

InteractOA: Showcasing the representation of knowledge from scientific literature in Wikidata

Muhammad Elhossary^a and Konrad U. Förstner^{a,b,*}

^a *Data Science and Services, ZB MED - Information Centre for Life Sciences, NRW, Germany*
E-mail: elhossary@zbmed.de

^b *Institute of Information Science, TH Köln – University of Applied Sciences, NRW, Germany*
E-mail: foerstner@zbmed.de

Abstract. Knowledge generated during the scientific process is still mostly stored in the form of scholarly articles. This lack of machine-readability hampers efforts to find, query, and reuse such findings efficiently and contributes to today's information overload. While attempts have been made to semantify journal articles, widespread adoption of such approaches is still a long way off. One way to demonstrate the usefulness of such approaches to the scientific community is by showcasing the use of freely available, open-access knowledge graphs such as Wikidata as sustainable storage and representation solutions. Here we present an example from the life sciences in which knowledge items from scholarly literature are represented in Wikidata, linked to their exact position in open-access articles. In this way, they become part of a rich knowledge graph while maintaining clear ties to their origins. As example entities, we chose small regulatory RNAs (sRNAs) that play an important role in bacterial and archaeal gene regulation. These post-transcriptional regulators can influence the activities of multiple genes in various manners, forming complex interaction networks. We stored the information on sRNA molecule interaction taken from open-access articles in Wikidata and built an intuitive web interface called *InteractOA*, which makes it easy to visualize, edit, and query information. The tool also links information on small RNAs to their reference articles from PubMed Central on the statement level. *InteractOA* encourages researchers to contribute, save, and curate their own similar findings. *InteractOA* is hosted at <https://tools.wmflabs.org/interactoa> and its code is available under a permissive open source licence. In principle, the approach presented here can be applied to any other field of research.

Keywords: Wikidata, interactions, regulatory networks, citations

1. Background and related work

1.1. Knowledge graphs and Wikidata

The term “knowledge graph” was coined by Google in 2012¹. Since then, knowledge graphs have attracted growing attention from researchers due to their robustness in many areas of science, besides application in numerous other fields [1]. Although several attempts have been made to define a knowledge graph, there is still no single, agreed-upon definition of what it entails [2]. As the name implies, a knowledge graph, also known as a semantic

*Corresponding author. E-mail: foerstner@zbmed.de.

¹<https://www.blog.google/products/search/introducing-knowledge-graph-things-not/>

1 network, is a means of storing knowledge in a graph-based model. It is a structured data model that represents a
2 network of real-world entities such as objects or concepts and illustrates the relationships between them. This infor-
3 mation is usually stored in a graph database, and can be visualized in a graph structure. Knowledge graphs can be
4 viewed as a network of nodes and edges, where nodes represent the entities and edges represent the relationships.
5 They can store vast amounts of heterogeneous data in a structured manner and handle complex relationships, mak-
6 ing them well suited for applications in various domains and research areas [3].
7

8 The Resource Description Framework (RDF) [4] is one model for implementing knowledge graph storage. The
9 RDF specification² highlights the simplicity of the RDF data model, which is represented at its fundamental level
10 in the form of subject-predicate-object triples. These triples can describe anything in a flexible and extensible way.
11 To query these RDF triples stores, a query language known as SPARQL, which is similar to SQL, is used to retrieve
12 and manipulate data stored in an RDF format.
13

14 Wikidata is one of the prime examples of knowledge graphs, and it combines two of their potential benefits: open-
15 ness (as the content is published under the Creative Commons Zero licence³) and ease of editability [5, 6]. It is based
16 on the Wikibase software, which allows RDF exports [7], and its content can be queried via the SPARQL endpoint
17 known as the *Wikidata Query Service*. Wikidata is maintained by the Wikimedia Foundation, which aims to provide
18 structured data about the world’s knowledge and make it available for anyone to use and extend collaboratively. A
19 number of projects are underway to analyse and increase the quality of Wikidata [8]. It is regarded as a key source
20 of identifiers [9] and has tremendous potential which remains largely untapped [10].
21

22 1.2. Insufficient management of data, information, and knowledge in the life sciences

23

24 Data, information, and knowledge are being generated at an ever-increasing pace in the field of biology. This is
25 due in large part to the widespread availability of high-throughput platforms in biological research, which can gen-
26 erate vast amounts of data, for example on genes, proteins, and other biological entities and their interactions. The
27 ability to collect, organize, and analyse this data—and the information and knowledge derived from it—is crucial
28 for future biological and biomedical research. While the FAIR principles are now widely acknowledged as a useful
29 framework for managing data [11], a large fraction of the knowledge generated on the basis of this data continues to
30 be stored in unstructured formats such as scholarly articles, which have formed the core of knowledge management
31 in research for centuries. Moreover, a significant proportion of these articles are not available in a machine-readable
32 formats (e.g. PDF), which in turn hinders the downstream information mining. Even the machine-readable formats
33 (e.g. in HTML, and XML) continue to lack the semantic enrichment. Attempts have been made to address this
34 [12, 13], but these have been largely ignored by the publishing industry and scholarly community. Recent projects
35 such as the Open Research Knowledge Graph (ORKG) [14] have tried to overcome this hurdle by offering a se-
36 mantic database to represent the research findings of publications in a separate machine-actionable form but are
37 referencing only on the article level and not to exact passages. Furthermore, domain-specific biological knowledge
38 graphs instances with SPARQL endpoints are also available such as UniProt [15] which is dedicated for storing
39 information of proteins and IntAct [16] for molecular interactions – both also referencing only on the article level.
40

41 Alongside their unstructured representation in academic literature, data, information, and knowledge are also
42 stored in biological databases. Even though these databases represent a valuable resource for larger or more spe-
43 cialized research communities, their development and maintenance is often limited to the lifetime of the respective
44 research project. There are numerous cases of valuable databases that were created as part of such projects, but
45 are no longer accessible. Examples of databases that cannot be accessed via their published URLs at the time of
46 writing this article include BSRD [17] and sRNAdb [18], which are databases for bacterial small RNAs, as well as
47 AANT [19], a database for amino acid-nucleotide interactions, and cpnDB [20], a database for bacterial Chaperonin
48 sequences. Several studies have examined this issue. Their results show that more than 30% of bioinformatical web
49

50 ²<https://www.w3.org/TR/2004/REC-rdf-concepts-20040210/>

51 ³<https://creativecommons.org/publicdomain/zero/1.0/>

1 services published over the last 23 years are currently unavailable, with lack of maintenance being the primary cause
2 of this decay [21–24]. Additionally, restrictive licensing can make it difficult for researchers to access and (re-)use
3 these resources, further hindering scientific advancement.

4
5 Another challenging aspect of knowledge management in research is the granularity of references. Claims are
6 cited on the level of full articles, which makes it very time-consuming for readers to find the exact location of a
7 particular statement in the referenced source. This poses a major obstacle to the verification and contextualization
8 of such references by readers.

9 10 *1.3. Wikidata as a knowledge-graph solution for biological data, information, and knowledge*

11
12 In the life sciences and in other research fields, knowledge graphs such as Wikidata are used to represent and
13 integrate information from a variety of sources, including genomic data, literature, and experimental results [25–
14 27]. Knowledge graphs can be employed to model complex biological systems and processes, and to facilitate data
15 mining and analysis. This may include the representation and integration of genomic data such as genes, proteins,
16 and pathways. Furthermore, they can be used to model the interactions of biomedical entities, for example linking
17 genes associated with antibiotic resistance in a pathogen [28].

18
19 Widespread use of knowledge graphs in the life sciences would facilitate the discovery of new biological insights
20 and relationships through data mining and visualization, and improve the interpretation and understanding of bi-
21 ological data through the integration of diverse data sources. One example of the implementation of knowledge
22 graphs for biological data is the Clinical Knowledge Graph (CKG) [29]. The CKG aims to assist in the delivery
23 of personalized medical treatment by using machine-learning methods to mine data from heterogeneous domains.
24 Wikidata pushes these capabilities further by facilitating the sharing and reuse of biological data through the use of
25 standardized data formats and open-access solutions. Further examples of existing solutions include WikiGenomes
26 [30], an openly editable knowledge graph for genomic annotations that is geared towards the molecular biology
27 community, ChlamBase [31], which is a central access point for genomic and proteomic information specifically
28 for the Chlamydia research community, and WikiPathways [32], an open, collaborative, community-based platform
29 dedicated to the curation of biological pathways. A number of articles have been published that encourage the use
30 of Wikidata-based solutions in biology, such as the Gene Wiki initiative [33]

31 32 *1.4. Small RNA regulatory networks*

33
34 In bacteria and archaea, gene expression is controlled by a variety of regulators. One class of regulators is the
35 small RNAs (also known as non-coding RNAs) [34], which are not translated into proteins to perform their regula-
36 tory functions. This class of RNAs is responsible for vital regulatory roles in gene expression. These small RNAs
37 are often expressed in their hundreds to control cellular functions such as the response to environmental changes,
38 and each small RNA can influence the activity of multiple targets of proteins or messenger RNAs. Small RNAs reg-
39 ulate their targets through various mechanisms [35, 36] such as down-regulation (by disrupting mRNA translation
40 through base-pairing to the ribosomal binding sites) and up-regulation (by inhibiting mRNA degradation). Despite
41 their importance in all known bacterial and archaea species, knowledge of small RNAs is comparatively limited,
42 and they often fail to be included in the creation of holistic models in systems biology.

43
44 The genomic locations of small RNAs are first identified through experiments and annotated in genomic reference
45 sequences; then, the interactions between these annotated small RNAs and genes are computationally predicted and
46 experimentally confirmed. The numerous interactions at the cellular level that occur under different environmental
47 conditions between different regulators, such as small RNAs and their target genes, can be represented as a network.
48 This network consists of nodes such as regulators and gene targets, with the edges between these nodes representing
49 the identified interactions. Information on these interactions can be represented and visualized in network graphs
50 to facilitate understanding of their complexity. Typically, researchers report their findings on bacterial small RNA
51 regulatory networks in various, usually unstructured formats that lack a consistent standard. This data is presented

1 in the main body of articles, in supplementary materials such as spreadsheets, or in flat files. For a limited number of
2 species, data is manually compiled by experts into web-based databases such as RegulonDB [37]. Without a uniform
3 format for reporting such information, it is difficult to gain a comprehensive understanding of these regulatory
4 networks.

7 2. Approach and Implementation

9 2.1. *InteractOA as a Wikidata-based application for small RNA regulatory networks*

11 With these needs and technological foundation in mind we modelled small RNA interactions in Wikidata and
12 built an application called *InteractOA* (OA for Open Access as is builds on Open Access scholarly articles) that
13 presents the data. The application features a graphical interface that allows users to query Wikidata for a chosen
14 organism and displays the network of interactions or citation information at the statement level for all referenced in-
15 teractions. The semantic nature of Wikidata, combined with its openness, can help to overcome many challenges in
16 the management of data, information, and knowledge in gene regulation research. Due to its easy-to-use web inter-
17 face, permissive license (CC0), already large corpus of data (including literature metadata) to which other data can
18 be connected to, the availability of a SPARQL endpoint for querying data, a platform to create applications building
19 upon the data (Toolforge⁴) as well as a strong (not only scientific) community backed by a sustainable organisation
20 (Wikimedia Foundation) Wikidata was chosen as platform to implement this solution to compile extracted scholarly
21 knowledge.

23 The Wikidata based approach of linking biological entities in order to represent regulatory networks is superior
24 to other methods in terms of both its openness and sustainability. Since intracellular interactions are often complex
25 and numerous, it is important to use a digital representation approach that is scalable, extensible, easily queryable,
26 and usable. Wikidata possesses all of these properties.

29 2.2. *Modelling small RNA interaction data in Wikidata*

31 Typically, data related to genomes and genes including small RNAs are stored in text files such as GFF3 (General
32 Feature Format)⁵ files. These files can be obtained from public repositories like NCBI RefSeq [38]. Each line
33 within a GFF file contains details about a single annotation, including the gene's location on the genome, its type,
34 function, and other attributes. To model small RNAs and their interactions, we translated each gene annotation in
35 these GFF files into Wikidata items. These items were then interconnected using specific properties to represent the
36 interactions. Selected, widely used bacterial strains were used for this modelling process due to the abundance of
37 their small RNA interaction data, including *Escherichia coli* strain *K-12 substr. MG1655* (NCBI genome assembly
38 [GCF_000005845.2](https://www.ncbi.nlm.nih.gov/assembly/GCF_000005845.2), RefSeq accession [NC_000913.3](https://www.ncbi.nlm.nih.gov/RefSeq/NC_000913.3)).

40 For the processing and integration of the data an ETL (extract-transform-load) workflow was established (see
41 Figure 1). The first step was to obtain genomic annotation data for these organisms from the NCBI RefSeq repos-
42 itory. Next, a Python tool was developed to automate the process of importing annotation records from GFF files
43 and linking them to Wikidata entries using defined IDs corresponding to specific strains. To do this, the tool extracts
44 information from each annotation record, generates Wikidata item definitions (including name, synonyms, and de-
45 scription), and handles communication with Wikidata's API to create new non-redundant Wikidata items. These
46 newly created Wikidata items were then expanded by importing the remaining annotation data, including entry type
47 (e.g. protein or RNA), genomic location, and external gene identifiers. This process was based on selected Wiki-
48 data items and properties defined for the biology domain, which were collected from Wikidata's listings for life

50 ⁴<https://wikitech.wikimedia.org/wiki/Portal:Toolforge>

51 ⁵<http://www.ensembl.org/info/website/upload/gff3.html>

Table 1
Table of properties and items used in the modelling of the annotations, interactions, and citations of small RNA interactions.

ID	Name	Type	Usage / Description
P703	found in taxon	Property	used to link entities of annotations to a certain Wikidata item that represents an organism at the strain level
P31	instance of	Property	assigns the type of GFF entry e.g. gene or ncRNA to a Wikidata item
P644	genomic start	Property	assigns the genomic start location of an annotation entry to a Wikidata item
P645	genomic end	Property	assigns the genomic end location of an annotation entry to a Wikidata item
P2548	strand orientation	Property	assigns the genomic strand of an annotation entry to a Wikidata item
P688	encodes	Property	to link a gene (Wikidata item) to its product (Wikidata item), e.g. protein, or RNA
P702	encoded by	Property	vice versa of P688
P361	part of	Property	used to describe a partial product of a gene, like genes that splice to multiple mRNAs
P527	has part(s)	Property	vice versa of P361
P351	Entrez Gene ID	Property	assigns an identifier for an annotation entry originated from NCBI Entrez database
P2249	RefSeq genome ID	Property	assigns an identifier qualifier the start, end, strand of an annotation entry originated from NCBI RefSeq database
P2393	NCBI locus tag	Property	assigns an identifier for an annotation locus tag originated from NCBI
P637	RefSeq protein ID	Property	assigns an identifier for protein GFF entry originated from by NCBI
P128	regulates (molecular biology)	Property	links 2 annotation entities based on the interaction of any type if type is unspecified
P3777	antisense inhibitor of	Property	similar to P128, when interaction type is known as inhibition by antisensing
P3771	activator of	Property	similar to P128, when interaction type is known as up-regulation by activation
P3774	blocker of	Property	similar to P128, when interaction type is known as mRNA translation blocking
P3773	antagonist of	Property	similar to P128, when interaction type is known as antagonizing
P3772	agonist of	Property	similar to P128, when interaction type is known as agonizing
P248	stated in	Property	used to add reference about an interaction claim, which is a Wikidata item that represents a scholarly article
P1683	quotation	Property	used to add quote statement to the reference about an interaction claim
P932	PMCID	Property	used to add the PubMed Central identifier for the article of reference about an interaction claim
Q427087	non-coding RNA	Item	class of RNA that is not translated into proteins
Q423832	antisense RNA	Item	RNA molecules hybridizing to complementary sequences in either RNA or DNA, altering the function of the latter
Q201448	transfer RNA	Item	adaptor molecule composed of RNA
Q285904	Transfer-messenger RNA	Item	bifunctional RNA that has properties of a tRNA and an mRNA
Q424665	signal recognition particle	Item	protein-RNA complex facilitating translocation of proteins across membranes
Q1012651	ribonuclease P activity	Item	catalysis of the endonucleolytic cleavage of RNA, removing 5' extra nucleotides from tRNA precursor.
Q11053	RNA	Item	family of large biological molecules
Q7187	Gene	Item	basic physical and functional unit of heredity
Q22809680	forward strand	Item	forward oriented strand in a double-stranded DNA molecule
Q22809711	reverse strand	Item	reverse oriented strand in a double-stranded DNA molecule
Q215980	ribosomal RNA	Item	RNA component of the ribosome, essential for protein synthesis in all living organisms
Q277338	pseudogene	Item	functionless relative of a gene

sciences⁶. See Table 1 for a list of items and properties used here.

⁶https://www.wikidata.org/wiki/Wikidata:List_of_properties/natural_science#Wikidata_property_related_to_biology

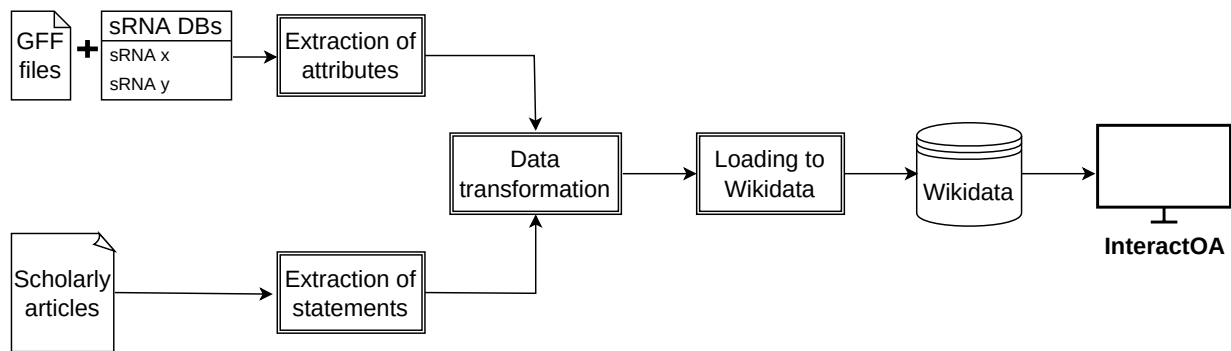


Fig. 1. Schema of the Extract-Transform-Load workflow used for integrating the data into Wikidata and presentation by *InteractOA*.

https://www.wikidata.org/wiki/Q50419231

regulates (molecular biology)

DNA-binding transcriptional dual regulator OxyR b3961 edit

▼ 1 reference

stated in	Small RNA GcvB Regulates Oxidative Stress Response of Escherichia coli
PMCID	8614746
quotation	As a result, it was most likely that the regulation of GcvB on OxyR existed at the post-transcriptional level (English)

[+ add reference](#)

Fig. 2. The sRNA encoded by *gcvB* (Q50419231) as an example of a citation at the statement level.

Next, the Wikidata items representing RNAs were linked to model the interactions. To do this, interaction data for numerous small RNAs was manually obtained from research articles, RegulonDB [37], and the Staphylococcal Regulatory RNA Database" (SRD) [39]. For each of the selected interactions taken from the databases, the underlying statements were manually extracted from the corresponding article and the statements collected in a dedicated file. The compiled information from this manual curation as well as the data automatically extracted from the GFF files was imported to Wikidata using a dedicated Python tool which first extracts the interaction from input files and then maps the names of interaction partners in the parsed file to the respective pre-imported annotations in Wikidata. The item of the pre-imported annotation is retrieved by the tool using a query template to get the corresponding item's ID. After that, the tool links the Wikidata item representing RNAs by using the selected properties. For example, the property "antisense inhibitor of" (P3777) was used to link the interaction between sRNA *omrA* (Item ID Q50419343) and gene *csgF* (Item ID Q23087296). For an example of such a link, see Figure 2. If the type of interaction has no corresponding property in Wikidata, the tool falls back to the more generic property (P128 *regulates*). To showcase the implementation until the point of writing this article, 776 small RNA annotations were imported and 253 connections between small RNAs and their targets were contributed specific to 3 species: *Escherichia coli*,

1 *Vibrio Cholerae*, and *Staphylococcus aureus*. 1

2
3 Both the Python tools used to import the data were built upon Wikidataintegrator⁷ [25] and pywikibot to facilitate 3
4 interaction with the Wikidata API during querying and importing. They are available on GitHub⁸, and they are 4
5 archived in Zenodo⁹. 5
6

7 2.3. Statement-level citations stored in Wikidata 7

8
9 A significant proportion of the sRNA interactions modelled using the method described above appear in open- 9
10 access articles or other articles freely accessible on PubMed Central¹⁰. For several of these interactions, the state- 10
11 ments describing the specific interactions were extracted manually. The Wikidata items of the source article, the 11
12 PubMed Central ID and statements (using the “quotation” property P1683) were added as interaction properties. 12
13 For example, as shown in Figure 2, the property “stated in” (P248) was used with the Wikidata item that corre- 13
14 sponds to a journal article entitled “Small RNA GcvB Regulates Oxidative Stress Response of *Escherichia coli*” 14
15 [40] (Q115652789) in order to link the interaction to the article in which it was mentioned. Moreover, the PubMed 15
16 Central ID of this article was linked with the property “PMCID” (P932), and the statement mentioning the inter- 16
17 action was also linked using the “quotation” property (P1683). This method of storing source statements of sRNA 17
18 interactions in Wikidata makes it easy to check and correct the features of interaction models. 18
19

20 2.4. *InteractOA* as a front end to relevant Wikidata items 21

22
23 The web front end *InteractOA* was developed to facilitate interactive exploration of the data stored in Wikidata, 23
24 as described above. *InteractOA* is implemented in Python and the Flask web framework. Its code is available on 24
25 GitHub¹¹ and its releases are archived at Zenodo¹². Wikidata features a query service¹³ that enables users to en- 25
26 ter SPARQL queries. The service generates tables as well as various types of interactive visualizations, including 26
27 bar charts and in this case, most importantly, network plots. The web front end uses these Wikidata capabilities 27
28 to visualize the regulatory interaction as an interactive network plot (see Figure 3 for an example). The interface 28
29 allows users to customize their queries using filters and to search using keywords without requiring any technical 29
30 knowledge. Once the user has selected the desired filters and keywords, *InteractOA* sends the generated SPARQL 30
31 query to the *Wikidata Query Service* and displays the results. The full landscape of small RNAs is displayed if no 31
32 filters or keywords were used. Figure 3 shows an example network of three small RNAs and their interactions with 32
33 other protein-coding genes based on a limited set of locus tag IDs as keywords. 33
34

35 Additionally, *InteractOA* provides a tabular view for all the extracted statements in Wikidata for a strain selected 35
36 by the user and presents them in a searchable table (see Figure 4 for an example). Its search function can filter 36
37 results by several criteria, for example by partners of interactions or by type of interaction. This solution enables 37
38 users to combine multiple statements from several studies at once and very likely shortens the time needed to 38
39 consult previous research on individual small RNAs. Moreover, the user can open the scholarly article from which 39
40 the statement originated, with the statement itself highlighted (see Figure 5 for an example). *InteractOA* functions 40
41 as an interface between the user and the Wikidata endpoint, with no data storage taking place on the platform itself. 41
42 Essentially, *InteractOA* operates as a middle layer implemented in the Python web framework flask¹⁴, designed to 42
43 facilitate dynamic data querying based on predefined SPARQL query templates. These templates contain parameters 43
44

45 ⁷<https://github.com/SuLab/WikidataIntegrator> 45

46 ⁸https://github.com/foerstner-lab/GFF_to_Wikidata_importer https://github.com/foerstner-lab/sRNA_Interactions_to_Wikidata_Importer 46

47 ⁹<https://zenodo.org/record/7638542> and <https://zenodo.org/record/7638552> 47

48 ¹⁰<https://pubmed.ncbi.nlm.nih.gov/> 48

49 ¹¹<https://github.com/foerstner-lab/InteractOA> 49

50 ¹²<https://doi.org/10.5281/zenodo.7638558> 50

51 ¹³<https://query.wikidata.org/> 51

¹⁴<https://flask.palletsprojects.com/>

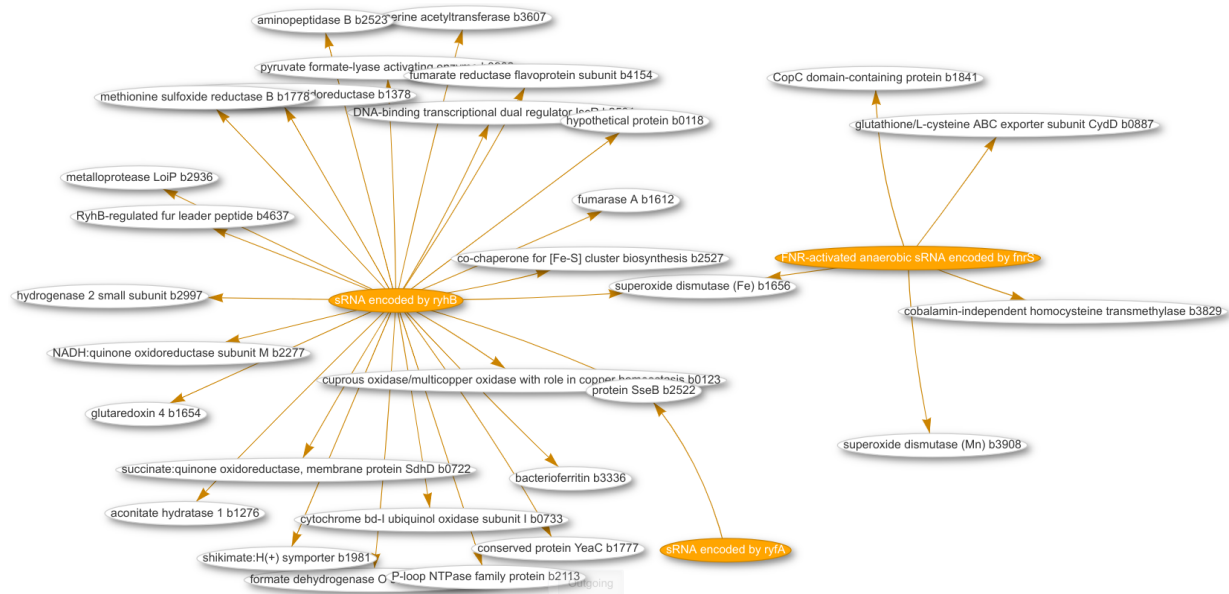


Fig. 3. Screenshot of a network visualization for 3 small RNAs (orange labels) and their interactions with the gene targets (white labels).

that correspond to specific fields and filters on the front end. When a user selects an organism and applies filters, the app processes these inputs and tailors the query from the template accordingly. In the case of visualizing small RNA networks, the specific query is channelled via URL parameters to be executed on Wikidata's query service. An iframe component on the webpage then employs this URL to display the results in the form of a network plot. For the highlighting of statements on small RNA interactions in scholarly articles, the app employs the user's selected organism to pass it to the corresponding parameter in the query template. The Python module WikidataIntegrator¹⁵ [25] is then utilized to execute the query, and the results are included in a tabular view which is then presented to the client. Before transferring the results to the user's web browser, they are enriched with links to the Wikidata items of the sRNAs as well as links to the corresponding article hosted at PubMed Central (PMC) where the selected statements can be highlighted. This colouring works by retrieving the HTML content of the PMCs articles to the server's memory based on the PMC ID and wrapping the content to be highlighted with marking tags, which is then returned to the user's browser and rendered there with the changed background colour. This process allows for aiding users in identifying and focusing on the relevant information.

3. Discussion













Nowadays, much of the research data on which scholarly articles are based is deposited in a structured format in dedicated repositories, yet the actual insights and knowledge derived from this research are often only accessible in an unstructured format within the confines of the article text. In this work, we have presented a solution to this dilemma based on the open-source Wikidata knowledge graph. As shown here, Wikidata provides a structured way to store knowledge generated within specific fields of research, for example by interleaving items of biological entities with bibliographic information while also providing links to the exact statements that are the source of each knowledge item in the corresponding open-access articles. In this work, we have showcased this approach by modelling the regulatory networks of small RNAs in bacteria and built a dedicated web tool that makes it easy to explore and visualize the data stored in Wikidata. The chosen approach also includes granular referencing of knowledge

¹⁵<https://github.com/SuLab/WikidataIntegrator>

Interactions and References

Referenced items: Escherichia coli str. K-12 substr. MG1655

Search:

#	sRNA	sRNA synonyms	Type of Regulation	Target Gene	Quote	Quote from	Wikidata
5	sRNA encoded by gcvB	gcvB, b4443, ECK2804, IS145, JWR0247, psrA11	antisense inhibitor of	serine/threonine:Na(+) symporter b3089	We compared RNA isolated from a wild-type strain and a gcvB deletion strain grown to mid-log phase in Luria-Bertani (LB) broth by microarray analysis to identify any additional regulatory targets of GcvB. One potential target identified by microarray analysis was sstT, which encodes a Na ⁺ /l-serine and l-threonine transport protein.	 	
11	sRNA encoded by gcvB	gcvB, b4443, ECK2804, IS145, JWR0247, psrA11	antisense inhibitor of	branched chain amino acid/phenylalanine ABC transporter periplasmic binding protein b3460	among the top candidate targets for the sRNA GcvB were mRNAs gtlI, livK, livK, yttT, aroP and argT, all genes encoding periplasmic transport proteins.	 	
12	sRNA encoded by gcvB	gcvB, b4443, ECK2804, IS145, JWR0247, psrA11	regulates (molecular biology)	DNA-binding transcriptional dual regulator OxyR b3961	As a result, it was most likely that the regulation of GcvB on OxyR existed at the post-transcriptional level	 	
14	sRNA encoded by gcvB	gcvB, b4443, ECK2804, IS145, JWR0247, psrA11	antisense inhibitor of	oligopeptide ABC transporter periplasmic binding protein b1243	The specific repression of dppA::gfp and oppA::gfp by pPLgcvB was evident from strongly reduced colony fluorescence of these strains on agar plates (Fig. 2B), which established that GcvB regulates dppA and oppA in the 5' mRNA region.	 	

Showing 1 to 12 of 12 entries (filtered from 53 total entries)

Fig. 4. Screenshot of all available statement level citations for a selected organism, which can be filtered with keywords.

We next explored how GcvB stimulated the expression of OxyR. The mRNA level of *oxyR* did not show significant changes in the two transcriptomes of the *gcvB* wild-type and knockout strains ([Supplementary Figure S3A](#)) and this finding was further demonstrated using the RT-qPCR assay ([Supplementary Figure S3B](#)). Moreover, we made an *oxyR* promoter with *lacZ* transcriptional fusion ([Supplementary Figure S3C](#)) in both the *gcvB* wild-type and knockout strains and observed that the β -galactosidase activity showed no significant changes in the two backgrounds ([Supplementary Figure S3D](#)). **As a result, it was most likely that the regulation of GcvB on OxyR existed at the post-transcriptional level.** To substantiate this hypothesis, we constructed the *oxyR* promoter with *lacZ* translational fusions in both the *gcvB* wild-type and knockout strains. We made two fusion constructions, with P1 and P2, respectively, carrying 99 and 45 nt after the translational start codon of *oxyR* ([Figure 4A](#)). Supporting the Western blot result ([Figure 3](#)), both translational fusions showed significantly decreased β -galactosidase activity in the *gcvB* knockout strain when being compared to that in the *gcvB* wild-type strain ([Figure 4B,C](#)), indicating GcvB activated the expression of OxyR at the translational level.

Fig. 5. The highlighted text in an article [40] is a statement level citation example for a claim about small RNA interaction.

sources. Storing the data in Wikidata ensures its long-term availability while opening up access to a large tool chain for imports, queries, and visualization. Moreover, it facilitates links to other relevant entities modelled in Wikidata.

Having demonstrated the usefulness of this approach to our own research field, we now intend to develop the application further to tap into its significant and, as yet, untapped potential. Currently, the steps required to extract statements from research articles are carried out manually. As the quantity of such manually curated article excerpts grows, we aim to train language models to assist with the human curators' extraction work. This text-mining-based approach will build upon related work conducted in our research group (Halder *et al.*, in preparation). For this

1 named-entity recognition (NER) of sRNA taken from databases followed by the extraction of potential interaction 1
2 statements combined with a language model will generate the foundation for manual curation of the data. The cu- 2
3 ration will be conducted by collaborating microbiologist with an expertise of certain bacterial species, especially 3
4 regarding their regulatory networks. For the curation, a web based tool developed by Halder *et al.* will be used. The 4
5 curated data will then be integrated into Wikidata as described above and by that can be explored in *InteractOA*. The 5
6 compiled data will also be used to improve the underlying language model and make the automatic extraction better. 6
7 Besides this, there is an opportunity to directly incorporate other cellular interactions such as regulatory proteins, 7
8 protein-protein interactions, and cellular sensing. 8

9
10 Despite the numerous useful features provided by Wikidata and its ecosystem of tools, there are challenges that 10
11 should be considered when choosing a similar approach. Wikidata's API is comparatively slow, which makes the 11
12 ingestion of larger data sets very time-consuming. Similarly, SPARQL queries are limited by the constraints of the 12
13 *Wikidata Query Service*. This latter issue could be solved by working with full data dumps provided by Wikidata. 13
14 Besides these technical issues, there is inevitably a risk of lower quality, or even significant vandalism, that comes 14
15 with choosing an openly-editable form of data storage in which anybody can add, remove, or modify entries. Thanks 15
16 to versioning, problematic edits do not pose a critical risk, however, and an interface to curate new edits could fur- 16
17 ther address this issue. 17

18
19 One additional concern pertains to the availability of item classes and properties within Wikidata. The imple- 19
20 mented model of *InteractOA* applies Wikidata's currently existing properties to the topic of small RNAs and their 20
21 interaction partners. These item classes and properties despite their sufficiency for the current implementation are 21
22 limited to relatively high-level descriptions such as "regulates" and "antisense inhibitor of". To further improve the 22
23 model, more item classes and properties would need to be agreed upon for Wikidata, for example specifying if the 23
24 regulation is positive or negative for the "regulates" property. There is also a need for further options similar to 24
25 the "antisense inhibitor of" property and for other interaction types such as "promoter of", "cis-acting", or "trans- 25
26 acting". Fortunately, there is currently an ongoing discussion on how to extend the pool of Wikidata properties for 26
27 biological data¹⁶. 27

28
29 Additionally, it has to be considered that currently only a fraction of the scholarly literature is covered by Wiki- 29
30 data¹⁷ [41] (estimated 22.5 Million of 389 Million articles listed in Google Scholar) and that our suggested approach 30
31 might have technical limitation in terms of scalability considering the current technical setup of Wikidata. In case 31
32 this approach would be applied in large scale to numerous fields and millions of articles, the capacity and perfor- 32
33 mance of Wikidata might not be sufficient. As solution for this would be domain specific Wikibase instances. This 33
34 approach could also be helpful to reduce the friction of introducing new item classes and properties which is usually 34
35 time-consuming as a consensus has to be found with all community members participating in the discussion and the 35
36 overall state of the ontology be considered. Still, for the time being, a solution built on Wikidata lowers the access 36
37 barriers to explore this approach by other research communities. 37

38
39 The here presented approach and its implementation stands at the intersection of Wikidata and biological scholars' 39
40 communities to bridge the gap between them, and to bring the advantages of Wikidata to the forefront, promoting 40
41 its use among biological scholars. Based on the showcase given by this, the prospective users would be encouraged 41
42 to contribute, curate, and save their findings as a result of highlighting Wikidata as a readily achievable centralized 42
43 solution for the preservation of their discoveries, as well a platform for promoting and citing their work. The in- 43
44 troduced statement-level citation in *InteractOA* as a feature that could potentially simplify the process of locating 44
45 and comparing information due to the amassed collective knowledge. Despite the challenges, we are convinced 45
46 that the approach showcased in this project can be applied to numerous other communities, even those that require 46
47 more complex data models. We hope the example provided here will motivate other research communities to make 47
48 knowledge and its sources available in a more structured fashion. 48

49
50 ¹⁶https://www.wikidata.org/wiki/Wikidata:WikiProject_Molecular_biology/Properties 50

51 ¹⁷<https://www.wikidata.org/wiki/Wikidata:Statistics> 51

4. Acknowledgments

This work was supported by the Bundesministerium für Bildung und Forschung (BMBF, grant number 16OA031Z).

References

- [1] P. Hitzler, A review of the semantic web field, *Communications of the ACM* **64**(2) (2021), 76–83. doi:10.1145/3397512.
- [2] A. Hogan, E. Blomqvist, M. Cochez, C. D'amato, G.D. Melo, C. Gutierrez, S. Kirrane, J.E.L. Gayo, R. Navigli, S. Neumaier, A.-C.N. Ngomo, A. Polleres, S.M. Rashid, A. Rula, L. Schmelzeisen, J. Sequeda, S. Staab and A. Zimmermann, Knowledge Graphs, *ACM Computing Surveys* **54**(4) (2021), 1–37. doi:10.1145/3447772.
- [3] M. Kejriwal, Knowledge Graphs: A Practical Review of the Research Landscape, *Information* **13**(4) (2022), 161. doi:10.3390/info13040161.
- [4] D. Brickley, R.V. Guha and B. McBride, RDF schema 1.1. W3C recommendation, *World Wide Web Consortium* **2** (2014).
- [5] D. Vrandečić, Wikidata: a new platform for collaborative data collection, in: *Proceedings of the 21st International Conference on World Wide Web*, ACM, 2012. doi:10.1145/2187980.2188242.
- [6] D. Vrandečić and M. Krötzsch, Wikidata: a free collaborative knowledgebase, *Communications of the ACM* **57**(10) (2014), 78–85. doi:10.1145/2629489.
- [7] F. Erxleben, M. Günther, M. Krötzsch, J. Mendez and D. Vrandečić, Introducing Wikidata to the Linked Data Web, in: *The Semantic Web – ISWC 2014*, Springer International Publishing, 2014, pp. 50–65. doi:10.1007/978-3-319-11964-9_4.
- [8] K. Shenoy, F. Ilievski, D. Garijo, D. Schwabe and P. Szekely, A study of the quality of Wikidata, *Journal of Web Semantics* **72** (2022), 100679. doi:10.1016/j.websem.2021.100679.
- [9] T.V. Veen, Wikidata - From “an” Identifier to “the” Identifier, *Information Technology and Libraries* **38**(2) (2019), 72–81. doi:10.6017/ital.v38i2.10886.
- [10] M. Mora-Cantalops, S. Sánchez-Alonso and E. García-Barriocanal, A systematic literature review on Wikidata, *Data Technologies and Applications* **53**(3) (2019), 250–268. doi:10.1108/dta-12-2018-0110.
- [11] M.D. Wilkinson, M. Dumontier, I.J. Aalbersberg, G. Appleton, M. Axton, A. Baak, N. Blomberg, J.-W. Boiten, L.B. da Silva Santos, P.E. Bourne, J. Bouwman, A.J. Brookes, T. Clark, M. Crosas, I. Dillo, O. Dumon, S. Edmunds, C.T. Evelo, R. Finkers, A. Gonzalez-Beltran, A.J.G. Gray, P. Groth, C. Goble, J.S. Grethe, J. Heringa, P.A.C. 't Hoen, R. Hooft, T. Kuhn, R. Kok, J. Kok, S.J. Lusher, M.E. Martone, A. Mons, A.L. Packer, B. Persson, P. Rocca-Serra, M. Roos, R. van Schaik, S.-A. Sansone, E. Schultes, T. Sengstag, T. Slater, G. Strawn, M.A. Swertz, M. Thompson, J. van der Lei, E. van Mulligen, J. Velterop, A. Waagmeester, P. Wittenburg, K. Wolstencroft, J. Zhao and B. Mons, The FAIR Guiding Principles for scientific data management and stewardship, *Scientific Data* **3**(1) (2016). doi:10.1038/sdata.2016.18.
- [12] D. Shotton, K. Portwin, G. Klyne and A. Miles, Adventures in Semantic Publishing: Exemplar Semantic Enhancements of a Research Article, *PLoS Computational Biology* **5**(4) (2009), e1000361. doi:10.1371/journal.pcbi.1000361.
- [13] A. Garcia, F. Lopez, L. Garcia, O. Giraldo, V. Bucheli and M. Dumontier, Biotea: semantics for Pubmed Central, *PeerJ* **6** (2018), e4201. doi:10.7717/peerj.4201.
- [14] S. Auer, A. Oelen, M. Haris, M. Stocker, J. D'Souza, K.E. Farfar, L. Vogt, M. Prinz, V. Wiens and M.Y. Jaradeh, Improving Access to Scientific Literature with Knowledge Graphs, *Bibliothek Forschung und Praxis* **44**(3) (2020), 516–529. doi:10.1515/bfp-2020-2042.
- [15] T.U. Consortium, UniProt: the universal protein knowledgebase, *Nucleic Acids Research* **46**(5) (2018), 2699–2699. doi:10.1093/nar/gky092.
- [16] S. Kerrien, Y. Alam-Faruque, B. Aranda, I. Bancarz, A. Bridge, C. Derow, E. Dimmer, M. Feuermann, A. Friedrichsen, R. Huntley, C. Kohler, J. Khadake, C. Leroy, A. Liban, C. Liefertink, L. Montecchi-Palazzi, S. Orchard, J. Risse, K. Robbe, B. Roechert, D. Thorneycroft, Y. Zhang, R. Apweiler and H. Hermjakob, IntAct—open source resource for molecular interaction data, *Nucleic Acids Research* **35**(Database) (2007), D561–D565. doi:10.1093/nar/gkl958.
- [17] L. Li, D. Huang, M.K. Cheung, W. Nong, Q. Huang and H.S. Kwan, BSRD: a repository for bacterial small regulatory RNA, *Nucleic Acids Research* **41**(D1) (2012), D233–D238. doi:10.1093/nar/gks1264.
- [18] J. Pischmarov, C. Kuenne, A. Billion, J. Hemberger, F. Cemič, T. Chakraborty and T. Hain, sRNADB: A small non-coding RNA database for gram-positive bacteria, *BMC Genomics* **13**(1) (2012). doi:10.1186/1471-2164-13-384.
- [19] M.M. Hoffman, AANT: the Amino Acid-Nucleotide Interaction Database, *Nucleic Acids Research* **32**(90001) (2004), 174D–181. doi:10.1093/nar/gkh128.
- [20] J.E. Hill, S.L. Penny, K.G. Crowell, S.H. Goh and S.M. Hemmingsen, cpnDB: A Chaperonin Sequence Database, *Genome Research* **14**(8) (2004), 1669–1675. doi:10.1101/gr.2649204.
- [21] Á. Ósz, L.S. Pongor, D. Szirmai and B. Györfy, A snapshot of 3649 Web-based services published between 1994 and 2017 shows a decrease in availability after 2 years, *Briefings in Bioinformatics* **20**(3) (2017), 1004–1010. doi:10.1093/bib/bbx159.
- [22] J.D. Wren, C. Georgescu, C.B. Giles and J. Hennessey, Use it or lose it: citations predict the continued online availability of published bioinformatics resources, *Nucleic Acids Research* **45**(7) (2017), 3627–3633. doi:10.1093/nar/gkx182.

- [23] S. Mangul, T. Mosqueiro, R.J. Abdill, D. Duong, K. Mitchell, V. Sarwal, B. Hill, J. Brito, R.J. Littman, B. Statz, A.K.-M. Lam, G. Dayama, L. Grieneisen, L.S. Martin, J. Flint, E. Eskin and R. Blekhman, Challenges and recommendations to improve the installability and archival stability of omics computational tools, *PLOS Biology* **17**(6) (2019), e3000333. doi:10.1371/journal.pbio.3000333.
- [24] T.K. Attwood, B. Agit and L.B.M. Ellis, Longevity of Biological Databases, *EMBNET Journal* **21**(0) (2015). doi:10.14806/ej.21.0.803.
- [25] A. Waagmeester, G. Stupp, S. Burgstaller-Muehlbacher, B.M. Good, M. Griffith, O.L. Griffith, K. Hanspers, H. Hermjakob, T.S. Hudson, K. Hybiske, S.M. Keating, M. Manske, M. Mayers, D. Mietchen, E. Mitraga, A.R. Pico, T. Putman, A. Riutta, N. Queralt-Rosinach, L.M. Schriml, T. Shafee, D. Slenter, R. Stephan, K. Thornton, G. Tsueng, R. Tu, S. Ul-Hasan, E. Willighagen, C. Wu and A.I. Su, Wikidata as a knowledge graph for the life sciences, *eLife* **9** (2020). doi:10.7554/elife.52614.
- [26] A. Waagmeester, E.L. Willighagen, A.I. Su, M. Kutmon, J.E.L. Gayo, D. Fernández-Álvarez, Q. Groom, P.J. Schaap, L.M. Verhagen and J.J. Koehorst, A protocol for adding knowledge to Wikidata: aligning resources on human coronaviruses, *BMC Biology* **19**(1) (2021). doi:10.1186/s12915-020-00940-y.
- [27] S. Bonner, I.P. Barrett, C. Ye, R. Swiers, O. Engkvist, A. Bender, C.T. Hoyt and W.L. Hamilton, A review of biomedical datasets relating to drug discovery: a knowledge graph perspective, *Briefings in Bioinformatics* **23**(6) (2022), bbac404. doi:10.1093/bib/bbac404.
- [28] J. Youn, N. Rai and I. Tagkopoulos, Knowledge integration and decision support for accelerated discovery of antibiotic resistance genes, *Nature Communications* **13**(1) (2022). doi:10.1038/s41467-022-29993-z.
- [29] A. Santos, A.R. Colaço, A.B. Nielsen, L. Niu, M. Strauss, P.E. Geyer, F. Coscia, N.J.W. Albrechtsen, F. Mundt, L.J. Jensen and M. Mann, A knowledge graph to interpret clinical proteomics data, *Nature Biotechnology* **40**(5) (2022), 692–702. doi:10.1038/s41587-021-01145-6.
- [30] T.E. Putman, S. Lelong, S. Burgstaller-Muehlbacher, A. Waagmeester, C. Diesh, N. Dunn, M. Munoz-Torres, G.S. Stupp, C. Wu, A.I. Su and B.M. Good, WikiGenomes: an open web application for community consumption and curation of gene annotation data in Wikidata, *Database* **2017** (2017). doi:10.1093/database/bax025.
- [31] T. Putman, K. Hybiske, D. Jow, C. Afrasiabi, S. Lelong, M.A. Cano, G.S. Stupp, A. Waagmeester, B.M. Good, C. Wu and A.I. Su, ChlamBase: a curated model organism database for the Chlamydia research community, *Database* **2019** (2019). doi:10.1093/database/baz041.
- [32] M. Martens, A. Ammar, A. Riutta, A. Waagmeester, D.N. Slenter, K. Hanspers, R.A. Miller, D. Digles, E.N. Lopes, F. Ehrhart, L.J. Dupuis, L.A. Winklers, S.L. Coort, E.L. Willighagen, C.T. Evelo, A.R. Pico and M. Kutmon, WikiPathways: connecting communities, *Nucleic Acids Research* **49**(D1) (2020), D613–D621. doi:10.1093/nar/gkaa1024.
- [33] S. Burgstaller-Muehlbacher, A. Waagmeester, E. Mitraga, J. Turner, T. Putman, J. Leong, C. Naik, P. Pavlidis, L. Schriml, B.M. Good and A.I. Su, Wikidata as a semantic framework for the Gene Wiki initiative, *Database* **2016** (2016), baw015. doi:10.1093/database/baw015.
- [34] E.G.H. Wagner and P. Romby, Small RNAs in Bacteria and Archaea, in: *Advances in Genetics*, Elsevier, 2015, pp. 133–208. doi:10.1016/bs.adgen.2015.05.001.
- [35] G. Storz, S. Altuvia and K.M. Wassarman, AN ABUNDANCE OF RNA REGULATORS, *Annual Review of Biochemistry* **74**(1) (2005), 199–217. doi:10.1146/annurev.biochem.74.082803.133136.
- [36] M.G. Jørgensen, J.S. Pettersen and B.H. Kallipolitis, sRNA-mediated control in bacteria: An increasing diversity of regulatory mechanisms, *Biochimica et Biophysica Acta (BBA) - Gene Regulatory Mechanisms* **1863**(5) (2020), 194504. doi:10.1016/j.bbagr.2020.194504.
- [37] A. Huerta, RegulonDB: a database on transcriptional regulation in Escherichia coli, *Nucleic Acids Research* **26**(1) (1998), 55–59. doi:10.1093/nar/26.1.55.
- [38] N.A. O'Leary, M.W. Wright, J.R. Brister, S. Ciufu, D. Haddad, R. McVeigh, B. Rajput, B. Robbertse, B. Smith-White, D. Ako-Adjei, A. Astashyn, A. Badretdin, Y. Bao, O. Blinkova, V. Brover, V. Chetvernin, J. Choi, E. Cox, O. Ermolaeva, C.M. Farrell, T. Goldfarb, T. Gupta, D. Haft, E. Hatcher, W. Hlavina, V.S. Joardar, V.K. Kodali, W. Li, D. Maglott, P. Masterson, K.M. McGarvey, M.R. Murphy, K. O'Neill, S. Pujar, S.H. Rangwala, D. Rausch, L.D. Riddick, C. Schoch, A. Shkeda, S.S. Storz, H. Sun, F. Thibaud-Nissen, I. Tolstoy, R.E. Tully, A.R. Vatsan, C. Wallin, D. Webb, W. Wu, M.J. Landrum, A. Kimchi, T. Tatusova, M. DiCuccio, P. Kitts, T.D. Murphy and K.D. Pruitt, Reference sequence (RefSeq) database at NCBI: current status, taxonomic expansion, and functional annotation, *Nucleic Acids Research* **44**(D1) (2015), D733–D745. doi:10.1093/nar/gkv1189.
- [39] M. Sassi, Y. Augagneur, T. Mauro, L. Ivain, S. Chabelskaya, M. Hallier, O. Sallou and B. Felden, SRD: a Staphylococcus regulatory RNA database, *RNA* **21**(5) (2015), 1005–1017. doi:10.1261/rna.049346.114.
- [40] X. Ju, X. Fang, Y. Xiao, B. Li, R. Shi, C. Wei and C. You, Small RNA GcvB Regulates Oxidative Stress Response of Escherichia coli, *Antioxidants* **10**(11) (2021), 1774. doi:10.3390/antiox10111774.
- [41] M. Gusenbauer, Google Scholar to overshadow them all? Comparing the sizes of 12 academic search engines and bibliographic databases, *Scientometrics* **118**(1) (2018), 177–214. doi:10.1007/s11192-018-2958-5.