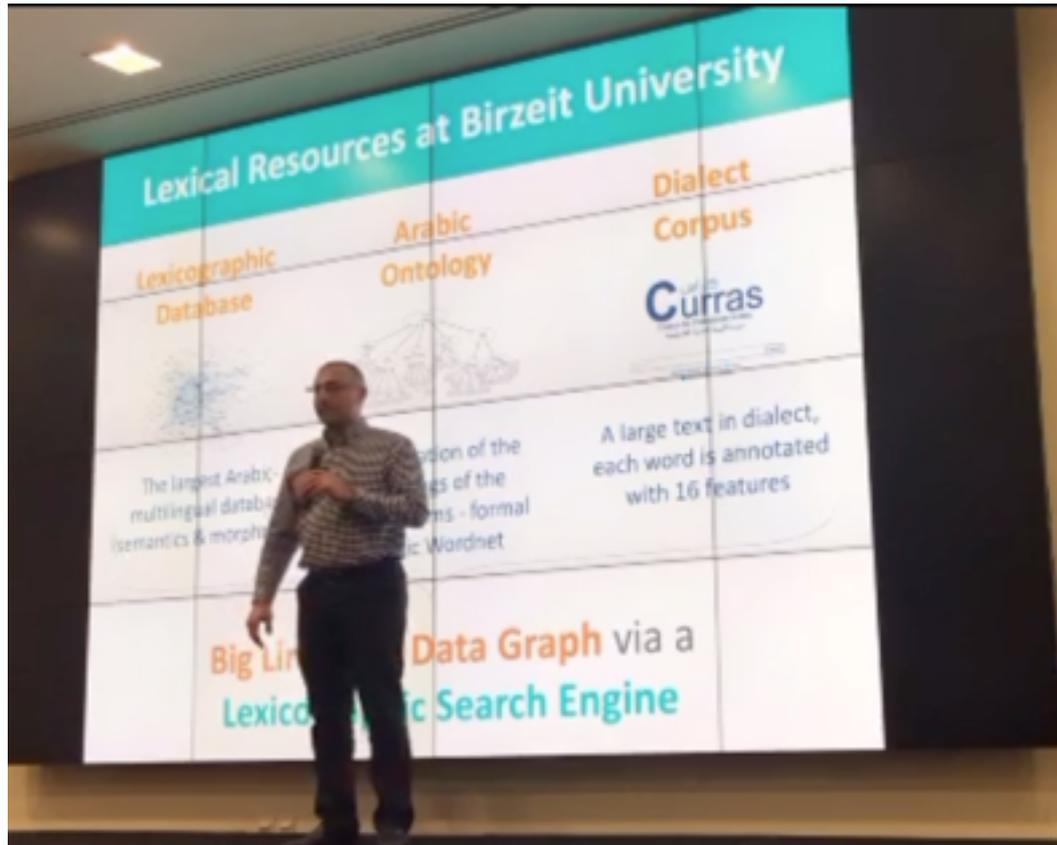


Towards building a **Big Linguistic Data Graph**

Mustafa Jarrar
Birzeit University
Palestine

Watch the lecture here



<https://twitter.com/ArabicSpeech/status/1122843057756475392?s=20>

Why Lexical Resources!

- ❖ The importance of lexical resources (dictionaries, thesauri, wordnets, linguistic ontologies) is increasing in many application areas, such as:
 - NLP tasks and applications
 - Information search and retrieval
 - Multilingual big data
 - Multilingual semantic web
 - Data integration
 - among many others.
- ❖ Lack of Arabic Lexical resources for human use!
- ❖ Lack of Arabic Lexical resources for NLP!

Digitize, Collect, Build, then clean and link

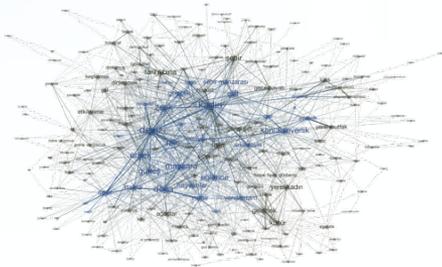
Our Solution



- **Make available online for people**
- **Make available through APIs for NLP applications**

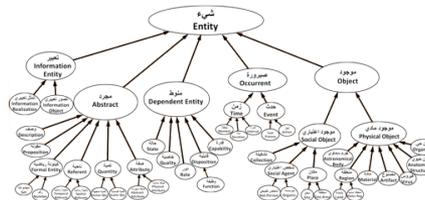
Lexical Resources at Birzeit University

Lexicographic Database



The largest Arabic-multilingual database (semantics & morphology)

Arabic Ontology



Classification of the meanings of the Arabic Terms - formal Arabic Wordnet

Dialect Corpus



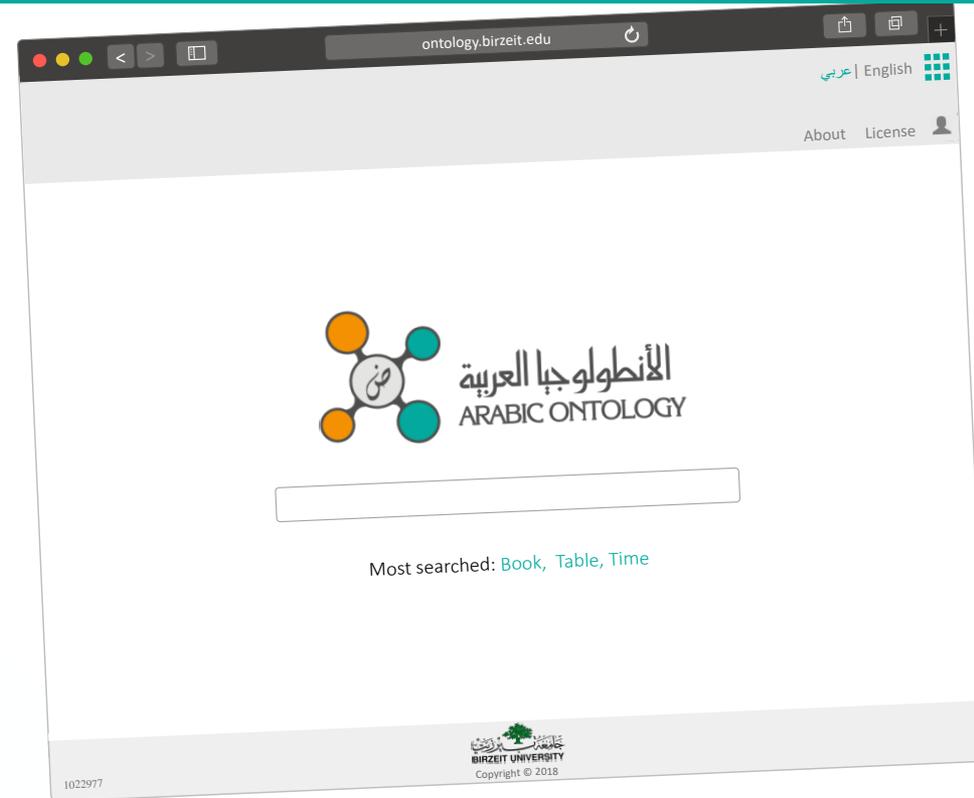
A large text in dialect, each word is annotated with 16 features

Big Linguistic Data Graph via a
Lexicographic Search Engine

Lexicographic Database

Lexicographic Search Engine

- **Free access to people:** students, translators, researchers, Arabic learners ...
- **API accessible** for NLP applications.



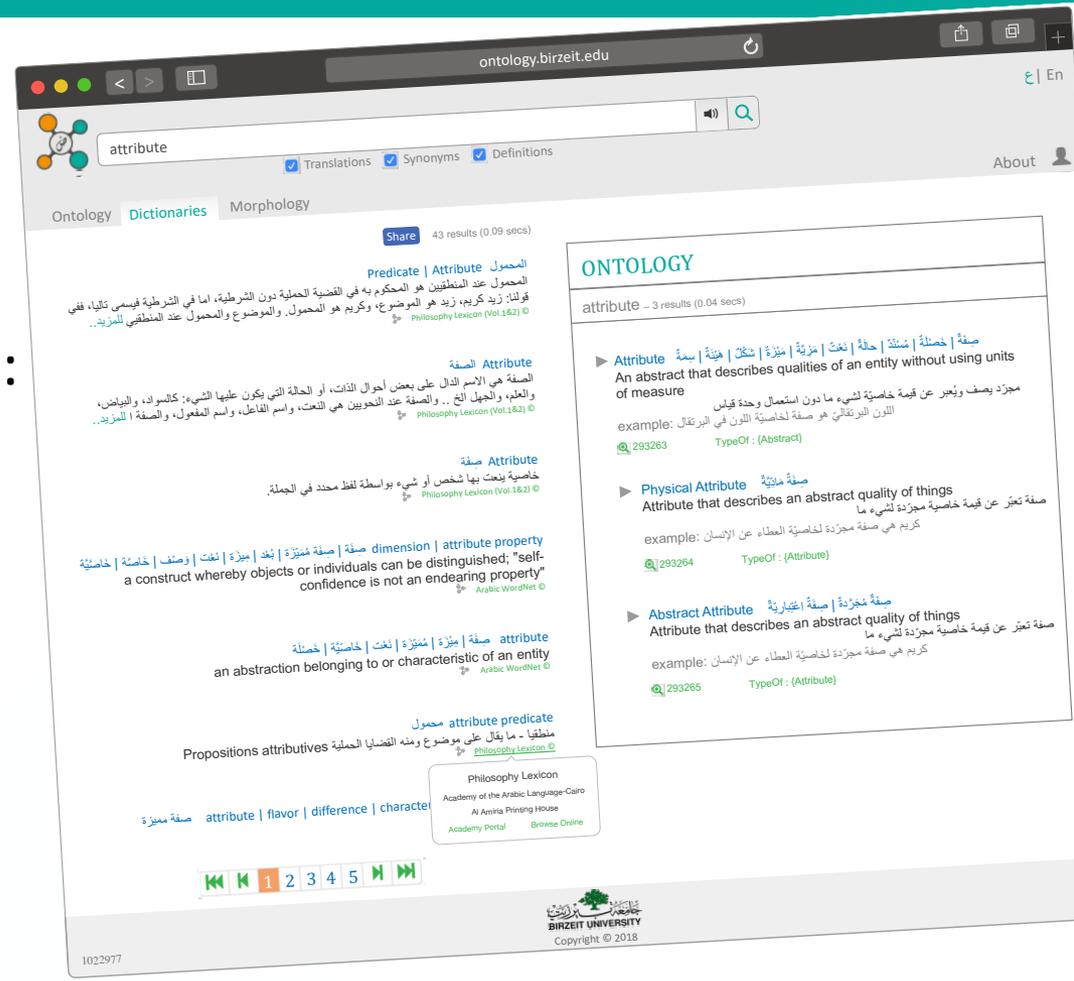
<https://ontology.birzeit.edu>

Based on:

Mustafa Jarrar, Hamzeh Amayreh: **An Arabic Multilingual Database with a Lexicographic Search Engine**. The 24th International Conference on Applications of Natural Language to Information Systems (NLDB 2019), Springer. Manchester, UK. 2019

Lexicographic Search Engine

- The largest lexicographic Arabic database
- Contains most lexicon types: glossaries, thesauri, bi/trilingual dictionaries, morph datasets, **Arabic Ontology**, and more.
- Covers most domains: science, technology, law, business, art, philosophy, ...



<https://ontology.birzeit.edu>

Lexicographic Search Engine

- Search **150 lexicons** for definitions, synonyms and **specialized translations**.
- **Accurate!** compared with machine translation.
- **The first of its kind!** e.g., there are no similar search engines for English lexicons

ontology.birzeit.edu

محرك بحث معجمي
أكبر قاعدة بيانات لغوية:
الأنطولوجيا العربية + 150 معجم.

150 معجم

معاجم عربية ومتعددة اللغات تم حوسبتها وتوحيدها. حاليا يمكن فقط استرجاع دلالة الكلمات (تعريف، مترادفات، وترجمة) وقريبا سيتم استرجاع التصريفات والاشتقاقات.

الأنطولوجيا العربية

شجرة المفاهيم العربية، تصنيف لمعاني الكلمات وليس الكلمات، مصنفة ومعرفة حسب ما توصلت إليه العلوم، وليس حسب ما شاع بين الناس كما المعاجم. انظر [1].

مفاهيم الأنطولوجيا العليا

المفاهيم الأعلى في شجرة الأنطولوجيا هي المفاهيم الأكثر تجزئاً وهي أمهات معاني الكلمات العربية، تم تقعيدها فلسفياً وتعريفها منطقياً. الشكل التالي يعرض أعلى ثلاثة مستويات في شجرة الأنطولوجيا.

```
graph TD
    Entity[شيء Entity] --> InformationEntity[معلومات Information Entity]
    Entity --> Abstract[مجرد Abstract]
    Entity --> DependentEntity[متموِّط Dependent Entity]
    Entity --> Occurrent[موجود Occurrent]
    Entity --> Object[موجود Object]
```

1022977

BIRZEIT UNIVERSITY
Copyright © 2018

DEMO

<https://ontology.birzeit.edu>

Lexicographic Search Engine

Conformance with W3C Standards:

- ✓ **W3C's Best Practices for Publishing Linked Data**
including the Cool URIs, simplicity, stability, and linking
- ✓ **W3C's RDF Lemon Model**

التسوية levelling | grading

تحريك التربة أثناء إعداد الأرض للري للوصول إلى سطح مستو أو سطح ذي انحدار منتظم.

 Hydrology Lexicon ©

```
...
@prefix aot: <http://ontology.birzeit.edu/term/>.
@prefix aoc: <http://ontology.birzeit.edu/lexicalconcept/>.
@prefix aor: <http://ontology.birzeit.edu/lexicon/>.

<aoc:1623> a ontolex:LexicalConcept;
ontolex:isEvokedBy <aot:Lex-grading>;
ontolex:isEvokedBy <aot:Lex-levelling>;
ontolex:isEvokedBy <aot:Lex-تسوية>;
skos:definition "تسوية" @ar;
skos:inScheme <aor:Hydrology_Lexicon_1>.

<aot:lex-grading> a ontolex:LexicalEntry, ontolex:Word;
ontolex:canonicalForm [ontolex:writtenRep "grading"@en];
skos:inScheme <aor:Hydrology_Lexicon_1>.

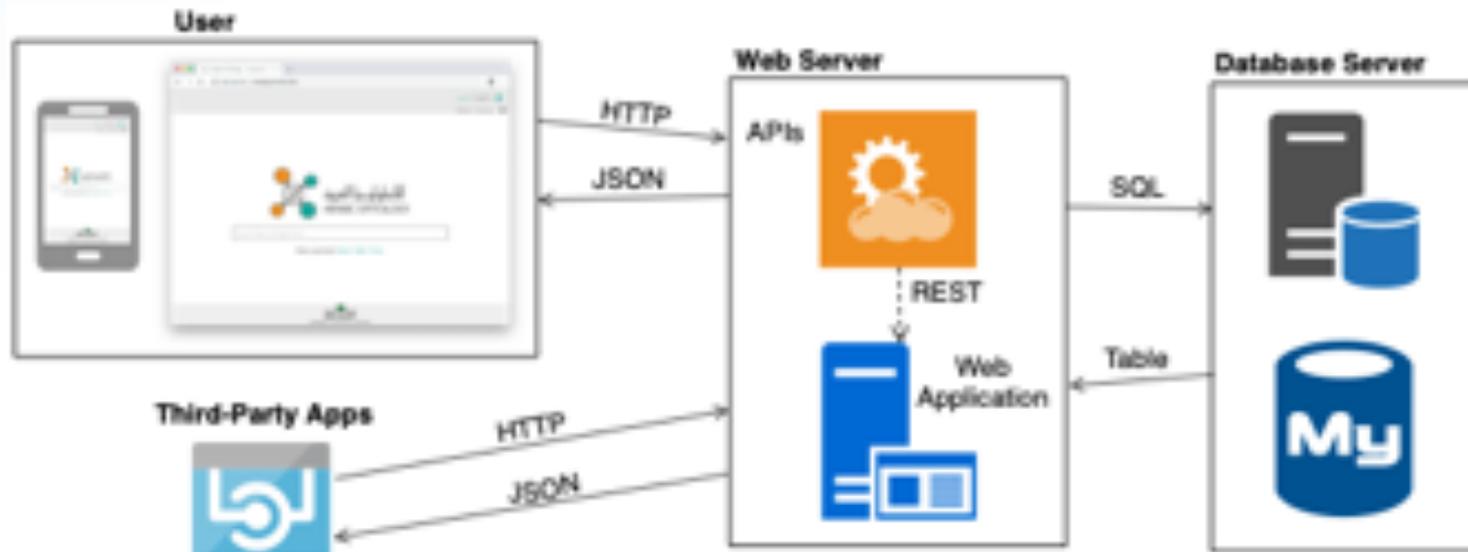
<aot:lex-levelling> a ontolex:LexicalEntry, ontolex:Word;
ontolex:canonicalForm [ontolex:writtenRep "levelling"@en];
skos:inScheme <aor:Hydrology_Lexicon_1>.

<aot:lex-تسوية> a ontolex:LexicalEntry, ontolex:Word;
ontolex:canonicalForm [ontolex:writtenRep "تسوية"@ar];
skos:inScheme <aor:Hydrology_Lexicon_1>.
```

Based On:

Mustafa Jarrar, Hamzeh Amayreh, John McCarae: **Progress on Representing Arabic Lexicons in Lemon**. The 2nd Conference on Language, Data and Knowledge (LDK 2019), Germany. 2019.

Lexicographic Search Engine



Search Engine Architecture

Lexicographic Search Engine

Some Statistics

Category	Lexical Concepts	Lexical entries	Synsets	Translations pairs	Glosses	Semantic relations
Total (Millions)	1.1 M	2.4 M	1.8 M	1.5 M	0.7 M	0.5 M
Sub Counts		1,100K Arabic 1,100K English 200K French 3K Others 1,300K Single-word 1,000K Multi-word	800K Arabic 800K English 200K French 50K Others	1,000K English-Arabic 300K English-French 200K French-Arabic	400K Arabic 300K English 1K Others	170K Sub-super links 29K Part-of links 260K Has-Domain links 30K Other links

Lexicographic Search Engine

API Access

RESTful web services

Ask us for an API Key!

LexAPI v1.0

LexAPI 1.0 is a set of RESTful webservice that all together form an API for other third-party software developers to retrieve linguistic data from the [lexicographic search engine](#).

This page explains APIs with example links on each. A click on one of the links will send the request to the corresponding API and the returned JSON object will appear inside the Output box on the right.

APIs:

- + Search Dictionaries for a term:
- + Search Arabic Ontology for a term:
- + Retrieve a lexical concept:
- + Retrieve an Arabic Ontology concept:
- + Retrieve Morphology information:
- + Autocomplete Service:
- + Retrieve subtypes of an Onotlogy concept:
- + Retrieve concepts part of another concept:

Output (JSON):

```
{
  "conceptID":1520039900,"arabicGloss":null,"englishGloss":"the fourth coordinate that is required (along with three spatial dimensions) to specify a physical event","tags":null,"example":null,"lang":null,"dataSourceId":152,"synsetFrequency":null,"dataSourceCacheAr":"شبكة المفردات العربية","dataSourceCacheEn":"Arabic WordNet","arabicWordsCache":"تُغد زايغ | وقت | زمن","englishWordsCache":"fourth dimension | time","superId":1520039870,"superOrder":0,"superTypeCacheAr":"تُغد","superTypeCacheEn":"dimension","categoryId":null,"area":null,"era":null,"rank":null,"status":null,"subTypesCount":0,"partOfCondiacritizedArabicWordsCache":null,"unacritizedArabicWordsCache":null,"normalizedEnglishWordsCache":"fourth dimension | time |","exactWord":null}
```


Arabic Ontology

Arabic Ontology

- Classification of the meanings of the Arabic terms, specified in D. Logic
- Can be used as a formal Arabic Wordnet -with ontologically-clean content.
- Linked with WordNet
- Benchmarked to scientific advances rather than to speakers' naïve beliefs as wordnets do.

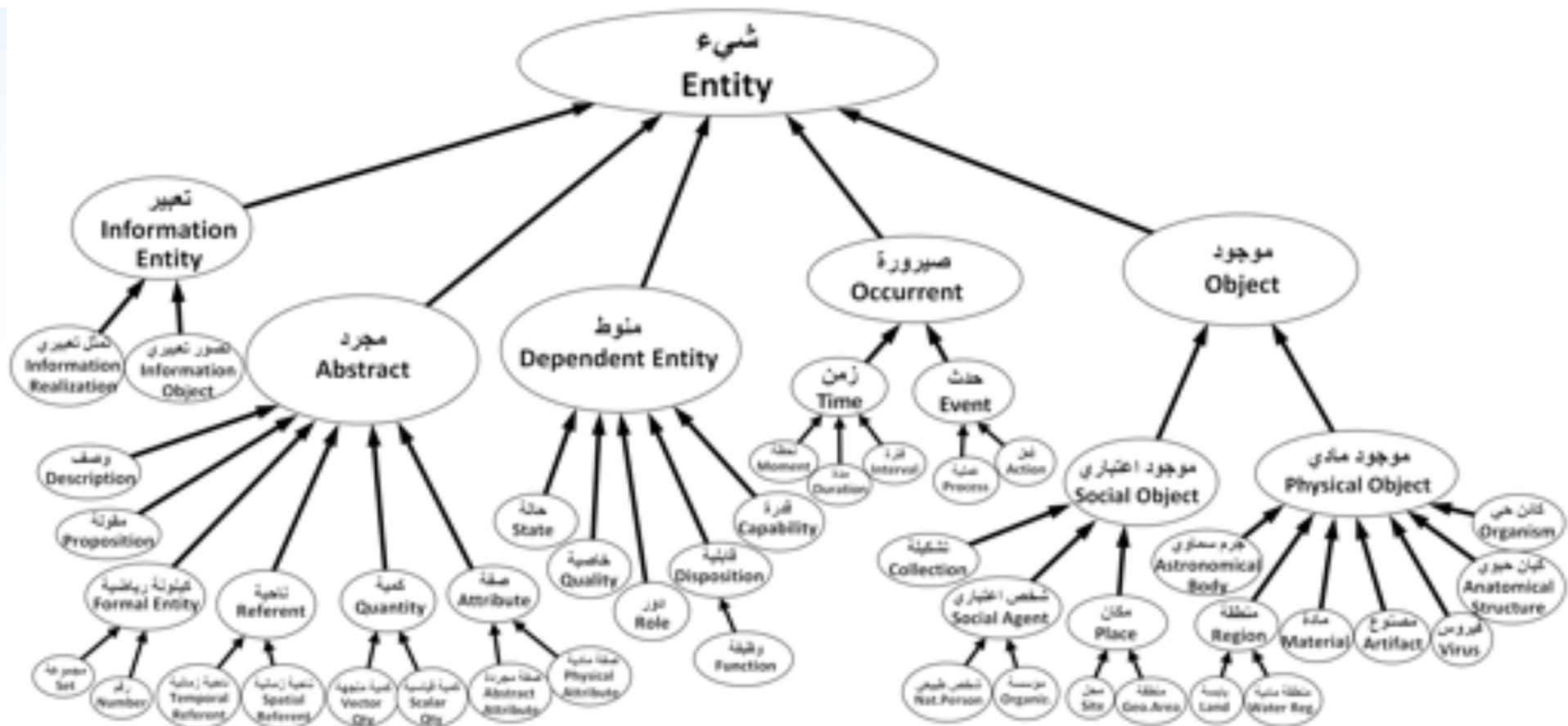
The screenshot shows the Arabic Ontology website interface. At the top, there is a search bar and navigation links for 'Translations', 'Synonyms', and 'Definitions'. Below the search bar, there are tabs for 'Ontology', 'Dictionaries', and 'Morphology'. The main content area displays the definition of 'Entity' (شيء | كَيْنُونَةٌ | كَائِنٌ) as 'Whatever existed or will exist, and can be realized or imagined'. Below this, there are four subtypes of Entity, each with its own definition and example:

- Object** (مَوْجُود | كَائِنٌ | قَائِمٌ | حَقِيقِيٌّ | واقِعِيٌّ | شيءٌ | ذاتٌ | فَيُومٌ): An entity that is wholly and independently present in time, and is realized either for its concrete or social existence. Example: كَلَّ شَيْءٌ عَلَى مَا يَرَامُ (293198).
- Occurrent** (صَنْبُورَةٌ | حَدَثٌ | حَادِثٌ | وَقَعٌ | أَمْرٌ): An entity realized by the time of its happening. Example: شَيْءٌ يَدْرِكُ ذَاتَهُ بِزَمَنِ حَدُوثِهِ لَا يُمْكِنُ فَهْمُ أَيِّ حَدَثٍ بِشَكْلِ مَنْفَصِلٍ عَنِ الْإِطَارِ الزَّمْنِيِّ لَهُ (293202).
- Dependent Entity** (مُنَوِّطٌ | مُعْتَمِدٌ | مُتَعَلِّقٌ | مَشْرُوطٌ): An entity realized by the time of its happening. Example: طُولُ الْمَبْنِيِّ مُنَوِّطٌ بِوُجُودِ الْمَبْنِيِّ وَإِلَّا فَلَا طُولَ لَهُ (293201).
- Abstract** (مُجَرَّدٌ | تَجْرِيدِيٌّ | غَيْرُ مَادِّيٍّ | نَظَرِيٌّ): An entity exists only in mind, cannot be measured or socially realized, and does not have a location. Example: (293198).

The footer of the website includes the Birzeit University logo and the text 'Copyright © 2018'.

<https://ontology.birzeit.edu/concept/293198>

Arabic Ontology



Based on:

Mustafa Jarrar: **The Arabic Ontology - An Arabic Wordnet with Ontologically Clean Content.** Applied Ontology Journal, IOS Press. (Forthcoming).

Dialect Corpus

Dialect Corpus

<http://portal.sina.birzeit.edu/curras>

- Collect a corpus written in Palestinian dialect.
- Describe and annotate each word with 16 tags.
- Stats:
 - 60k tokens,
 - 16K unique tokens



Based on:

Mustafa Jarrar, Nizar Habash, Faeq Alrimawi, Diyam Akra, Nasser Zalmout: Curras: **An Annotated Corpus for the Palestinian Arabic Dialect**. Journal Language Resources and Evaluation. Volume(51), Issue(3). Springer. 2017

Dialect Corpus

Middle size!

Size	Resource	Tokens (%)	Documents
	Facebook	4852 (8.6)	35 threads
	Twitter	4694 (8.3)	38 threads
	Blogs	11,245 (19.8)	37 threads
	Forums	1027 (1.8)	33 threads
	Palestinian stories	3149 (5.6)	6 stories
	Palestinian terms	1468 (2.6)	1 doc
	TV Shows: وطن ع وئر <i>Watan Aa Watar</i>	30,265 (53.4)	41 episodes
	Total	56,700	190

Accuracy of annotators

Category	A1 accuracy	A2 accuracy
Complex-prefix	97.4	97.6
Stem	86.3	93.1
Complex-suffix	92.5	96.0
Surface (CODA)	90.8	97.8
Lemma	93.4	93.1
Lemma _{MSA}	87.5	87.2
Gloss	74.6	80.9
POS	94.1	95.7
Person	97.8	98.8
Aspect	99.3	99.1
Gender	86.2	96.9
Number	94.8	96.8

Accurate



Connecting
lexical resources

W3C Lemon RDF Model Standard

Lexicon Model for Ontologies: Community Report, 10 May 2016



Final Community Group Report 10 May 2016

Editors:

[Philipp Cimiano](#) ([Cognitive Interaction Technology Excellence Center, Bielefeld University](#))

[John P. McCrae](#) ([Insight Centre for Data Analytics, National University of Ireland, Galway](#))

[Paul Buitelaar](#) ([Insight Centre for Data Analytics, National University of Ireland, Galway](#))

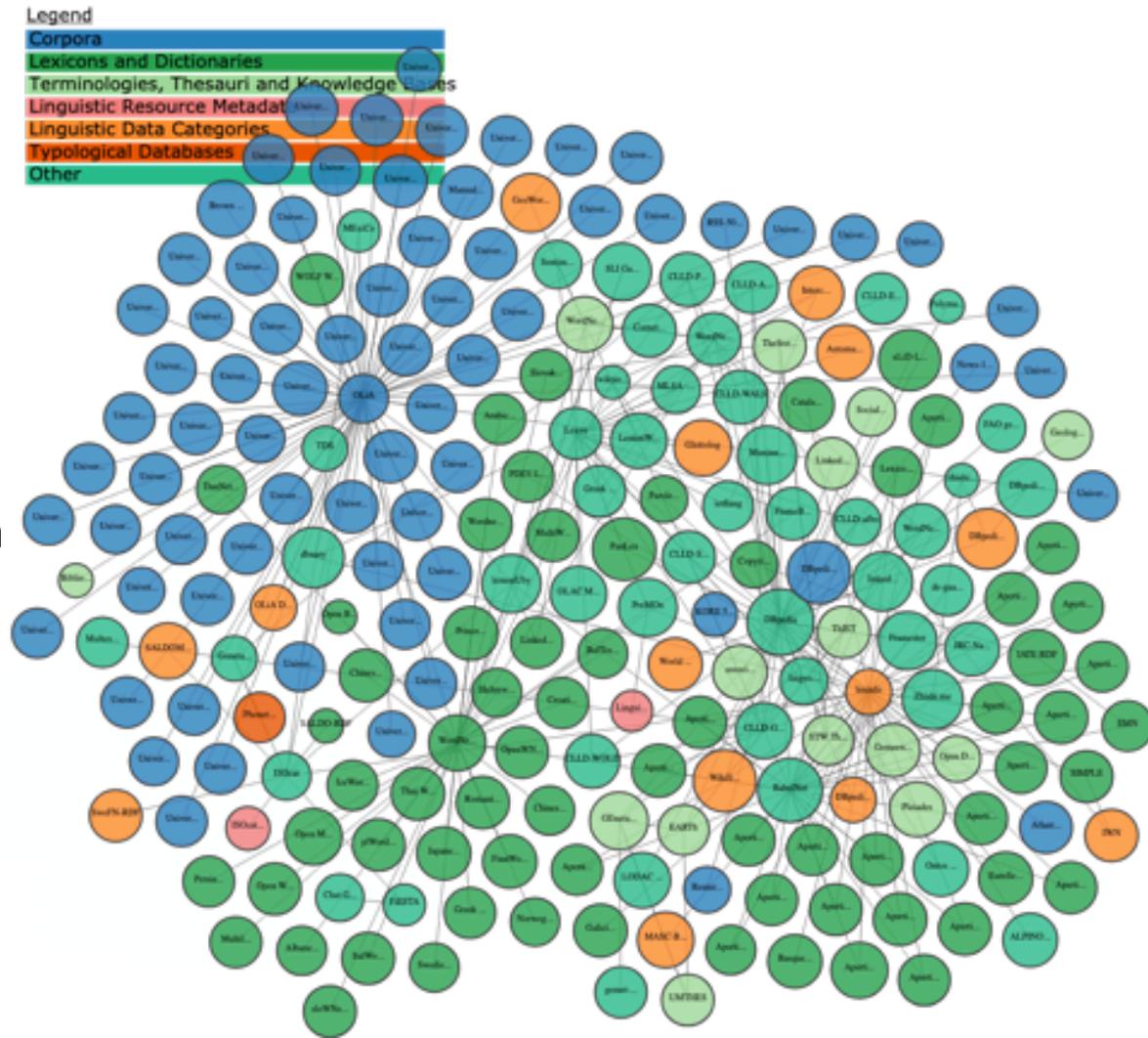
[Copyright](#) © 2016 the Contributors to the Lexicon Model for Ontologies: Community Report, 10 May 2016 Specification, published by the [Ontology-Lexicon Community Group](#) under the [W3C Community Final Specification Agreement \(FSA\)](#). A human-readable [summary](#) is available.

Abstract

This document describes the lexicon model for ontologies (*lemon*) as a main outcome of the work of the Ontology Lexicon (Ontolex) community group.

The Linguistic Linked Open Data Cloud

- A collaborative effort to develop a Linked Open Data (sub-)cloud of linguistic resources.
- Represent (lexical entries, concepts, synsets, and other) using Lemon RDF model, then interlinked.



Work in Progress

- Lemon-izing and Interlinking Arabic resources with the Linguistic Linked Open Data Cloud
- Building an Arabic Knowledge Graph
- Collect and interlink dialect corpora

References

1. Mustafa Jarrar. The Arabic Ontology - An Arabic Wordnet with Ontologically Clean Content. Applied Ontology Journal, 2019 [Forthcoming].
2. Mustafa Jarrar and Hamzeh Amayreh. An Arabic-Multilingual Database with a Lexicographic Search Engine. Proceedings of the 24th International Conference on Applications of Natural Language to Information Systems (NLDB), 2019.
Mustafa Jarrar, Hamzeh Amayreh: An Arabic Multilingual Database with a Lexicographic Search Engine. The 24th International Conference on Applications of Natural Language to Information Systems (NLDB 2019), Springer. Manchester, UK. 2019
3. Mustafa Jarrar: Search Engine for Arabic Lexicons. Proceedings of the 5th Conference on Translation and the Problematics of Cross-cultural Understanding. The Forum for Arab and International Relations. Doha, Qatar. 2018
4. Hamzeh Amayreh, Mohammad Dwaikat, and Mustafa Jarrar. Lexicon Digitization -A Framework for Structuring, Normalizing and Cleaning Lexical Entries, 2018.
5. Mustafa Jarrar, Fadi Zaraket, Rami Asia, Hamzeh Amayreh: Diacritic-Based Matching of Arabic Words. ACM Asian and Low-Resource Language Information Processing. Volume 18, No 2, Pages(10:1-10:21), ACM, December 2018.
6. Mustafa Jarrar, Nizar Habash, Faeq Alrimawi, Diyam Akra, and Nasser Zalmout. Curras: An Annotated Corpus for the Palestinian Arabic Dialect. Journal Language Resources and Evaluation, 51(3):745–775, 2017.
7. Mustafa Jarrar, Nizar Habash, Diyam Akra, Nasser Zalmout: Building a Corpus for Palestinian Arabic: a Preliminary Study. In proceedings of the EMNLP 2014 Workshop on Arabic Natural Language Processing. Association for Computational Linguistics (ACL), Pages (18-27). October 25, 2014, Doha, Qatar. ISBN: 978-1-937284-96-1
8. Mustafa Jarrar. Building a Formal Arabic Ontology (Invited Paper). In Proceedings of the Experts Meeting on Arabic Ontologies and Semantic Networks. ALECSO, Arab League, 2011