

Danish in Wikidata Lexemes

Finn Årup Nielsen

Cognitive Systems, DTU Compute, Technical University of Denmark

27 July 2019

Wikidata

Article **Berlin** edit
From Wikidata

Capital of Germany edit

Also known as: City of Berlin edit x

Continent	Europe <small>[3 sources]</small>
Country	Germany <small>[2 sources]</small>
Population	3,490,445 <small>[1 source]</small>
	3,500,000 <small>[2 sources]</small>
	<small>[other values]</small>
Calling code	030 <small>[2 sources]</small>
Mayor	Klaus WJ <small>[0 sources]</small>
Vehicle registration	Klaus Wowerit <small>[1 source]</small> <i>German politician</i>
	Klaus Wunderlich <small>[2 sources]</small> <i>German musician</i>
Area	Klaus Waldeck <small>[3 sources]</small> <i>Austrian musician and former lawyer</i>
Twin city	Klaus Wagner <small>[3 sources]</small> <i>German mathematician</i>
	Klaus Wagner <small>[3 sources]</small> <i>Stalker of the British Royal Family</i>

[\[new fact\]](#)

“Wikidata: Verifiable, Linked Open Knowledge That Anyone Can edit” (Dario Taraborelli)

Open data accessible from website, API, dump files and SPARQL endpoint.

Every page is an “item” with label, aliases, properties, property values, and Wikipedia et al. links (Vrandečić and Krötzsch, 2014).

Wikidata-site mockup from 2012 for Berlin (Q64): <https://www.wikidata.org/wiki/Q64>.

Wikidata Lexemes

(L117) **gentagelse** bearbeiten
da

Sprache **Dänisch**
Lexikalische Kategorie **Substantiv**

Aussagen

Genus	Utrum	bearbeiten
	▼ 0 Fundstellen	+ Fundstelle hinzufügen
		+ Wert hinzufügen

Kompositum aus	gentage	bearbeiten
	Ordnungsnummer 1	
	▼ 0 Fundstellen	+ Fundstelle hinzufügen
		+ Wert hinzufügen
	-else	bearbeiten
	Ordnungsnummer 2	
	▼ 0 Fundstellen	+ Fundstelle hinzufügen
		+ Wert hinzufügen

Radikal	tage	bearbeiten
	▼ 0 Fundstellen	+ Fundstelle hinzufügen
		+ Wert hinzufügen

DanNet 2.2 word ID <small>Englisch</small>	11017802	bearbeiten
	▼ 0 Fundstellen	+ Fundstelle hinzufügen
		+ Wert hinzufügen

In 2018, Wikidata introduced a new type of entities for lexemes, their form(s) and sense(s).

Lexemes are prefixed with 'L', e.g., **L117** for the Danish word *gentagelse* (repetition).

On the same page: the sense(s) and form(s) of the lexeme

The lexeme, form and sense (L-Wikidata) may be described by links to the ordinary Wikidata (Q-Wikidata).

Wikidata lexeme basic model

- lemma (wikibase:lemma)
- language (dct:language)
- lexical category (wikibase:lexicalCategory)
- Zero or more statements with property and property values
- Zero or more Senses (ontolex:sense)
 - gloss (skos:definition)
 - Zero or more statements with property and property values
- Zero or more Form (ontolex:lexicalForm)
 - representation (ontolex:representation)
 - Zero or more grammatical features (wikibase:grammaticalFeature)
 - Zero or more statements with property and property values

Wikidata Lexemes RDF

```
wd:L117-F3 a wikibase:Form ,
            ontolex:Form ;
  rdfs:label "gentagelser"@da ;
  ontolex:representation "gentagelser"@da ;
  wikibase:grammaticalFeature wd:Q146786 ,
                               wd:Q53997857 ;
  wdt:P5279 "gen·ta·gel·ser" ;
  wdt:P2859 "\"gEn$%ta:?$@l$s6" ;
  p:P5279 wds:L117-F3-83a6b790-4e1d-56b2-2511-95f1877d046e

wds:L117-F3-83a6b790-4e1d-56b2-2511-95f1877d046e a wikibase:StableForm ;
  wikibase:BestRank ;
  wikibase:rank wikibase:NormalRank ;
  ps:P5279 "gen·ta·gel·ser" .
```

Wikidata lexeme statistics

Count	Description	Query
815724	Number of grammatical feature links	<code>[] wikibase:grammaticalFeature []</code>
356681	Number of forms	<code>[] a ontolex:Form</code>
56518	Number of lexemes	<code>[] a ontolex:LexicalEntry</code>
56518	Number of language links	<code>[] dct:language []</code>
56518	Number of lexical category links	<code>[] wikibase:lexicalCategory []</code>
14196	Number of senses	<code>[] a ontolex:LexicalSense</code>
14196	Number of sense links	<code>[] ontolex:sense []</code>
7586	Number of sense to item links	<code>[] wdt:P5137 []</code>

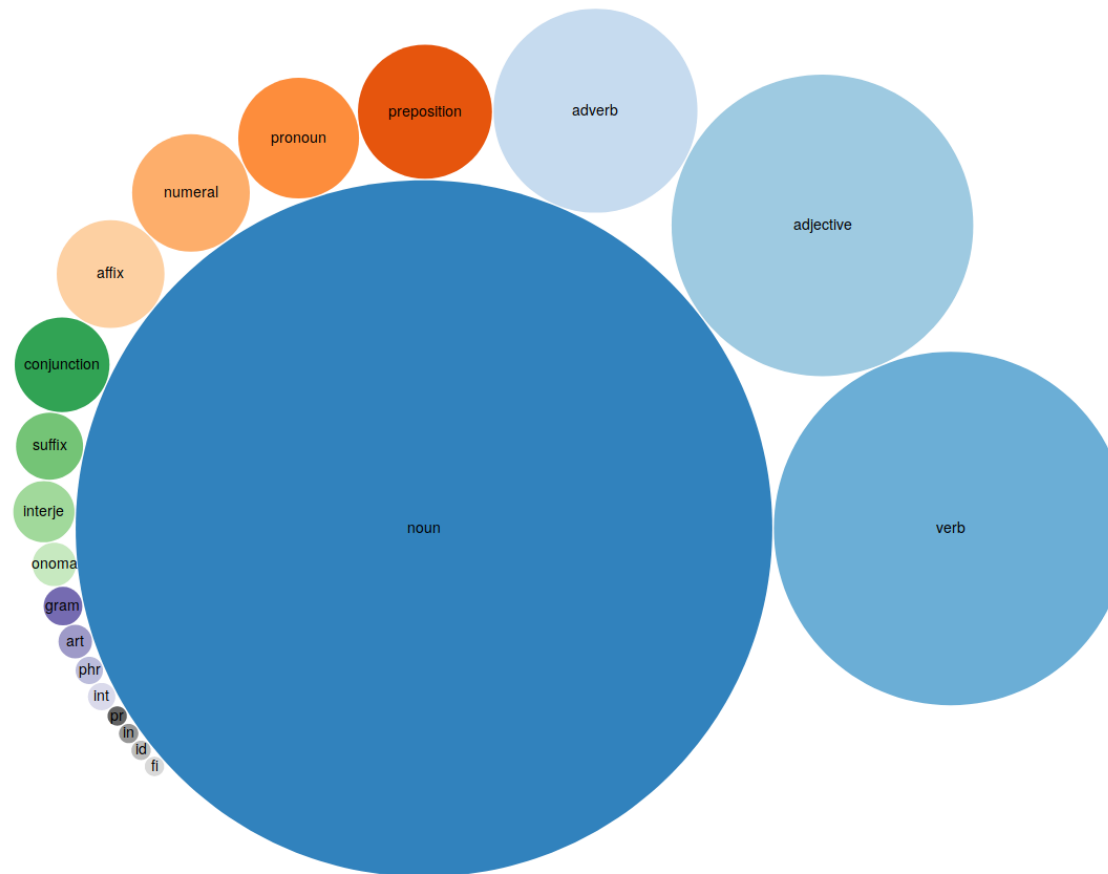
Ordia's statistics: <https://tools.wmflabs.org/ordia/statistics/>

Wikidata lexeme language statistics

Number of lexemes	Language
15980	English
10448	French
7620	Swedish
3021	Basque
2807	Nynorsk
2614	Czech
2453	Polish
2279	German
2207	Danish
721	Japanese

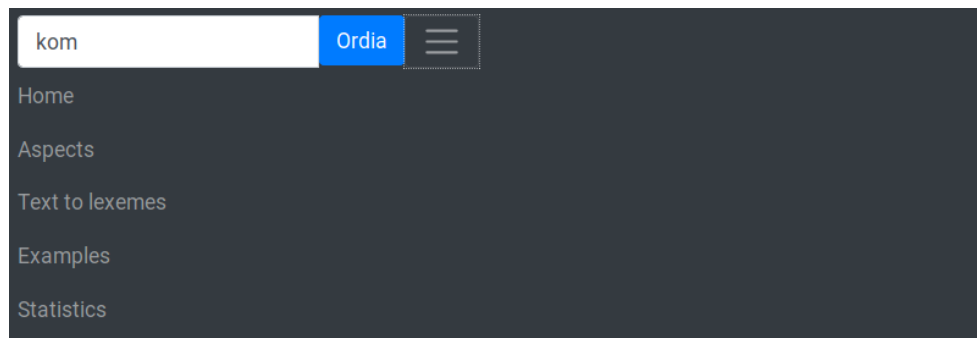
Ordia's language statistics <https://tools.wmflabs.org/ordia/language/>

Wikidata lexeme Danish statistics



<https://tools.wmflabs.org/ordia/language/Q9035>

Wikidata lexeme tools



Search results

- [komme \(L3065\)](#) – Danish, verb
komme (da)
- [komma \(L38134\)](#) – Swedish, verb
komma (sv)
- [ktoś \(L13354\)](#) – Polish, pronoun
komuś (pl)
- [kto \(L23890\)](#) – Polish, pronoun
komu (pl)
- [komst \(L45528\)](#) – Danish, noun
komst (da)
- [kommen \(L46000\)](#) – Danish, noun
kommen (da)
- [komentacja \(L2787\)](#) – Polish, noun
komentacja (pl)

Create new lexeme in Wikidata

Data from [Wikidata](#) | Code from [GitHub repository](#) | Hosted on [Wikimedia Toolforge](#), a [Wikimedia Foundation](#) service | License for content: CC0 for data, CC-BY-SA for text and media | Report technical problems at Ordia's [Issues](#) GitHub page.

Several tools have been built on top of Wikidata: [Wikidata:Tools/Lexicographical_data](#)

Ordia ([Nielsen, 2019](#)) is SPARQL-based webservice at <https://tools.wmflabs.org/ordia/>

Lea Lacroix' DerDieDas game <http://auregann.fr/derdiedas/>

Lucas Werkmeister's form input <https://tools.wmflabs.org/lexeme-forms/>

Alicia Fagerving's Haiki <https://tools.wmflabs.org/hauki/>

Wikidata lexeme properties

Lexemes: “instance of” (values, e.g., countable noun, mass noun, compound word, unadapted loan word), auxiliary verb, DanNet 2.2 word-ID, “derived from”, word stem, image file link, usage example.

Senses: “item for this sense” (very important because it links to the rest of the Wikidata graph with, e.g., hypernym information), image file link, audio file link, translation.

Forms: pronunciation audio file link, X-SAMPA, hyphenation (<https://tools.wmflabs.org/ordia/hyphenation/>).

Wikidata and Linguistic Linked Open Data

Wikidata has deep links to items in the Linguistic Linked Open Data cloud.

Identifier	Count	Property
ILI	27	P5063
BabelNet	60'478	P2581
DanNet word	1'514	P6140

For DanNet we have specified when a Wikidata lexeme does *not* have and associated DanNet 2.2 word identifier with *no value*, see, for instance, *investeringsforvaltningsholdingvirksomhed*.

The generic property *exact match* (P2888) has been used to link to some hundred Princeton WordNet URI identifiers, especially the ImageNet WordNet 3.0 identifiers (Nielsen, 2018).

Text-to-lexemes

Text to lexemes

Blue cars, green bikes and red motorcycles must stop at the crossing.

Language: English

Submit

Extraction

Search:

Word	Form	Lexeme	Lexical category	Features	Sense	Sense image
and	and	and	conjunction		L1385-S1	
at	at	at	preposition			
bikes	bikes	bike	verb	simple present // third-person singular		
bikes	bikes	bike	noun	plural	L10698-S1	
blue	blue	blue	noun	singular		
blue	blue	blue	adjective	positive	L3269-S1	
cars	cars	car	noun	plural	L3648-S1	
crossing	crossing	cross	verb	present participle		
crossing	crossing	crossing	noun	singular	L31093-S1	

Lexeme extraction in Ordia

For a Danish example: “Regeringen spiser grønne æbler om vinteren”.

For “Blue cars, green bikes and red motorcycles must stop at the crossing.”

Validation

Wikidata's property constraint system can indicate that, e.g., an identifier is used twice.

Shape Expressions (ShEx) for lexemes (Nielsen et al., 2019) can specify constraints can specify that Danish noun in definite plural should end with e(r)ne and possible exceptions.

ShEx example with DanNet:

```
<dannet-statement> EXTRA rdf:type {  
  # DanNet identifier should either be novalue or a string  
  ( rdf:type [ wdno:P6140 ] | ps:P6140 xsd:string )  
}
```

Available in a special name space in Wikidata: <https://www.wikidata.org/wiki/EntitySchema:E15>.

Wikidata license

Wikidata uses *Creative Commons Zero* (public domain)

The license of other linguistics resources may create limitations on what we can include in Wikidata.

... use of the Europarl corpus, NST data, ...

Summary

Wikidata lexemes is a very young project.

Expanding from zero lexemes in 2018 to 56528 in July 2019 in multiple language.

Rich annotation scheme for lexemes, forms and senses.

SPARQL queries with visualization possible.

Wikidata Lexeme data validation is possible

CC0 licens for linguistic resources could Wikidata to include the information.

References

- Nielsen, F. Å. (2018). [Linking ImageNet WordNet Synsets with Wikidata](#). *WWW '18 Companion: The 2018 Web Conference Companion, April 23–27, 2018, Lyon, France*, pages 1809–1814. DOI: [10.1145/3184558.3191645](#).
- Nielsen, F. Å. (2019). [Ordia: A Web application for Wikidata lexemes](#).
- Nielsen, F. Å., Thornton, K., and Gayo, J. E. L. (2019). [Validating Danish Wikidata lexemes](#).
- Vrandečić, D. and Krötzsch, M. (2014). [Wikidata: a free collaborative knowledgebase](#). *Communications of the ACM*, 57:78–85. DOI: [10.1145/2629489](#).