



Lexicographical data

Getting Started with Modeling Words in Your Language

Lydia Pintscher

Mohammed Abdulai

Wikidata Community
Communication Manager,
Wikimedia Deutschland

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

Agenda

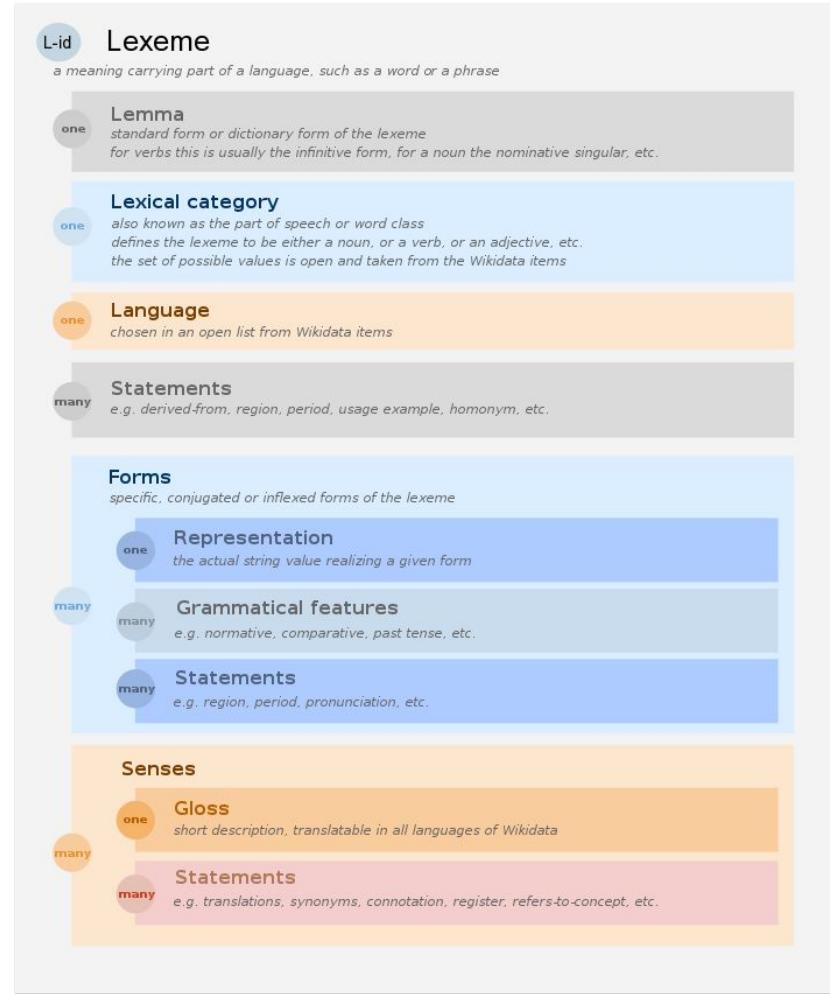
- Lexeme Data Model: A Quick Recap
- Creating Your First Lexeme
- Best Practices For Modeling Words
- Resources And Tools



Lexeme data model: A quick recap

Glossary: Wikidata:Lexicographical data/Glossary

Data model: MediaWiki:Extension:Wiki baseLexeme/Data Model





Creating Your First Lexeme

Search for Lexemes

- Type “L:word” in the search box and press Enter
- Works with lemmas and Forms
- No suggestion appears, don’t worry

Search results

To search for Wikidata items by their title on a given site, use [Special:ItemByTitle](#).

L:Mouse

Advanced search:

Search in: (Main) Property

[mouse \(L1119\)](#)
English, noun...
14 statements, 4 forms - 13:23, 3 June 2024

[mouse \(L332274\)](#)
English, verb...
4 statements, 5 forms - 17:49, 2 May 2024

[moucha \(L19721\)](#)
Czech, noun...
1 statement, 14 forms - 20:06, 22 December 2019



Special page

Search Wikidata



Create a new Lexeme

You can check whether a Lexeme already exists by using [the search](#). You can also learn more about Lexemes in the help box below.

Lemma *

Base form of a word, e.g. 'cat'

Lexeme's language *

The Lexeme's language, e.g. 'English'

Lexical category *

The Lexeme's lexical category, e.g. 'noun'

By clicking "Create Lexeme", you agree to the [terms of use](#), and you irrevocably agree to release your contribution under the [Creative Commons CC0 License](#).

[Create Lexeme](#)

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](#)
[Special pages](#)
[Printable version](#)
[Page information](#)
[Get shortened URL](#)
[Download QR code](#)

About Lexemes

Lexemes contain **lexicographical data** which is data about words or phrases, such as language, etymology, inflections, etc. Here is an example:

cat (L7) cat Language: English
en Lexical category: noun

Lexemes don't contain general data (date of birth, opening date, author, country, coordinates, website, etc.) about the entity or concept to which they refer. If you want to submit general data, you need to [create an Item](#) instead.

[Privacy policy](#) [About Wikidata](#) [Disclaimers](#) [Code of Conduct](#) [Developers](#) [Statistics](#) [Cookie statement](#) [Mobile view](#) [Data access](#)



Go to <https://www.wikidata.org/wiki/Special:NewLexeme>

Create a new Lexeme

You can check whether a Lexeme already exists by using [the search](#). You can also learn more about Lexemes in the help box below.

Lemma *

Lexeme's language *

Lexical category *

By clicking "Create Lexeme", you agree to the [terms of use](#), and you irrevocably agree to release your contribution under the [Creative Commons CC0 License](#).

Create Lexeme

About Lexemes

Lexemes contain [lexicographical data](#) which is data about words or phrases, such as language, etymology, inflections, etc. Here is an example:

cat (L7)	cat	Language: English
	en	Lexical category: noun

Lexemes don't contain general data (date of birth, opening date, author, country, coordinates, website, etc.) about the entity or concept to which they refer. If you want to submit general data, you need to [create an Item](#) instead.

Hit the “Create Lexeme” button



(L4223)

mouse

en



edit

Language [English](#)

Lexical category [noun](#)

Statements

+ add statement

Senses

+ add Sense

Forms

+ add Form

Take note:

Special page

Create a new Lexeme

You can check whether a Lexeme already exists by using [the search](#). You can also learn more about Lexemes in the [Lexeme help page](#).

Not all
languages are
supported for
Lexemes yet!

Lemma *

made-up-word

Lexeme's language *

Tae



Spelling variant of the Lemma * [\(Help\)](#)

Language code for the Lemma's spelling variant, e.g. 'mis-x-Q36790'

Lexical category *

noun

By clicking "Create Lexeme", you agree to the [terms of use](#), and you irrevocably agree to release your contribution under the [Creative Commons CC0 License](#).

[Create Lexeme](#)

Create a new Lexeme

You can check whether a Lexeme already exists by using [the search](#). You can also learn more about Lexemes in the [Lexeme page](#).

Lemma *

made-up-word

Lexeme's language *

Tae

Spelling variant of the Lemma * [\(Help\)](#)

Tae



Lexical category *

noun

By clicking "Create Lexeme", you agree to the [terms of use](#), and you irrevocably agree to release your contribution under the [Creative Commons CC0 License](#).

Create Lexeme

For unsupported languages:

Special page

Create a new Lexeme

You can check whether a Lexeme already exists by using [the search](#). You can also learn more about Lexemes in the [Lexeme page](#).

Select the
“mis”
language
variant

Lemma *

made-up-word

Lexeme's language *

Tae

Spelling variant of the Lemma * [\(Help\)](#)

unsupported language (mis)



Lexical category *

noun

By clicking "Create Lexeme", you agree to the [terms of use](#), and you irrevocably agree to release your contribution under the [Creative Commons CC0 License](#).

[Create Lexeme](#)

(L4222)

made-up-word

mis



edit

Language [Tae](#)

Lexical category [noun](#)

Statements

+ add statement

Senses

+ add Sense

Forms

+ add Form

Now you can...

- Add statements (eg. auxiliary verb)
- Add a Sense
- Add translation (needs another Sense)
- Add a Form (present, first-person singular, etc)

(L1119)

mouse

en



Language English

Lexical category noun

Statements



instance of

count noun



▼ 0 references

+ add reference

+ add value

described by source



Merriam-Webster online dictionary



▼ 0 references

+ add reference

+ add value

usage example

When the mouse laughs at the cat, there's a hole nearby.
(English)

▼ 0 references

Senses

L1119-S1	English	any small rodent of the genus <i>Mus</i>	 edit
	Spanish	mamífero roedor de pequeño tamaño, del género <i>Mus</i>	
	Persian	پستانداری کوچک از راسته جوندگان	
	French	petit rongeur du genre <i>Mus</i>	
	Hebrew	מכרסם קטן מהסוג <i>Mus</i>	
	Kannada	ಇಲೀಯು ದಂಡಕರ್ಗಳ ಗೆಣಕ್ಕೆ ಸೇರಿದ್ದ	
	Portuguese	qualquer pequeno roedor do gênero <i>Mus</i>	

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

Senses

L1119-S1

English	any small rodent of the genus <i>Mus</i>	 edit
Spanish	mamífero roedor de pequeño tamaño, del género <i>Mus</i>	
Persian	بسناداری کوچک از راسته جوندگان	
French	petit rongeur du genre <i>Mus</i>	
Hebrew	מָסֵעַ קָטָן מִהְסָוג <i>Mus</i>	
Kannada	ಇಲೀಯು ದಂಡಕರ್ಗಳ ಗೆಣಕ್ಕೆ ಸೇರಿದ್ದ	
Portuguese	qualquer pequeno roedor do gênero <i>Mus</i>	

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



L1119-F2

mice
en



Grammatical features plural

Statements about L1119-F2

+ add statement

L1119-F3

mouse's
en



Grammatical features genitive case

Statements about L1119-F3

+ add statement

L1119-F4

mouses
en



Grammatical features plural



Pro tips!

- Look at what already exists before creating
- Find an example for inspiration
- Look for the most used properties in Ordia
<https://tools.wmflabs.org/ordia/property/>
- Try, don't be afraid of breaking anything :)
- In doubt: ask the other editors for help ([LexData talk page](#),
[Telegram](#))



Best Practices for Modeling Words

Should There Be a Lexeme for It?

- Evidence of a Lexeme's existence in a language is crucial when creating it on Wikidata.
- For well-documented languages (e.g., English, French), this is obligatory.
- For less documented languages, it's a strong recommendation.
- Evidence can be external identifiers, described-by sources, reference URLs, usage examples, or gloss quotes.

Lemmata

- The lemma of a lexeme should ideally be the representation found in a dictionary.
- This will generally depend on the Lexeme's language and lexical category.
- For languages with multiple scripts, the lemma should represent each script.

Lexical Categories (aka “Parts of Speech”)

- Refer to different types of words in a language
- That help us understand how words function in sentences
- instance of (P31) value on a Lexeme should be more specific than the Lexeme's lexical category.

Senses

- Senses describe the different meanings of a Lexeme
- Therefore, each Lexeme should have at least 1 sense
- Qualify each sense with an item for this sense (P5137) statement where an Item exists

Forms

- Reference at least 1 form. Preferably on a usage example (P5831) statement
- or on another statement, eg.
 - described by source (P1343)
 - attested as (P7855)
 - attested in (P5323) are possible other properties



Frequently asked questions

Can I enter dialects of a Lexeme?

Yes, simply edit
to add a new
lemma.

The screenshot shows a Lexeme page for the entry L1347. It displays three forms: 'colour' (en-gb), 'colour' (en-ca), and 'color' (en-us). Each form has an 'edit' button next to it. Below the forms, there are language and category details: English noun. The 'Statements' section includes an 'instance of' statement pointing to 'loanword', which has 0 references and buttons for adding a reference or a value. The 'derived from lexeme' section lists 'colour' as the mode of derivation, inheritance, with 1 reference, and buttons for adding a value.

Lexeme [Discussion](#) Read View View history [Page](#) More

(L1347) colour colour color
en-gb en-ca en-us [edit](#)

Language English
Lexical category noun

Statements

instance of loanword [edit](#)
▼ 0 references [+ add reference](#)
[+ add value](#)

derived from lexeme colour [edit](#)
mode of derivation inheritance
► 1 reference [+ add value](#)

Can I enter different scripts of a Lexeme?

Yes!

You can add
spellings of the
same word in
multiple
scripts.

Lexeme [Discussion](#) Read [View](#) [View history](#)

(L8303) ruwa | روا | ha-arab [edit](#)

Language Hausa
Lexical category noun

Statements

grammatical gender	masculine edit
	▼ 0 references + add reference
	+ add value

usage example

Ruwa (ana amfanidashi domin inganta rayuwani dan adam (Hausa))	edit
▼ 0 references	+ add reference

What if the dialect or script's language code isn't supported?

Lexeme [Discussion](#)

Read

[View](#)

[View history](#)



(L940072)

مس

bal

mis

bal-x-Q123750890

edit

МИС

bal-x-Q123734863

Use mis-x-Q...
as language of
the lemma

Language [Balochi](#)

Lexical category [noun](#)

Statements

[described by source](#)

Baluchi Glossary

edit

page(s)

111

[1 reference](#)

[+ add value](#)

Does QuickStatements work for Lexemes?

- Not really :(
- Possible: add new statements on the Lexeme level for existing Lexemes
- Not possible: create new Lexemes, Forms or Senses

Is it possible to query Lexemes?

- Yes! The Query Service supports Lexemes
- All words in Indonesian <https://w.wiki/AHRR>
- Longest words in English <https://w.wiki/45o>
- More queries & ideas:
[Wikidata:Lexicographical_data/Ideas_of_queries](#)
- Get help: [Wikidata:Request_a_query](#)



Resources and Tools

Tools to edit Lexemes: Wikidata Lexeme Forms

- Creates Lexemes and generate a lot of Forms
<https://tools.wmflabs.org/lexeme-forms/>
- Possibility to create new templates (eg. conjugations, declinations)

English verb

present

They every day.

third-person singular

He every day.

simple past

He every day last week.

present participle

They are right now.

past participle

We have for hours.

Create

Advanced

Tools to edit Lexemes: Luthor

Finds examples from Wikisource as potential usage examples to Lexemes

Lexeme: bahasa (L6539)

noun

Senses:

1. id: sistem komunikasi tertentu, biasanya dinamai untuk wilayah atau penutur yang menggunakannya

Possible usage examples:

Al-Qur'an

Terjemahan **bahasa** Indonesia dari القرآن (Al-Qur'an) 3496 Terjemahan **bahasa** Indonesia dari القرآن (Al-Qur'an) Departemen Agama RI Terjemahan **bahasa** Indonesia

Undang-Undang Republik Indonesia Nomor 28 Tahun 2014

Program Komputer adalah seperangkat instruksi yang diekspresikan dalam bentuk **bahasa**, kode, skema, atau dalam bentuk apapun yang ditujukan agar komputer bekerja

Citra Manusia Dalam Puisi Indonesia Modern 1920-1960

Tim Penyusun "Citra" Pusat Pembinaan dan Pengembangan **Bahasa** Pusat Pembinaan dan Pengembangan **Bahasa** Departemen Pendidikan dan Kebudayaan Jakarta DAFTAR

Undang-Undang Dasar Negara Republik Indonesia Tahun 1945

Bahasa Negara ialah **Bahasa** Indonesia.

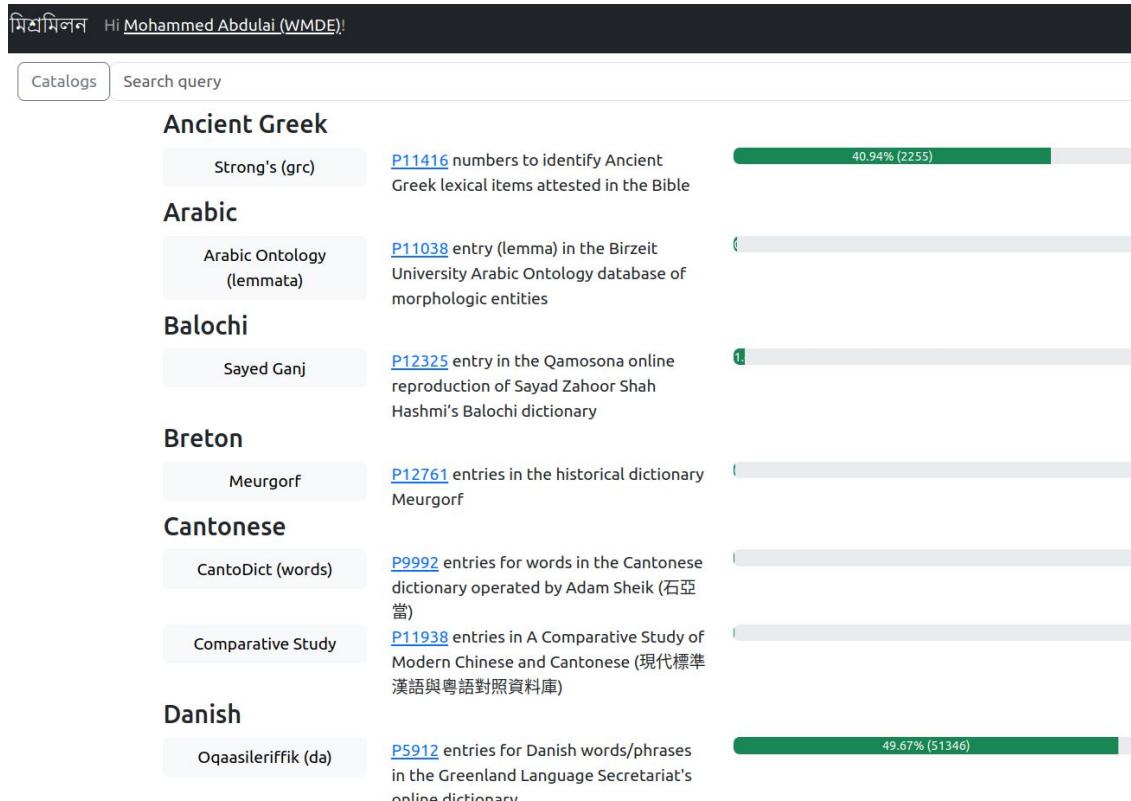
Al-Qur'an/Fussilat

Maha Penyayang. 3 Kitab yang dijelaskan ayat-ayatnya, yakni bacaan dalam **bahasa** Arab, untuk kaum yang mengetahui, 4 yang membawa berita gembira dan yang

<https://luthor.toolforge.org/>

Tools to edit Lexemes: Mishramilan

- Lists entries of some external databases
- Allows users to match those entries against Wikidata lexemes.



<https://mishramilan.toolforge.org/#/>

Tools to edit Lexemes: text-to-lexemes

- Ordia text-to-Lexeme transforms a text into Lexemes and helps with creation <https://tools.wmflabs.org/ordia/text-to-lexemes>

ro
ró

Text to lexemes

A Barata diz que tem sete saias de filó
É mentira da barata, ela tem é uma só
Ah ra ra, iá ro ró, ela tem é uma só !

verb third-person singular // present indicative L41087-S1

ser verb third-person singular // present indicative L39470-S1

é é ser verb third-person singular // present indicative L39470-S2

ro	verb	third-person singular // present indicative	L41087-S1
ró	verb	third-person singular // present indicative	L41087-S1
Ah ra ra, iá ro ró, ela tem é uma só !	verb	third-person singular // present indicative	L41087-S1
é	verb	third-person singular // present indicative	L39470-S1
é	verb	third-person singular // present indicative	L39470-S2

Tools to edit Lexemes: Lingua Libre

- Record words in your language
https://lingualibre.fr/wiki/LinguaLibre:Main_Page
- Generate lists of words from Lexemes
https://lingualibre.fr/wiki/Help:Create_your_own_lists
- Automatically adds the pronunciation as a statement



Lingua Libre

Tools to reuse Lexemes

- The **Query Service**
 - Lexemes don't appear in the documentation... yet ^^
- The **API**

wbladdform

wbladdsense

wbleditformeelements

wbleditsenseelements

wblinktitles

wblmergelexemes

wblremoveform

wblremovesense

How to help

There are many ways to improve lexicographical data on Wikidata!

- 1 First steps, regardless of language
 - 1.1 Check what lexemes already exist in your language
 - 1.2 Check any language-specific documentation for your language
- 2 Tasks for all languages
 - 2.1 Add pronunciation audio to existing lexemes
 - 2.2 Find more specific issues to resolve with Wikidata queries
 - 2.3 Make it easier to add specific types of lexemes
 - 2.4 Add usage examples to existing lexemes
 - 2.5 See if text can be generated using the lexemes in your language
- 3 Tasks for higher-resourced languages
 - 3.1 Add external identifiers to existing lexemes
 - 3.2 Add senses to lexemes that currently don't have them
- 4 Tasks for lower-resourced languages
 - 4.1 Add lexemes for concepts commonly encountered in linguistics
 - 4.2 Add lexemes for concepts from the weekly Lexemes Challenge

[d:Wikidata:Lexicographical_data/How_to_help](#)



Thanks for your attention!

Get in touch with us:

Lydia Pintscher

lydia.pintscher@wikimedia.de
@nightrose

Mohammed Abdulai

mohammed.abdulai@wikimedia.de
@masssly

Credits

This presentation is licensed under CC BY-SA 4.0.