# Introduction to Gene Wiki: Using wikidata as a proxy to scholarly articles

Andra Waagmeester[1]

1)    Micelio, Antwerp, Belgium | Email: andra@micelio.be, Twitter: @andrawaag

# The Gene Wiki project, circa 2008

**Summarized knowledge via crowdsourcing**

**Data imported from structured databases**

Huss, PLoS Biol, 2008

# Reelin

From Wikipedia, the free encyclopedia

**Reelin** is a large secreted extracellular matrix glycoprotein that helps regulate processes of neuronal migration and positioning in the developing brain by controlling cell–cell interactions. Besides this important role in early development, reelin continues to work in the adult brain. It modulates synaptic plasticity by [2][3] It also stimulates dendrite[4] migration of neuroblasts genera zones. It is found not only in the tissues.

Reelin has been suggested to b expression of the protein has be bipolar disorder, but the cause of this observation remains uncertain as studies show that psychotropic medication itself affects reelin expression. Moreover, epigenetic hypotheses aimed at explaining the changed levels of reelin expression[6] are controversial.[7][8] Total lack of reelin causes a form of lissencephaly. Reelin may also play a role in Alzheimer's disease, temporal lobe epilepsy and autism.

Reelin's name comes from the abnormal reeling gait of *reeler* mice,[9] which were later found to have a deficiency of this brain protein and were homozygous for mutation of the RELN gene. The

---

Overlay text box:

Reelin has been suggested to be implicated in pathogenesis of several brain diseases. **The expression of the protein has been found to be significantly lower in schizophrenia and psychotic bipolar disorder,** but the cause of this observation remains uncertain as studies show that psychotropic medication itself

---

**Reelin**

...hic structure of the third reelin repeat domain.[1]

**Available structures**

| PDB | Ortholog search: PDBe, RCSB | |
|---|---|---|
| | List of PDB id codes | [show] |

**Identifiers**

**Symbols** RELN ; LIS2; PRO1598; RL

**External** OMIM: 600514 MGI: 103022

# Wikipedia: Maintained independently by >300 language communities



Dutch

Greek

English



Dutch

Greek

English

# Wikidata is to data as Wikipedia is to text

Wikidata is a collaboratively edited knowledge base operated by the Wikimedia Foundation

- Completely free, even for commercial usage (CC0)
- Anybody can contribute
- Covers all domains of knowledge
- Extensive item history, talk pages, projects, users
- Integration with the semantic web
- High performance query engine (SPARQL)

- Stable! Long term support not dictated by funding cycles
- Actively developed
- Already has large number of active users, editors, contributors!

**WIKIDATA**

**A giant graph of knowledge!**

# Getting data in..

# Community engagement and model discussion

# Formally capture and describe model and community consensus

**Model development**

- Legacy review – develop punch lists for existing data issues that needs fixing

- Documentation – terse, human-readable representation helping contributers and maintainers quickly grok the model

- Client pre-submission – submitters test their data before submission to make sure they're saying what they want to say and that the receiving schema can accommodate all of their data

- Server pre-ingestion – submission process checks data as it comes in and either rejects or warns about non-conformant data



```
Data (Turtle)
<samples9298996>
  rdf:type bf:Text ;
  rdf:type bf:Work ;
  bf:title "Oliver Twist." ;
  bf:class <id.loc.gov/…/PZ3> ;
  bf:creator [
    rdf:type bf:Person ;
    bf:label "Dickens, Charles, 1812-1870." ;
  ] .

<id.loc.gov/…/PZ3>
  rdf:type bf:LCC ;
  bf:label "PZ3.D55O165PR4567" .
```

| pt | gene humano | | ✏edit |

```
# E108: genome_assembly
IMPORT <https://www.wikidata.org/wiki/Special:EntitySchemaText/E108>
PREFIX E108: <https://www.wikidata.org/wiki/Special:EntitySchemaText/E108#>

# E109: human chromosome
IMPORT <https://www.wikidata.org/wiki/Special:EntitySchemaText/E109>
```

```
p:P31 @<#P31_instance_of_gene> ;
```

```
PREFIX wdt: <http://www.wikidata.org/prop/direct/>
PREFIX p: <http://www.wikidata.org/prop/>
```

```
<#P31_instance_of_gene> {
    ps:P31 @<#gene_types> ;      # Instance of [P31] gene types
    prov:wasDerivedFrom @<#ncbi-gene-reference> OR @<#ensembl-gene-reference>
}
```

```
start = @<#wikidata-human-gene>
```

```
(
        p:P644 @<#P644_genomic_start> ; # Its genomic start location
        p:P645 @<#P645_genomic_end> ; # Its genomic end location
)* ; # Zero or more start and end locations.
```

```
# Value statements contain either actual values, or pointers to other Wikidata items.
Identifier statements capture
# external identifiers, erroneous statements are those that are errors.
```

check entities against this Schema⧉ | ✏edit |

Enter an entity to check e.g.Q42    [Check]

# 🐝ShEx2 — Simple Online Validator

```
# Shape Expression for Human genes in Wikidata
PREFIX wd: <http://www.wikidata.org/entity/>
PREFIX wdt: <http://www.wikidata.org/prop/direct/>
PREFIX p: <http://www.wikidata.org/prop/>
PREFIX prov: <http://www.w3.org/ns/prov#>
PREFIX pq: <http://www.wikidata.org/prop/qualifier/>
PREFIX xsd: <http://www.w3.org/2001/XMLSchema#>
PREFIX prv: <http://www.wikidata.org/prop/reference/value/>
PREFIX pr:  <http://www.wikidata.org/prop/reference/>
PREFIX ps: <http://www.wikidata.org/prop/statement/>

BASE <http://www.wikidata.org/entity/>

start = @<#wikidata-human-gene>

# Query with results
# SELECT * WHERE {?item wdt:P31 wd:Q7187 ; wdt:P703 wd:Q15978631 .} LIMIT 10

# Indicates which shape to use to start iterating over the graph if none is
provided.

# wikidata-human gene is the main shape for a human gene data model in Wikidata.
Each line between the brackets
# represents the structure than can be enforced to validate human gene annotations
in Wikidata
```

abort (ctl-enter)

| Query | Entities to check |

| | | | |
|---|---|---|---|
| <http://www.wikidata.org/entity/Q414043> | @ | START | - ✓ |
| <http://www.wikidata.org/entity/Q415594> | @ | START | - ✓ |
| <http://www.wikidata.org/entity/Q416426> | @ | START | - ✓ |
| <http://www.wikidata.org/entity/Q417169> | @ | START | - ✓ |
| <http://www.wikidata.org/entity/Q417743> | @ | START | - ✓ |
| <http://www.wikidata.org/entity/Q418553> | @ | START | - ✓ |

# Seeding with data

- Model structure of items (genes, drugs, diseases, .. etc) & relationships between items
- Import data from many sources and ontologies
- Linked to many identifiers from external databases
- Architecture for maintaining data from external sources

<> Code      ⓘ Issues  4      ⑃ Pull requests  1      🗐 Projects  0      ⚡ Pulse      ⊪ Graphs

A Wikidata Python module integrating the MediaWiki API and the Wikidata SPARQL endpoint

🕐 **397** commits      ⑂ **2** branches      🏷 **1** release      👥 **7** contributors      ⚖ MIT

Branch: **master** ▾      New pull request            Find file      **Clone or download** ▾

🔲🔲 **sebotic** fixed an omission where new items don't get created when domain not s...   …      Latest commit `2f5d2fd` 22 hours ago

📁 doc                          Wikidata to Wikipedia mapping prototype for diseases added.                          2 years ago

📁 wikidataintegrator           fixed an omission where new items don't get created when domain not s…                 22 hours ago

🔲 **Jenkins**

Jenkins  ▸  Running  ▸

New Item

People

Build History

Edit View

Delete View

Manage Jenkins

My Views

Credentials

**Build Queue**

Running Bots

| All | **Running** | + |

| S | Name | Last Success ↑ | Last Failure |
|---|---|---|---|
| 🔵 | ProteinBot_homo_sapiens | 1 day 21 hr - #12 | N/A |
| 🔵 | GOBot_bigmem | 2 days 15 hr - #15 | 9 days 15 hr - #14 |
| 🔵 | GeneBot_Homo_sapiens | 2 days 19 hr - #25 | 2 days 20 hr - #24 |
| 🔵 | Disease_Ontology | 2 days 23 hr - #11 | 4 days 13 hr - #8 |
| 🔵 | GeneDiseaseBot | 2 days 23 hr - #9 | 1 mo 6 days - #2 |

Feedback loop

Examples  Help  More tools  文A English

```
1  SELECT * WHERE {
2      ?item wdt:P356 ?Doi}
```

Server error: Unexpected end of JSON input

# Wikibase and WBStack



github.com/wmde/wikibase-docker

wbstack.com



wmde / **wikibase-docker**

Code | Pull requests 6 | Actions | Security | Insights

🐳 Docker images and example compose file for Wikibase and surrounding services



WbStack Alpha

Hi andra@micel.io    DASHBOARD    ACCOUNT    LOGOUT

## Your open data..

..starting with a Wikibase stack

**Blog posts:**
- Infrastructure overview
- November review
- October introduction

**What is Wikibase?**

Wikibase is an open source software suite for running a collaborative knowledge base. One installation of it is Wikidata.

— learningwikibase.com CC-BY 4.0

**Community**          **Organization**
Mediawiki              About Us
Wikibase               Privacy

# Introducing Wikiproject Biodiversity



https://www.wikidata.org/wiki/Wikidata:WikiProject_Biodiversity

```
1  #species observed through iNaturalist not being described by an English Wikipedia article
2  PREFIX wbt: <http://biodiversity.wiki.opencura.com/prop/direct/>
3  PREFIX wb: <http://biodiversity.wiki.opencura.com/entity/>
4  SELECT DISTINCT ?taxon ?taxonLabel ?wikidata WHERE {
5      ?taxon wbt:P13 wb:Q131918 ;
6             wbt:P11 ?wikidata .
7    SERVICE <https://query.wikidata.org/sparql> {
8        ?wikidata wdt:P31 [] .
9        MINUS {?article schema:about ?wikidata ;schema:isPartOf <https://en.wikipedia.org/> ;schema:inLanguage "en" .}
10      }
11   SERVICE wikibase:label { bd:serviceParam wikibase:language "[AUTO_LANGUAGE],en". }
12  }
```

👁▾  ❓                                                                    3753 results in 49728 ms

| taxon | taxonLabel |
| --- | --- |
| 🔍 <http://biodiversity.wiki.opencura.com/entity/Q101721> | Pyrochroa serraticornis |
| 🔍 <http://biodiversity.wiki.opencura.com/entity/Q101862> | Plagionotus arcuatus |
| 🔍 <http://biodiversity.wiki.opencura.com/entity/Q102061> | Panorpa germanica |
| 🔍 <http://biodiversity.wiki.opencura.com/entity/Q103399> | Cryptocephalus moraei |
| 🔍 <http://biodiversity.wiki.opencura.com/entity/Q103704> | Ctenicera pyrrhos |

https://tinyurl.com/y5mqqkgl

# Incomplete... but doesn't need to be complete

Example: Get drugs that act as channel blockers from Wikidata, get the pathways that these drugs are part of from Wikipathways

```
Wikidata Query    Examples    Prefixes    Tools    Help
1 PREFIX bd: <http://www.bigdata.com/rdf#>
2 PREFIX wp:      <http://vocabularies.wikipathways.org/wp#>
3 PREFIX dcterms:  <http://purl.org/dc/terms/>
4 PREFIX dc:       <http://purl.org/dc/elements/1.1/>
5
6 SELECT DISTINCT ?metabolite ?wikidatadrug ?wikidatadrugLabel
7                 ?title ?wpIdentifier WHERE {
8   ?protein wdt:P279* wd:Q422500 .
9   ?protein wdt:P279|wdt:P31 wd:Q8054 .
10  ?wikidatadrug wdt:P129 ?protein .
11  ?wikidatadrug p:P129/pq:P794 wd:Q389934 .
12  SERVICE wikibase:label { bd:serviceParam wikibase:language "en" }
13  SERVICE <http://sparql.wikipathways.org/> {
14    ?metabolite a wp:Metabolite ;
15      wp:bdbWikidata ?wikidatadrug ;
16      dcterms:isPartOf ?pathway .
17    ?pathway a wp:Pathway .
18      ?pathway dc:title ?title .
19      ?pathway dc:identifier ?wpIdentifier .
20   }
21 } LIMIT 100
```

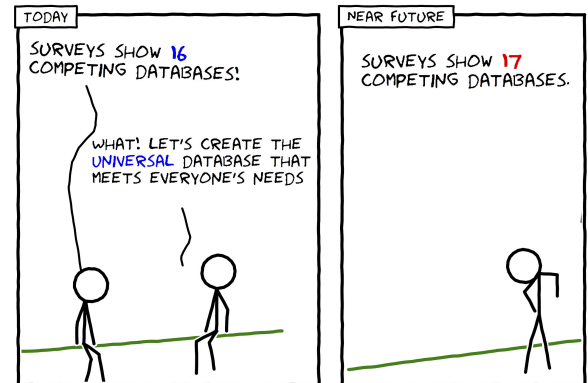http://tinyurl.com/m8g7q7p

Not trying to create the

Can be linked to greater semantic web

**Federated queries**: The ability to make one query across multiple sources

*Interoperability!*

| metabolite | wikidatadrug | wikidatadrugLabel | title | wpIdentifier |
|---|---|---|---|---|
| <http://identifiers.org/hmdb/HMDB01852> | wd:Q29417 | tretinoin | Nuclear receptors in lipid metabolism and toxicity | <http://identifiers.org/wikipathways/WP1099> |
| <http://identifiers.org/hmdb/HMDB01852> | wd:Q29417 | tretinoin | Nuclear receptors in lipid metabolism and toxicity | <http://identifiers.org/wikipathways/WP1326> |
| <http://identifiers.org/hmdb/HMDB01852> | wd:Q29417 | tretinoin | Cardiac Progenitor Differentiation | <http://identifiers.org/wikipathways/WP3127> |
| <http://identifiers.org/hmdb/HMDB01852> | wd:Q29417 | tretinoin | Dopaminergic Neurogenesis | <http://identifiers.org/wikipathways/WP3196> |
| <http://identifiers.org/hmdb/HMDB01852> | wd:Q29417 | tretinoin | Nuclear Receptors in Lipid Metabolism and Toxicity | <http://identifiers.org/wikipathways/WP299> |



HOW BIOLOGICAL DATABASES PROLIFERATE

TODAY
SURVEYS SHOW 16 COMPETING DATABASES!
WHAT! LET'S CREATE THE UNIVERSAL DATABASE THAT MEETS EVERYONE'S NEEDS

NEAR FUTURE
SURVEYS SHOW 17 COMPETING DATABASES.

https://doi.org/10.1371/journal.pcbi.1005128.g001

# Acknowledgements

**Wikidata as a FAIR knowledge graph for the life sciences**

iD Andra Waagmeester, iD Gregory Stupp, iD Sebastian Burgstaller-Muehlbacher, iD Benjamin M Good, iD Malachi Griffith, iD Obi Griffith, iD Kristina Hanspers, iD Henning Hermjakob, iD Kevin Hybiske, iD Sarah M. Keating, iD Magnus Manske, iD Michael Mayers, iD Elvira Mitraka, iD Alexander R. Pico, iD Timothy Putman, iD Anders Riutta, iD Núria Queralt-Rosinach, iD Lynn M. Schriml, iD Denise Slenter, iD Ginger Tsueng, iD Roger Tu, iD Egon Willighagen, iD Chunlei Wu, iD Andrew I Su

**doi:** https://doi.org/10.1101/799684

**Thousands of Wikidatans**



By Helpameout - Own work, CC BY-SA 3.0, https://commons.wikimedia.org/w/index.php?curid=20337311

# References

- [A protocol for adding knowledge to Wikidata, a case report (Q90557988)](#)
- [Wikidata as a knowledge graph for the life sciences (Q87830400)](#)
- [Using Shape Expressions (ShEx) to Share RDF Data Models and to Guide Curation with Rigorous Validation (Q64225211)](#)
- [Wikidata as an intuitive resource towards semantic data modeling in data FAIRification (Q59925812)](#)
- [Wikidata as a semantic framework for the Gene Wiki initiative (Q23712646)](#)

# Hands-on

- Search for a Wikipedia article that cites a DOI; (e.g. https://en.wikipedia.org/wiki/Cetiosauriscus )
- Check if that DOI already exists in Wikidata (e.g. https://w.wiki/j9q )
- in Wikidata on a scholarly article using a DOI; https://boiling-ridge-97667.herokuapp.com/
- Annotate the wikidata item with metadata e.g main subject (e.g. https://www.wikidata.org/wiki/Q100968025 )
- Disambiguate its authors (i.e. resolve its author string to author): https://sourcemd.toolforge.org/new_resolve_authors.php
- Explore that article using Scholia: https://scholia.toolforge.org/work/Q100968025