These issues have led us to build gene-level entries into RNAcentral.

We could not simply use gene objects from existing resources, such as Ensembl [8], because RNAcentral’s data includes many sequences that are absent from other collections. However, the goal was to produce gene objects that are comparable to those in Ensembl for all organisms in RNAcentral without human intervention. Finally, the gene-building pipeline must be able to accommodate the changes that occur between RNAcentral releases.

As each RNAcentral release may add thousands of new sequences to an organism the pipeline must be able to detect when two genes in two different releases are similar enough to be revised versions of the same gene. RNAcentral genes should have stable identifiers across releases even in the presence of changes. We note that in RNAcentral a single release adds around 1 million new sequences and may add several thousand new ones to highly studied organisms such as humans. The pipeline must tolerate large changes in the number of transcripts within a gene.

Gene building pipeline

The new RNAcentral gene building pipeline solves these problems using two major steps. First, for each release we run a graph clustering algorithm to build transcript clusters, and then we compare transcript clusters to the previous release to assign gene identifiers. Each transcript is treated as a node, and



Homo sapiens (human) mitochondrially encoded 16S rRNA (MT-RNR2)

URS000047A7F4_9606

1,559 nucleotides 6 databases (ENA, Ensembl, Expression Atlas, GeneCards, HGNC, MalaCards) Found in 0 other species 429 publications

rRNA / cytosolic rRNA / cytosolic LSU rRNA

Caution, this is an AI generated summary based on literature. This may have errors, see [here for more](#). Please share your [feedback](#) with us.

> **mt-rnr2**: MT-RNR2 encodes for 16S ribosomal RNA, a component of the mitochondrial ribosome, essential for protein synthesis within the mitochondria [PMC4357308]. This gene exhibits methylation differences across various tissues, with notably high methylation levels in the brain [PMC4357308]. While MT-RNR2 is associated with Alzheimer's disease, it is regulated by a set of microRNAs, including hsa-miR-mit3, hsa-miR-mit6, and hsa-miR-mit4; however, the association does not necessarily imply causation for the disease [PMC4324738]. The drug PU-91 upregulates MT-RNR2, which is associated with improved mitochondrial health and suggests a potential role in mitochondrial protection and cytoprotection in AMD ARPE-19 transmittochondrial cybrid cells [PMC6756897]. Furthermore, MT-RNR2 has isoforms classified as pseudogenes—MTRNR2L1 and MTRNR2L2—adding to the complexity of its gene family and suggesting possible functional implications [PMC9028063].

Figure 2. Image of an RNAcentral sequence page for human MT-RNR2 (URS000047A7F4_9606), with the LitSumm summary highlighted.

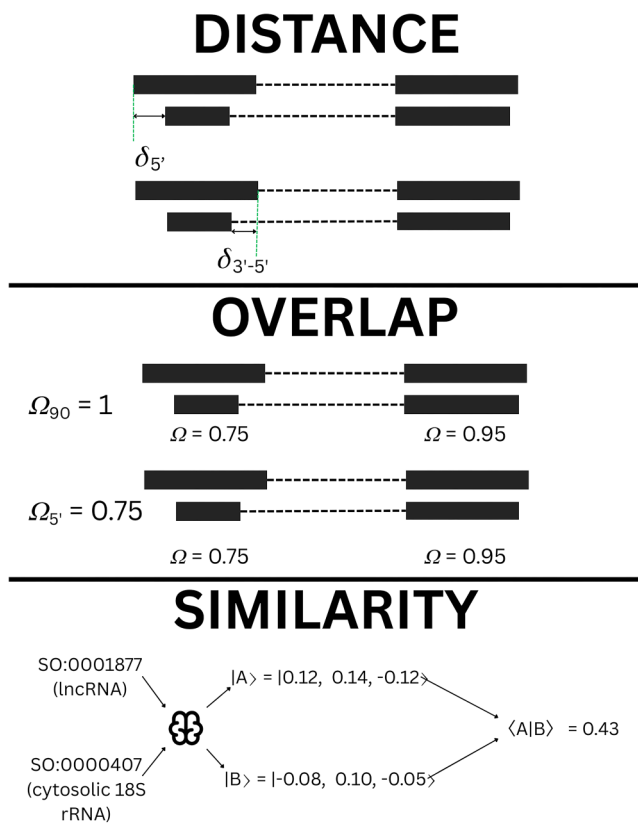


Figure 3. Illustration of feature calculation based on mapped transcript coordinates. Distance metrics primarily focus on the 5' exon and are simply the distance in number of nucleotides between the two transcripts. Overlap metrics consider the 5' exon as well as the global exon structure, measuring the overlap of the 5' exons of two transcripts and the number of exons within the transcripts having >90% overlap. Finally, the type similarity uses a Node2Vec model trained on the sequence ontology terms to convert types into normalized vectors, then calculates the inner product to give a numerical similarity between RNA types.

all pairs of transcripts whose start site is within 1 kb are compared using a random forest model to determine which pairs should be connected. We then find communities [27] in the resulting graph and treat those communities as genes. Communities were used to prevent the tendency of rare but large transcripts to join several otherwise distinct genes. We exclude

Table 2. Cross validation performance of the gene building model built and tested on human genes from Ensembl

Fold	Accuracy	F1	AUC	AP
1	0.968	0.989	0.994	0.998
2	0.970	0.989	0.995	0.998
3	0.970	0.989	0.995	0.998
4	0.968	0.989	0.994	0.998
5	0.968	0.989	0.994	0.997
Average	0.969	0.989	0.994	0.997

piRNAs from this gene-building effort, as these ncRNAs have a unique biogenesis pathway [28].

To guide the development of the random forest model, we asked ~30 consortium members to manually cluster 10–20 transcripts in three example regions. We then discussed with them why certain transcripts were clustered together and what features they considered. We used these features and trained a random forest on the existing Ensembl/GENCODE [9] transcripts from the human genome. We selected the human genome as this has undergone the most manual curation. We use three classes of measures: distance, overlap, and similarity. The distance measures the proximity of the start and end of the 5' most exon. For overlap we look at how many exons have >90% overlap and the percentage overlap of the 5' most exon. Finally, for similarity, we use node2vec [29] to compute similarity between the SO RNA types between transcripts as assigned by the source expert databases. Node2vec is a method that allows the calculation of vector similarity within a graph structure. As the SO is a directed acyclic graph this allows us to easily compute similarity between any pair of SO terms.

The random forest model is used to infer the probability that any given pair of transcripts come from the same gene. This pairwise comparison is then used to provide an edge weight between nodes in the graph, which is used when calculating the communities that become genes.

We used five-fold cross validation when training the procedure on human genes and found excellent performance, with an average Area Under the Receiver-Operator Curve (AUC) of 0.994 (Table 2)—where an AUC of 1.0 would be a perfect classifier. As this performance could be due to imperfect splits in the data, we applied the model to Ensembl annotations of a similar well-annotated species, mouse, and observed consistently high performance (average AUC 0.974, Table 3).

Table 3. Cross validation performance of the gene building model built on human genes but on mouse genes from Ensembl

Fold	Accuracy	F1	AUC	AP
1	0.887	0.959	0.975	0.989
2	0.887	0.959	0.975	0.989
3	0.887	0.959	0.975	0.990
4	0.888	0.960	0.976	0.990
5	0.888	0.960	0.969	0.980
Average	0.887	0.959	0.974	0.988

Inspecting the random forest model, we find the feature importances behave as expected from our discussions with consortium members (Fig. 4); the most important features are to do with the distances between 5' coordinates of the transcripts, both the simple distance between the two 5' transcript start coordinates, and the distance between the coordinates of the 3' end of the 5' exon. Transcripts being located on the same strand was obviously also important, but sequence ontology type similarity was less crucial than we expected. The relative unimportance of sequence type similarity may be a reflection of the training data used (coming from known genes) having relatively low diversity.

We then developed logic to merge clusters between releases. Because the build pipeline uses only data available from RNAcentral's GFF files, we were able to run it for all sequences from release 12, from September 2019, onward. This allowed us to inspect how different merging criteria would change the resulting genes. After inspecting the data, we found that treating clusters with a start site within 1 kb between two releases as being the same gene produced stable results.

Gene identifiers and metadata

RNAcentral genes are given accessions with the following pattern: `RNACG<species-prefix><11-digit hash>.<version>`.

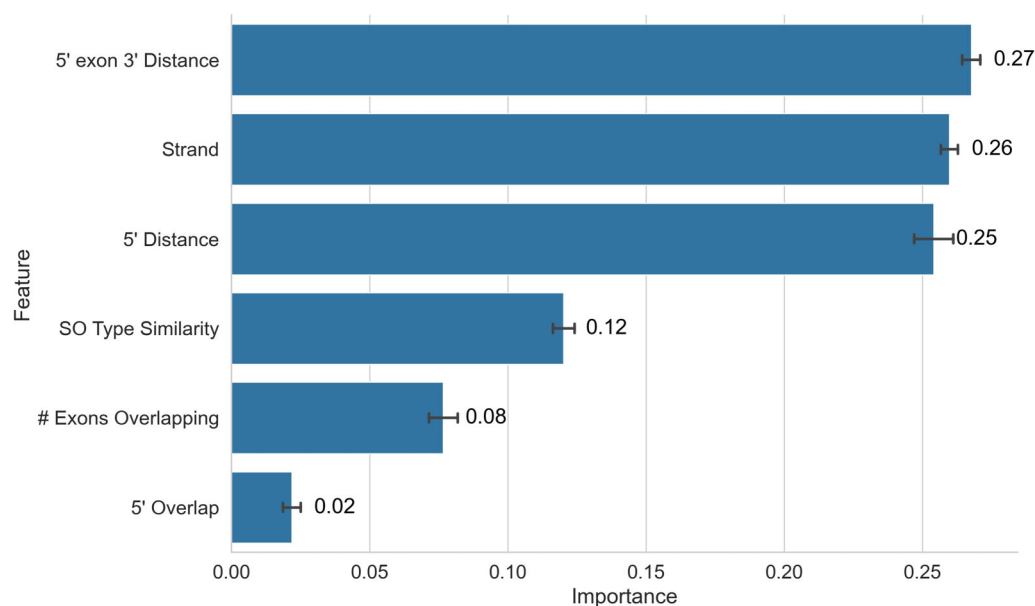


Figure 4. Feature importances from the random forest model, averaged across the five folds of cross-validation. Feature importances were calculated using the mean decrease in Gini impurity, aggregated over all trees in the forest, with confidence intervals derived from the standard deviation of importance across folds. These match what was expected from discussions with experts, with measures of distance at the 5' end of transcripts being most informative, alongside the transcripts being mapped to the same strand. SO type similarity was relatively unimportant for the model, possibly reflecting a lack of diversity in the types contained in the training data.

This allows tracking the organism, identity, and version of these genes over time. The species prefix is the same as the Ensembl species prefix, and for *Homo sapiens* there is no species prefix. The 11-digit hash is computed using SHA256 from the gene coordinates, chromosome, and assembly ID and is intended to be similar to gene names from Ensembl. The version number of each gene increments by one each time there is any change to the gene. These changes include adding or removing transcripts, or changing the start or end points of a gene. If no changes are made, then a gene will retain the same accession between releases. Some genes have remained unchanged since creation in release 12, while others have changed with every release.

Alongside the identifier, each gene is given an RNA type based on SO terms and a description. These are assigned from the RNA types and descriptions provided by expert databases and the Rfam and R2DT annotations. The R2DT template assignments supplemented Rfam family annotations to identify specific SO types, such as tRNA and rRNA types. We used a weighted voting system where data from curated resources such as HGNC [3], FlyBase [10], and GENCODE are weighted more than non-curated resources such as ENA. Similarly, we only consider results from Rfam or R2DT if they match over 90% of the longest sequence in the gene. Manual inspection of the gene names and types indicates this system performs well overall; however, we plan to continue to improve and extend the system to better deal with edge cases and uncertainty.

Current human gene set

RNAcentral contains 103 814 human ncRNA genes from 600 225 transcripts. These genes cover 56 different SO terms, with the most common being lncRNA, 65 187, and anti-sense lncRNA, 16 790, followed by pre-miRNA, 8560. The remaining 13 277 genes are SRP RNA, generic ncRNAs,

RNAcentral

(TAXONOMY:"562" OR lineage_path:"562")

Examples: human RN7SL, Rat, lncRNA, PDBe, ENSG00000228630

v26 Databases Tools API Downloads Browse About Help Feedback

Results 15 out of 374,833 sequences

Sort by: Popular species, Length ↓

Entry type

- Sequence (374,833)
- Gene (0)

RNA types

- ncRNA (373,385)
- rRNA (193,691)
 - cytosolic rRNA (87,953)
 - cytosolic LSU rRNA (56,401)
 - cytosolic 23S rRNA (19,533)
 - cytosolic 5S rRNA (11,755)
 - cytosolic 28S rRNA (41)
 - cytosolic 5 8S rRNA (1)
 - cytosolic SSU rRNA (31,552)

Escherichia coli str. K-12 substr. MG1655 23S ribosomal RNA

URS00019ABAEF_511145
2,933 nucleotides ENA

Escherichia coli str. K-12 substr. MG1655 23S ribosomal RNA

URS00019F16E6_511145
2,932 nucleotides ENA

Escherichia coli str. K-12 substr. MG1655 bacterial LSU rRNA

URS000077EFF8_511145
2,932 nucleotides ENA silva

Figure 5. An example of a search for all subspecies of *Escherichia coli* in RNAcentral.
[https://rnacentral.org/search?q=\(TAXONOMY:%22562%22%20OR%20lineage_path:%22562%22\)](https://rnacentral.org/search?q=(TAXONOMY:%22562%22%20OR%20lineage_path:%22562%22))

and a variety of specific terms, e.g. tRNA subtypes. The average ncRNA gene in RNAcentral contains 6 transcripts, with the minimum being 2 and maximum being 4001 for a mitochondrial small subunit rRNA with many fragments. Comparing the genes to Ensembl shows that, on average, RNAcentral genes map to 1.6 Ensembl genes, indicating our pipeline will, on average, merge genes that Ensembl considers separate.

We examined several specific cases, including MALAT1, NEAT1, XIST, and the clustering of a miRNA sequence URS000075CC93_9606. We found that both MALAT1 and NEAT1 were each built into a single gene as expected. On the other hand, the miRNA URS000075CC93_9606 represents several different miRNAs, hsa-mir-1302-2 and hsa-mir-1302 9 to 11. Ensembl and miRBase list this sequence as occurring in four different genes, and RNAcentral's gene building reflects that. On the other hand, some genes are not built correctly. For example XIST, ENSG00000229807, is split across five different genes in RNAcentral. This is due to XIST's complex splicing structure, where there is a large distance between possible 5' exons of XIST. In general, the RNAs that are not built correctly are lncRNAs with this complex alternative splicing pattern. We are looking into correcting this issue, as well as alternative approaches to solving this problem.

RNAcentral provides access to the genes in several ways. First, the text search now includes genes; users can select the "Genes" value in the Entry Type facet and see only gene-level results. Users are able to use any facet. Additionally, all sequence pages link to the genes, if any, the sequence is a part of. Finally, genes have been added to our GFF files as a "predicted_gene" entry.

Continued development

The gene pipeline performs well overall; however, there is scope for improvement. The identified feature importances indicate that the model will group ncRNA clusters into single genes, rather than as might be expected grouping each ncRNA into its own gene. ncRNA clusters occur for several types, including miRNA, tRNA, and snoRNA; we are working with type-specific expert databases to refine the gene model for these types. Another limitation arises for the "inside-out" genes such as SNHG1, where the ncRNA product is coded in an intron. For these RNAs, our model is unlikely to perform well, due to limitations in the training data. For training, we use genes and their transcripts from Ensembl, but these do not split ncRNAs in transcript introns into their own transcript; as a result the random forest model has not seen how to handle these correctly. We are working to improve the preprocessing of our training data to remove this limitation.

The set of genes identified by the system now can be considered a superset of the ncRNA genes that exist; this method has high recall, but relatively low precision. While we cannot know for sure which of the predicted genes are real and therefore cannot train directly to improve precision, we can apply further methods to improve precision by filtering out predictions that are demonstrably wrong. We are exploring the use of sequence homology within the predicted gene set to filter out genes we predict have no evidence of evolutionary conservation; we expect this additional filtering will remove many genes and improve the quality of those that remain. Alongside these post-prediction filtering steps, we are implementing measures to exclude transcripts with excess protein coding potential from being considered for gene building by incorporat-

ing additional QA steps at import, such as Pfam, stopFree [11], and tcode [11].

Technical infrastructure and tools

Since its last publication, RNAcentral has improved several key components: we extended phylogenetic search functionality, migrated to a modern genome browser, updated R2DT to version 2.1, integrated SwissBioPics visualization, and incorporated the Expression Atlas viewer.

Search and discovery enhancements

In release 24, RNAcentral extended phylogenetic search functionality to include subspecies filtering. Users can now toggle between searching only the selected species (e.g. *Escherichia coli*, NCBI taxonomy ID 562) or including all subspecies variants. Now, the user is able to select “Include subspecies” to also find subspecies of *E. coli*. An example of using this feature is shown in Fig. 5. This enhancement addresses user requests for more granular taxonomic control and has been applied to all taxonomic entries in RNAcentral.

Improved visualization

We have improved several aspects of our visualizations. First, our genome browser was migrated from Genoverse to igv.js (Fig. 6B) [30]. This modern genome browser allows users to track uploads, is faster, is less error-prone, and is easier to develop with.

R2DT version 2.1 includes several new templates, including mitochondrial tRNAs, tmRNAs, and new rRNA models [31]. As a result of the update, our coverage increased to 30 million sequences from 13 million previously, now covering 55% of sequences.

Third, we have added SwissBioPic integration, which is a visualization of the subcellular localization of molecules (Fig. 6C) [32]. The locations use the “located_in” Gene Ontology annotations imported into RNAcentral. This provides users a simple way to see where a molecule is found in the cell.

The Expression Atlas integration queries data for RNA expression across 4000+ studies and 60+ species as represented in Expression Atlas. Users can access tissue-specific and developmental stage expression data where available, with direct links to source experiments. The integration covers nearly 300 000 sequences in RNAcentral. This visualization is developed and maintained by the Expression Atlas team at EMBL-EBI. We reuse it to allow users easy access to browse and visualize the expression of RNAs in a well-tested and familiar tool, as shown in Fig. 6A.

API and programmatic access

RNAcentral continues to provide programmatic access to our data and services. As previously, we have an FTP site with data dumps, an API for programmatic access, and a publicly accessible Postgres database. Since our last publication, we have added a Swagger UI to our API documentation at <https://rnacentral.org/api/schema/swagger-ui/>. We note that LitSumm summaries are also now available via the API. We now provide complete PostgreSQL database dumps via `pg_dump`, updated each release and compressed to ~42 Gb. These dumps enable local installations and high-throughput analyses without API rate limits. We would like to emphasize that in the future it may be necessary to remove old dumps if

the overall size grows too large; we will only commit to providing the previous release’s data dump. We encourage any users who are interested in our data at scale, particularly those who scrape our pages to use this resource.

Additionally, our services provide APIs for programmatic use. For example, sequence search and R2DT both allow automated submissions. Interested parties may find the documentation in our help pages.

Finally, RNAcentral provides a series of embeddable widgets for our services, including the sequence search, R2DT, and LitScan. These widgets have been embedded into several external resources, including Rfam (LitScan, sequence search), FlyBase (R2DT), GtRNAdb (sequence search), NAKB (R2DT), PomBase (R2DT), Ribocentre, Ribocentre-switch (sequence search), and SGD (R2DT). Documentation for how to integrate resources is available in our help pages, and we invite any interested resources to contact us (<https://rnacentral.org/contact>) to discuss integrating our widgets into their site.

Licensing and data accessibility

In release 20, RNAcentral transitioned to a CC0 license. The CC0 license allows reuse without any attribution and is considered best practice for biological knowledgebases. We welcome anyone to reuse our data as they see fit, as we place our data into the public domain as much as possible. While CC0 licensing requires no attribution, we encourage citations when RNAcentral data contributes to published research, as these metrics support continued funding and development.

Future directions

RNAcentral plans to continue its growth in sequences and number of data types. Our focus in the near future is to integrate more useful data types, provide more tools for community use, and continue to develop our gene model and enhance our gene-level annotations. To support this, we will expand the quality assurance pipeline, applying additional testing for coding potential using tools such as Pfam [12] at each release. Data we are particularly interested in is RNA modifications, and so we plan to integrate thousands of nCRNA transcripts with information about posttranscriptionally modified residues from MODOMICS [33], following the integration of high-throughput modification mapping data from the Sci-ModoM database [34]. One such tool will be an MCP server for LLM access to RNAcentral data. Additionally, we plan to extend LitSumm summaries to all sequences with publications. Users interested in particular data types or researchers who are interested in collaborating are encouraged to contact us at <https://rnacentral.org/contact> and discuss their needs.

Acknowledgements

Author contributions: All RNAcentral consortium members represented here contributed to the development of their respective databases and have provided new data and features to RNAcentral since last publication. RNAcentral has been prepared by Andrew Green, Carlos Eduardo Ribas, Isaac Jandalala, Philippa Muston (RNAcentral), Colman O’Cathail, Guy Cochrane (ENA), Christina Ernst, Lingyun Zhao, Pedro Madrigal (Expression Atlas), Helen Attrill, Steven Marygold (FlyBase), Doron Lancet, Niv Dobzinski (GeneCards), Patricia P. Chan, Todd M. Lowe (GtRNAdb), Elspeth A. Bruford,

Ruth L. Seal (HGNC), Henning Hermjakob, Kalpana Panneerselvam (IntAct), Robert D. Finn, Tatiana A. Gurbich (MGnify), Sam Griffiths-Jones (miRBase), Bastian Fromm, Kevin J. Peterson (MirGeneDB), Dominik Sordyl, Janusz M. Bujnicki (MODOMICS), Sameer Velankar, Sri Devan Appasamy, Sudakshina Ganguly (PDBe), Peng Zhang, Shunmin He (piRBase), Kim M. Rutherford, Valerie Wood (PomBase), Ruth C. Lovering (QuickGO), Ernesto Picardi (REDportal), Nancy Ontiveros (Rfam), Lin Huang, Zhichao Miao (Ribocentre), Anton S. Petrov, Holly McCann (RiboVision), Emanuele Cavalleri, Marco Mesiti (RNA-KG), Elena Rivas, Marcell Szikszai, Marcin Magnus (RNA3DB), Jan Gerken, Maria Chuvpochina (SILVA), Danny Bergeron, Michelle Scott (snoDB), Kelly P. Williams (tmRNA Website), Robin R. Gutell, Cheong Xin Chan (CRW), Mark Quinton-Tulloch, Stavros Diamantakis (WormBase), Anton I. Petrov (R2DT), Alex Bateman, Blake A. Sweeney (RNAcentral)

Conflict of interest

A.B. is a member of the Nucleic Acids Research Editorial Board. J.B. is an editor of Nucleic Acids Research.

Funding

This work was supported by Biotechnology and Biological Sciences Research Council (BBSRC) [BB/J019231/1, BB/J019232/1, and BB/N019199/1] and Wellcome Trust [218302/Z/19/Z and WT/Biomedical Resource: 218302/A/19/Z]. Anton S. Petrov was supported by NASA Grant 80NSSC24K0344 and NSF Grant 2243706. Funding to pay the Open Access publication charges for this article was provided by EMBL-EBI's.

Data availability

All data are freely available at <https://rnacentral.org>. The data can be accessed in the FTP archive, as well as through an API and a public Postgres database (see <https://rnacentral.org/help> for instructions) under the CC0 license. The code is available at <https://github.com/rnacentral> under the Apache 2.0 license.

References

1. RNAcentral Consortium RNAcentral: an international database of ncRNA sequences. *Nucleic Acids Res* 2015;43:D123–9. <https://doi.org/10.1093/nar/gku991>
2. Ontology Consortium Gene, Aleksander SA, Balhoff J *et al.* The gene ontology knowledgebase in 2023. *Genetics* 2023;224:iyad031.
3. Eilbeck K, Lewis SE, Mungall CJ *et al.* The Sequence Ontology: a tool for the unification of genome annotations. *Genome Biol* 2005;6:R44. <https://doi.org/10.1186/gb-2005-6-5-r44>
4. Del Toro N, Shrivastava A, Ragueneau E *et al.* The IntAct database: efficient access to fine-grained molecular interaction data. *Nucleic Acids Res* 2022;50:D648–53. <https://doi.org/10.1093/nar/gkab1006>
5. Cavalleri E, Cabri A, Soto-Gomez M *et al.* An ontology-based knowledge graph for representing interactions involving RNA molecules. *Sci Data* 2024;11:906. <https://doi.org/10.1038/s41597-024-03673-7>
6. Cavalleri E, Perlasca P, Mesiti M. RNA-KG v2.0: an RNA-centered knowledge graph with properties. arXiv:2508.07427, 10 August 2025.
7. Bashir Ruth Fishilevich Simon, Iny-Stein Tsippi *et al.*, GeneCaRNA: A Comprehensive Gene-centric Database of Human Non-coding RNAs in the GeneCards Suite. *Journal of Molecular Biology* 2021;433. <https://doi.org/10.1016/j.jmb.2021.166913>
8. Consortium RNAcentral. RNAcentral 2021: secondary structure integration, improved sequence search and new member databases. *Nucleic Acids Res* 2021;49:D212–20. <https://doi.org/10.1093/nar/gkaa921>
9. Wang J, Shi Y, Zhou H *et al.* piRBase: integrating piRNA annotation in all aspects. *Nucleic Acids Res* 2022;50:D265–72. <https://doi.org/10.1093/nar/gkab1012>
10. Aranda B, Blankenburg H, Kerrien S *et al.* PSICQUIC and PSISCORE: accessing and scoring molecular interactions. *Nat Methods* 2011;8:528–9. <https://doi.org/10.1038/nmeth.1637>
11. Jin J, Lu P, Xu Y *et al.* PLncDB V2.0: a comprehensive encyclopedia of plant long noncoding RNAs. *Nucleic Acids Res* 2021;49:D1489–95. <https://doi.org/10.1093/nar/gkaa910>
12. Zhou B, Ji B, Shen C *et al.* EVLncRNAs 3.0: an updated comprehensive database for manually curated functional long non-coding RNAs validated by low-throughput experiments. *Nucleic Acids Res* 2024;52:D98–106. <https://doi.org/10.1093/nar/gkad1057>
13. Deng J, Shi Y, Peng X *et al.* Ribocentre: a database of ribozymes. *Nucleic Acids Res* 2023;51:D262–8. <https://doi.org/10.1093/nar/gkac840>
14. Bu F, Lin X, Liao W *et al.* Ribocentre-switch: a database of riboswitches. *Nucleic Acids Res* 2024;52:D265–72. <https://doi.org/10.1093/nar/gkad891>
15. Gurbich TA, Almeida A, Beracochea M *et al.* MGnify Genomes: a resource for biome-specific microbial genome catalogues. *J Mol Biol* 2023;435:168016. <https://doi.org/10.1016/j.jmb.2023.168016>
16. McCann HM, Meade CD, Banerjee B *et al.* RiboVision2: a web server for advanced visualization of ribosomal RNAs. *J Mol Biol* 2024;436:168556. <https://doi.org/10.1016/j.jmb.2024.168556>
17. D'Addabbo P, Cohen-Fultheim R, Twersky I *et al.* REDportal: toward an integrated view of the A-to-I editing. *Nucleic Acids Res* 2025;53:D233–42.
18. Madrigal P, Thanki AS, Fexova S *et al.* Expression Atlas in 2026: enabling FAIR and open expression data through community collaboration and integration. *Nucl. Acids Res.* 2025. <https://doi.org/10.1093/nar/gkaf1238>
19. Ontiveros-Palacios N, Cooke E, Nawrocki EP *et al.* Rfam 15: RNA families database in 2025. *Nucleic Acids Res* 2025;53:D258–67. <https://doi.org/10.1093/nar/gkae1023>
20. Skoufos G, Kakoulidis P, Tastsoglou S *et al.* TarBase-v9.0 extends experimentally supported miRNA-gene interactions to cell-types and virally encoded miRNAs. *Nucleic Acids Res* 2024;52:D304–10. <https://doi.org/10.1093/nar/gkad1071>
21. Hudson CM, Williams KP. The tmRNA website. *Nucleic Acids Res* 2015;43:D138–40. <https://doi.org/10.1093/nar/gku1109>
22. Nawrocki EP, Petrov AI, Williams KP. Expansion of the tmRNA sequence database and new tools for search and visualization. *NAR Genom Bioinform* 2025;7:lqaf019. <https://doi.org/10.1093/nargab/lqaf019>
23. Schäffer AA, McVeigh R, Robbertse B *et al.* Ribovore: ribosomal RNA sequence analysis for GenBank submissions and database curation. *BMC Bioinformatics* 2021;22:400. <https://doi.org/10.1186/s12859-021-04316-z>
24. Wang L, Park HJ, Dasari S *et al.* CPAT: coding-potential assessment tool using an alignment-free logistic regression model. *Nucleic Acids Res* 2013;41:e74. <https://doi.org/10.1093/nar/gkt006>
25. Green A, Ribas CE, Ontiveros-Palacios N *et al.* LitSumm: large language models for literature summarization of noncoding

- RNAs. *Database (Oxford)* 2025;2025:baaf006. <https://doi.org/10.1093/database/baaf006>
26. DeYoung J, Martinez SC, Marshall IJ *et al.* Do multi-document summarization models synthesize? arXiv:2301.13844, 31 January 2023.
 27. Blondel VD, Guillaume J-L, Lambiotte R *et al.* Fast unfolding of communities in large networks. *J Stat Mech* 2008;2008:P10008.
 28. Zhang J, Chen S, Liu K. Structural insights into piRNA biogenesis. *Biochim Biophys Acta Gene Regul Mech* 2022;1865:194799. <https://doi.org/10.1016/j.bbagr.2022.194799>
 29. Grover A, Leskovec J. Node2vec: scalable feature learning for networks. arXiv:1607.00653, 3 July 2016. <https://doi.org/10.48550/ARXIV.1607.00653>
 30. Robinson JT, Thorvaldsdottir H, Turner D *et al.* igv.js: an embeddable JavaScript implementation of the Integrative Genomics Viewer (IGV). *Bioinformatics* 2023;39:btac830. <https://doi.org/10.1093/bioinformatics/btac830>
 31. McCann H, Meade CD, Williams LD *et al.* R2DT: a comprehensive platform for visualizing RNA secondary structure. *Nucleic Acids Res* 2025;53:gkaf032. <https://doi.org/10.1093/nar/gkaf032>
 32. Le Mercier P, Bolleman J, de Castro E *et al.* SwissBioPics-an interactive library of cell images for the visualization of subcellular location data. *Database (Oxford)* 2022;2022:baac026. <https://doi.org/10.1093/database/baac026>
 33. Cappannini A, Ray A, Purta E *et al.* MODOMICS: a database of RNA modifications and related information. 2023 update. *Nucleic Acids Res* 2024;52:D239–44. <https://doi.org/10.1093/nar/gkad1083>
 34. Boileau E, Wilhelmi H, Busch A *et al.* Sci-ModoM: a quantitative database of transcriptome-wide high-throughput RNA modification sites. *Nucleic Acids Res* 2025;53:D310–7. <https://doi.org/10.1093/nar/gkae972>