

Ordia: A Web application for Wikidata lexemes

Finn Årup Nielsen
Cognitive Systems, DTU Compute, Technical University of Denmark, Kongens Lyngby



Summary

Ordia is a web application for Wikidata lexemes querying the *Wikidata Query Service* on-the-fly when the user navigates the Ordia interface. A user may view individual lexeme, their forms and senses as well as its compound and derivation graph.

Other webpages on Ordia aggregates lexemes for language, lexical category, grammatical feature, properties or references.

Ordia has also special features: text-to-lexemes and integration of knowledge graph embedding similarity results.

Wikidata lexemes

Wikidata has as of May 2019 over 46,000 lexemes. Ordia has a dedicated page that shows live statistics:

Statistics

Count	Description	Query
23624	Number of grammatical feature links	<code>[[wikibase:grammaticalFeature]]</code>
12779	Number of forms	<code>[[en:wikibase:form]]</code>
46193	Number of lexical category links	<code>[[wikibase:lexicalCategory]]</code>
46176	Number of language links	<code>[[en:wikibase:language]]</code>
46168	Number of lemmas	<code>[[en:wikibase:lemma]]</code>
12265	Number of sense links	<code>[[en:wikibase:sense]]</code>
12172	Number of senses	<code>[[en:wikibase:sense]]</code>
6775	Number of sense to form links	<code>[[en:wikibase:sense]]</code>



Ordia's Stack

Python, Flask and Jinja: Standard Python-based web stack.

Wikidata Query Service: Extended SPARQL endpoint at <https://query.wikidata.org/> with visualization capabilities.

JavaScript, JQuery, Data Tables: Formatting the better from the Wikidata Query Service, e.g., for better clickable tables.

GitHub: Repository at <https://github.com/fnielsen/ordia>.

Wikimedia Toolforge: The cloud infrastructure provided by the Wikimedia Foundation hosts the canonical website of Ordia at <https://tools.wmflabs.org/ordia/>.

Ordia is inspired by Scholia.¹

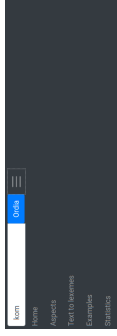
Conceptual choices

The user should easily be able to perform powerful SPARQL queries by navigating the Ordia interface URLs should be predictable, e.g., `/lexical-category/Q24905` shows the Wikidata item for verb (Q24905) as a lexical category, `/property/P31/value/Q639512` is the aspect for instance of (P31) cranberry morpheme (Q639512).

Each page should link to other pages and let the user discover new lexicographic relations. The interface should use graphics whenever possible, e.g., graphs for lexeme and concept relations and for displaying images associated with senses of lexemes.

Searching

Ordia may be used to search Wikidata for lexemes and their forms. Ordia uses the Wikidata API for this function (this is also possible in Wikidata itself via the "L_" prefix).



Search results

- [kornm \(Q3045\)](#) – Danish, verb
- [kornm \(Q4\)](#)
- [korn \(Q133584\)](#) – Polish, pronoun
- [korn \(Q132819\)](#) – Polish, pronoun
- [korn \(Q\)](#)
- [korn \(Q\)](#)
- [korn \(Q\)](#)
- [kornm \(Q4009\)](#) – Danish, noun
- [korn \(Q108\)](#) – Polish, noun
- [kornm \(Q27197\)](#) – Polish, noun
- [kornm \(Q\)](#)

[Create new instance of Wikidata](#)

Data from Wikidata. Code from [Wikidata:Toolforge](#) hosted on [Wikimedia Toolforge](#), a Wikimedia Foundation service. License for content: CC0 for data, CC BY-SA for text and media. Report feedback problems or request content page.



The Javascript *DataTable* package allows the user to filter results displayed in tables in Ordia on the client side.

Knowledge graph embedding

Ordia can show any of the Q-items of Wikidata. These pages show the linked lexemes and senses as well as a list of related Wikidata items obtained through the RDF2Vec-inspired² Wembedder web-service³ which implement a similarity measure via knowledge graph embedding. A table is constructed by interpolating the Wembedder API into a SPARQL query sent to Wikidata Query Service. Here it is for January (Q108), where June and March are the most similar concepts.

Q108

Lexemes	Similarity
January (Q108)	0.999999
June (Q5466)	0.999999
March (Q108)	0.999999
February (Q5466)	0.999999
April (Q5466)	0.999999
May (Q5466)	0.999999
July (Q5466)	0.999999
August (Q5466)	0.999999
September (Q5466)	0.999999
October (Q5466)	0.999999
November (Q5466)	0.999999
December (Q5466)	0.999999

Most similar items

Similarity	Concept	Lexeme
0.999999	January (Q108)	January (Q108)
0.999999	June (Q5466)	June (Q5466)
0.999999	March (Q108)	March (Q108)
0.999999	February (Q5466)	February (Q5466)
0.999999	April (Q5466)	April (Q5466)
0.999999	May (Q5466)	May (Q5466)
0.999999	July (Q5466)	July (Q5466)
0.999999	August (Q5466)	August (Q5466)
0.999999	September (Q5466)	September (Q5466)
0.999999	October (Q5466)	October (Q5466)
0.999999	November (Q5466)	November (Q5466)
0.999999	December (Q5466)	December (Q5466)



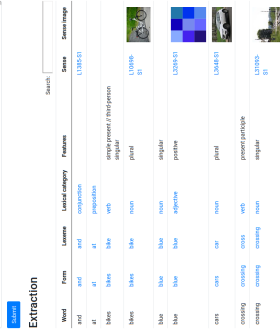
Text-to-lexemes

The text-to-lexemes facility in Ordia enables a user to enter a text in the web interface and let Ordia extract sentences and words.

Extracted words are sent for matching to the Wikidata Query Service and Ordia displays the result of the matching in a table.

Text to lexemes

Enter one or more lines and extract lexemes from stop at all crossing



Here the sentence for text-to-lexemes is "Blue cars, green bikes and red motorcycles must stop at the crossing.", and the first part of the returned table is shown.

In this case all words are matched to lexemes in Wikidata. If some words were not matched then Ordia creates convenient links for entry of the lexemes in Wikidata.



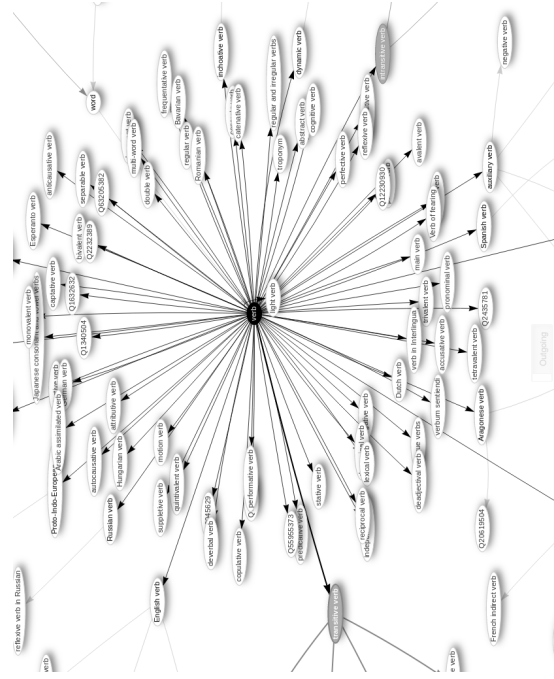
<https://tools.wmflabs.org/ordia/text-to-lexemes?text=Blue+cars+green+bikes+and+red+motorcycles+must+stop+at+the+crossing.&text-language=en>

Aspects

A number of so-called *aspects* show specific information related to: languages, lexical categories, grammatical features, properties or references.

Each of these aspects has a dynamically updated index page (e.g., `/language/Q9035` for language) and dynamically constructed pages for individual Wikidata items (e.g., `/language/Q9035` for the Danish language).

Here is the ontology graph that is part of the *lexical category* aspect for verbs (Q24905)



The black node is the concept of interest (verb) while arrows point to subconcepts (hyponyms).

"Transitive verb" and "intransitive verb" concepts are greyed to indicate that they are used as lexical categories for lexemes in Wikidata.

Current use

Copy-and-paste of a text into Ordia's text-to-lexeme facility to see if all words match Wikidata lexemes, and if not enable the user to quickly create the missing lexemes.

Get an overview of the Wikidata linguistic ontology and its usage, from, e.g., the lexical category aspect.

Searching for lexemes

View statistics, e.g., number of lexemes, number of forms, lexemes per language and lexemes per lexical category.

Future

Current annotation for Wikidata lexemes is too sparse for Ordia to make sense as a cross-language dictionary, etymological dictionary or synonym dictionary.

Possible future extensions could include internationalization of the user interface and Wikidata entry directly in Ordia.

Acknowledgment

This work is funded by the Innovation Fund Denmark through the projects Danish Center for Big Data Analytics driven Innovation (DABAI) and Teaching platform for developing and automatically tracking early stage literacy skills (ATEL).

References

- [1] Nielsen, F.A., Mieschen, D. and Wulffhagen, E. Scholia, Scipitomeres and Wikidata: A Web Application for Wikidata Lexemes
- [2] Ristoski, P. and Paulheim, H. RDF2Vec: RDF Graph Embeddings for Data Mining. The Semantic Web – ISWC 2016, 3016 pages, 498–514
- [3] Nielsen, F.A., Wembedder: Wikidata entity embedding web service, 2017.