# Lexicographical data

## An Introduction to Lexemes on Wikidata

**Lydia Pintscher**
Wikidata Portfolio Lead,
Wikimedia Deutschland

**Mohammed Abdulai**
Wikidata Community
Communication Manager,
Wikimedia Deutschland

This session is recorded: Please mute your microphone and camera when you're not speaking.

**Wikidata Lexicodays** 28-30 June 2024

# Agenda

- Wikidata: A Quick Recap

- Words In Wikidata? How Does It Work?

- Why Is Lexicographical Data On Wikidata Important?

- What Could We Do With Lexemes?

- The State Of Lexicographical Data On Wikidata

- Useful Resources

# Wikidata: A quick recap

# Wikidata: the basics

- A knowledge base
- Part of the Wikimedia projects
- Structured data
- Linked to other databases
- Multilingual
- Collaborative
- Released under public domain (CC0)
- Based on facts and references
- Made for humans and machines

# Wikidata Item

# ¿Words in Wikidata: How does it work?

Wikidata
Lexicodays
28-30 June 2024

# Wait, what's the difference?

## Concept "mouse"

- Species of mammal
- Taxon name
- Average size
- Picture
- Encyclopedia of Life ID



File:House mouse.jpg, public domain

## Lexeme "mouse"

- Language: English
- Lexical category: noun
- Plural form: mice (irregular)
- Etymology: Proto-Germanic *mūs
- Senses: animal, computer device
- Translations: souris (fr), rato (pt)
- Audio pronunciation ▶

/maʊs/

**Wikidata Lexicodays** 28-30 June 2024

Glossary:
[Wikidata:Lexicographical_data/Glossary](Wikidata:Lexicographical_data/Glossary)

Data model:
[MediaWiki:Extension:WikibaseLexeme/Data Model](MediaWiki:Extension:WikibaseLexeme/Data Model)

L-id **Lexeme**
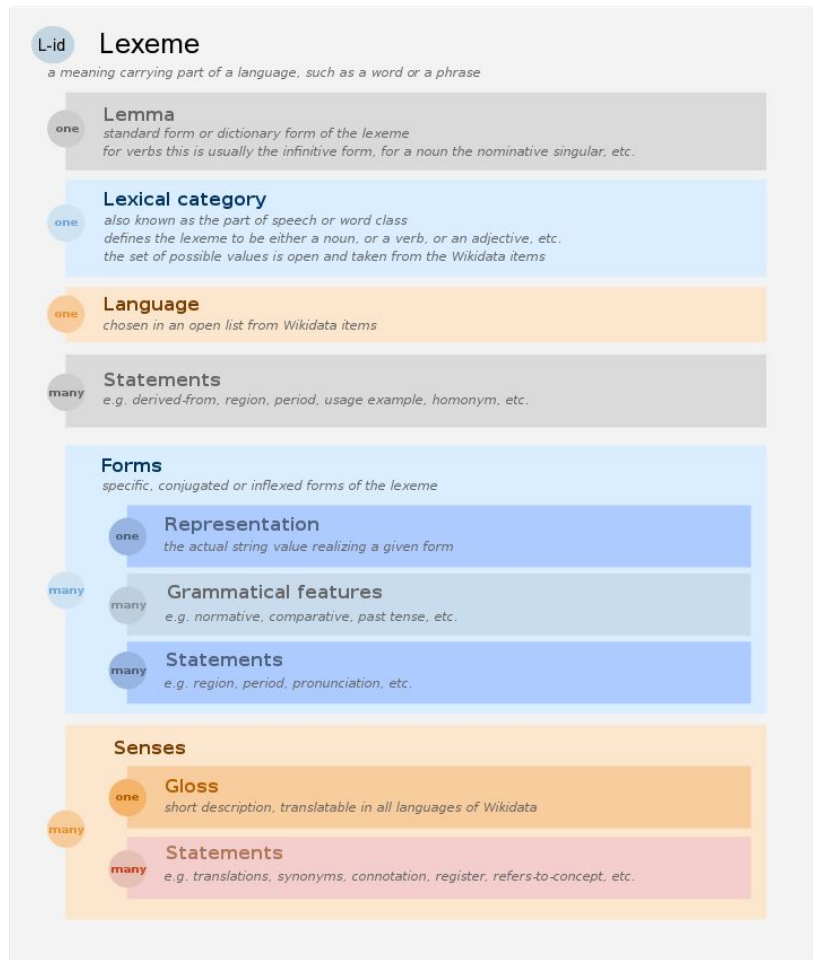*a meaning carrying part of a language, such as a word or a phrase*

one **Lemma**
standard form or dictionary form of the lexeme
for verbs this is usually the infinitive form, for a noun the nominative singular, etc.

one **Lexical category**
also known as the part of speech or word class
defines the lexeme to be either a noun, or a verb, or an adjective, etc.
the set of possible values is open and taken from the Wikidata items

one **Language**
chosen in an open list from Wikidata items

many **Statements**
e.g. derived-from, region, period, usage example, homonym, etc.

**Forms**
specific, conjugated or inflexed forms of the lexeme

many one **Representation**
*the actual string value realizing a given form*

many **Grammatical features**
e.g. normative, comparative, past tense, etc.

many **Statements**
e.g. region, period, pronunciation, etc.

**Senses**

many one **Gloss**
short description, translatable in all languages of Wikidata

many **Statements**
e.g. translations, synonyms, connotation, register, refers-to-concept, etc.

**Wikidata Lexicodays**
⬚⬚⬚⬚ **28-30 June 2024**

**WIKIDATA**

Lexeme    Discussion                                                                      Read    View history    ☆

(L1119)   # mouse
          ✏️edit

          en

Language English
Lexical category noun

## Statements

| instance of | ⬍ | count noun | ✏️edit |
|---|---|---|---|
| | | ▾ 0 references | |
| | | | ➕ add reference |
| | | | ➕ add value |

| described by source | ⬍ | Merriam-Webster online dictionary | ✏️edit |
|---|---|---|---|
| | | ▾ 0 references | |
| | | | ➕ add reference |
| | | | ➕ add value |

| usage example | ⬍ | When the mouse laughs at the cat, there's a hole nearby. (English) ✏️edit | |
|---|---|---|---|
| | | ▾ 0 references | |
| | | | ➕ add reference |

### Navigation / sidebar

Main page
Community portal
Project chat
Create a new Item
Recent changes
Random Item
Query Service
Nearby
Help
Donate

Lexicographical data

Create a new Lexeme
Recent changes
Random Lexeme

Tools

What links here
Related changes
Special pages
Permanent link
Page information
Concept URI
Cite this page
Get shortened URL
Download QR code

Print/export

Download as PDF
Printable version

**Wikidata Lexicodays**
⬙⬙⬙⬙⬙ 28-30 June 2024

# Senses

| | | | |
|---|---|---|---|
| | English | any small rodent of the genus Mus | ✎ edit |
| | Portuguese | qualquer pequeno roedor do gênero Mus | |
| | Persian | پستانداری کوچک از راستهٔ جوندگان | |
| L1119-S1 | Kannada | ಇಲಿಯು ದಂತಕಗಳ ಗಣಕ್ಕೆ ಸೇರಿದ | |
| | Spanish | mamifero roedor de pequeño tamaño, del género Mus | |
| | Hebrew | מכרסם קטן מהסוג Mus | |
| | French | petit rongeur du genre Mus | |

## Statements about L1119-S1

| image |  ✎ edit |
|---|---|
| | Apodemus sylvaticus bosmuis.jpg |
| | 2,160 × 1,524; 2.53 MB |
| | ▾ 0 references |
| | + add reference |
| | + add value |

| item for this sense | mouse | ✎ edit |
|---|---|---|
| | ▾ 0 references | |
| | | + add reference |
| | Mus | ✎ edit |
| | ▾ 0 references | |
| | | + add reference |

| L1119-S2 | English | input device | | edit |
| | Swedish | inmatningsenhet | | |
| | Russian | компьютерная мышь | | |
| | Spanish | dispositivo de entrada | | |
| | Hebrew | | טלק רזיבא | |
| | Kannada | ಪರದೆಯ ಮೇಲ ಕರ್ಸರ್ ಅನ್ನು ನಿಯಂತ್ರಿಸುವ ವಸ್ತು | | |
| | French | périphérique d'entrée informatique | | |

## Statements about L1119-S2

**image**



ComputerMouseCloseup3.jpg
4,928 × 3,264; 3.16 MB

▾ 0 references

╋ add reference

╋ add value

**item for this sense**    mouse      ✎ edit

▾ 0 references

╋ add reference

╋ add value

**translation**    mouse (English) - input device      ✎ edit

▾ 0 references

╋ add reference

**Wikidata
Lexicodays**
▨▨▨ **28-30 June 2024**

# Forms

L1119-F1    mouse
en

Grammatical features   singular

✏ edit

## Statements about L1119-F1

pronunciation audio

▶ 0:00 / 0:00   — 🔊 ⋮
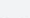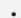
✏ edit

En-uk-a mouse.ogg
1.6 s; 18 KB

pronunciation variety    British English

▾ 0 references

＋ add reference

▶ 0:00 / 0:00   — 🔊 ⋮

✏ edit

En-us-mouse.ogg
0.8 s; 12 KB

pronunciation variety    American English

▾ 0 references

＋ add reference

＋ add value

＋ add statement

L1119-F2    mice
en

✏ edit

Grammatical features   plural

## Statements about L1119-F2

# Why is Lexicographical Data on Wikidata Important?

# Why is it interesting?

- Structured data = machine readable
- Can be reused by tools, research, dictionaries, translation services
- CC0 = open knowledge, can be reused by all
- Huge variety of languages, including underserved ones
- International community = more people to help

# What's the difference to Wiktionary?

- Wiktionary = plain text + templates, Wikidata = structured data
- Wikidata can be easily parsed and reused
- Wikidata works with Lexemes, Wiktionary combines different Lexemes in the same article
- Wikidata = CC0, Wiktionary = CC-BY SA
- Wikidata aims to support Wiktionaries (if they want to)

# What's the difference to other services?

- We're providing the background data to build anything on top of it
- We're doing much more than translation: we help machines understand languages
- We give access to the data in CC0
- We include all languages, not only the most profitable ones
- We empower people to contribute to the data

# What could we
# do with Lexemes?

# Support Wikimedia projects

- Provide structured data to be reused on pages

- Enable working together on the same data

- New tools to make contributing easier

- Text generation, eg. for Abstract Wikipedia

# Abstract Wikipedia

- Human-readable text generation for Wikipedias

- …in different languages

- Based on available data

# Components of Abstract Wikipedia

| Wikidata Items | → | Concepts and the relationships between them | |
| Wikidata Lexemes | → | Words and their meanings and forms | → | Abstract Wikipedia — Abstract content representation | → | Wikipedia article |
| Wikifunctions functions | → | Functions | |

# Dictionary applications

- Looking up definitions and translations

- Special purpose dictionaries (rhyme, specific topics)

- Thesauri and synonym dictionaries

- Build translation tools (especially for underserved languages that don't have any yet)

# The Surrounding Ocean



**Animalia** — *English, proper noun*
**animal** — *English, noun*

**mourningly** — *English, adverb*

**mousepad** — *English, noun*
マウス — *Japanese, noun*

**Maus** — *Bavarian, noun* — small gray or brown mammal with predominantly long tail, rat-like, in the breeding form also white or black

**mouse** — *English, noun*
1. any small rodent of the genus Mus
2. input device

**jengbariga** — *Dagbani, noun* — any small rodent of the genus Mus

**mouselike** — *English, adverb*

Explore words and their meanings, translations, and synonyms: d:Wikidata:The_Surrounding_Ocean

**Wikidata Lexicodays** 28-30 June 2024

# Wikidata Powered Language Keyboards

- Provide translations when typing in a second language

- Assist with Noun genders , verb conjugations & preposition cases

- Word annotation

# Scribe Keyboard



- Keyboards for second language learners
- Currently in French, German, Italian, Portuguese, Russian, Spanish, Swedish, (more languages to be added)



**Noun Genders**
Feminine, masculine, neutral ...

**Translation**
English to keyboard language

**Verb Conjugation**
In-keyboard dictionary

https://github.com/scribe-org

# Other Language learning tools

- Creating word lists and lessons

- Illustrating words

- Creating games and exercises

# Research

- How do languages evolve over time, social class and more?

- Do classes of words change their meaning over time?

- Localizing words on maps

# Text analysis

- Sentiment analysis

- Part of speech tagging

- Named entity recognition

# The State of Lexicographical Data on Wikidata
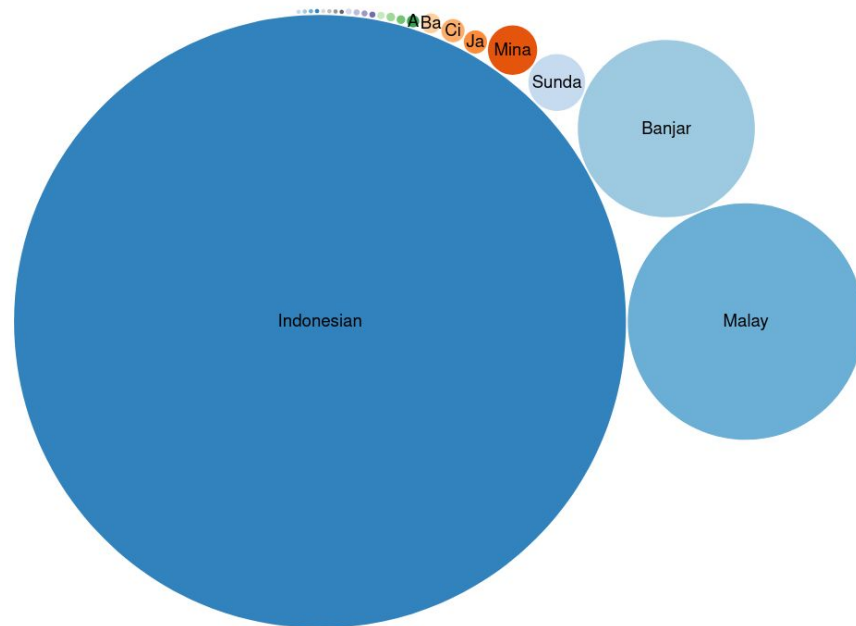
# Distinct languages of Wikidata Lexemes



Try it: https://w.wiki/6RiP

Wikidata
Lexicodays
28-30 June 2024

# Distinct Indonesian languages of Wikidata Lexemes



Try it: https://w.wiki/AJ5U

# Indonesian languages Excl. En & Ar



Try it: https://w.wiki/AJ5u

Wikidata
Lexicodays
28-30 June 2024

# Explore LexData coverage in more detail



Ordia:
https://ordia.toolforge.org/

# Explore LexData coverage in more detail

- Wikidata:Lexicographical coverage:
  [d:Wikidata:Lexicographical_coverage](d:Wikidata:Lexicographical_coverage)

- Wikidata:Lexicographical data/Statistics:
  [d:Wikidata:Lexicographical_data/Statistics](d:Wikidata:Lexicographical_data/Statistics)

- Wikimedia Grafana:
  [https://grafana.wikimedia.org/goto/INWiGyyIg?orgId=1](https://grafana.wikimedia.org/goto/INWiGyyIg?orgId=1)

*To enable truly meaningful applications we need more data (depth and breadth) and people to take care of it!*

# Useful Resources

# Useful links

- Lexicographical Data Overview: Wikidata:Lexicographical_data
  - Try the existing tools Wikidata:Tools/Lexicographical_data
  - Suggest ideas of tools Wikidata:Lexicographical_data/Ideas_of_tools
  - Ideas & examples of queries Wikidata:Lexicographical_data/Ideas_of_queries
- Create and edit Lexemes Special:NewLexeme
- Suggest or discuss new properties Wikidata:Property_proposal/Lexemes
- Discuss with the community about how to model words
  - Wiki talk page: Wikidata_talk:Lexicographical_data
  - Telegram: https://bit.ly/4e8quK8

# Thanks for your attention!

Get in touch with us:

**Lydia Pintscher**
lydia.pintscher@wikimedia.de
@nightrose everywhere

**Mohammed Abdulai**
mohammed.abdulai@wikimedia.de
@masssly everywhere

**Wikidata Lexicodays**
28-30 June 2024