# An Arabic-Multilingual Database with
# a Lexicographic Search Engine

**Mustafa Jarrar**
Birzeit University
Palestine

# An Arabic-Multilingual Database with a Lexicographic Search Engine

Mustafa Jarrar(✉) and Hamzeh Amayreh

Birzeit University, Birzeit, Palestine
mjarrar@birzeit.edu, hamayreh@staff.birzeit.edu

**Abstract.** We present a lexicographic search engine built on top of the largest Arabic multilingual database, allowing people to search and retrieve translations, synonyms, definitions, and more. The database currently contains about 150 Arabic multilingual lexicons that we have been digitizing, restructuring, and normalizing over 9 years. It comprises most types of lexical resources, such as modern and classical lexicons, thesauri, glossaries, lexicographic datasets, and (bi/)tri-lingual dictionaries. This is in addition to the Arabic Ontology – an Arabic WordNet with ontologically cleaned content, which is being used to reference and interlink lexical concepts. The search engine was developed with the state-of-the-art design features and according to the W3C's recommendation and best practices for publishing data on the web, as well as the W3C's Lemon RDF model. The search engine is publicly available at (https://ontology.birzeit.edu).

## 1 Introduction and Motivation

The increasing demands to use and reuse dictionaries (of all types) in modern appli-

❖ 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!
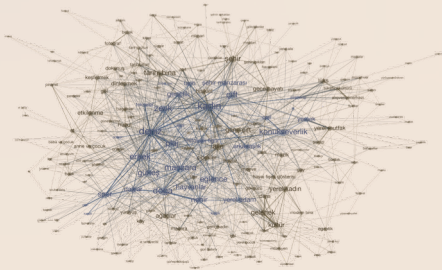
# Digitize, Collect, Build, then clean and link

**Solution**



➤ **Make available online for people**

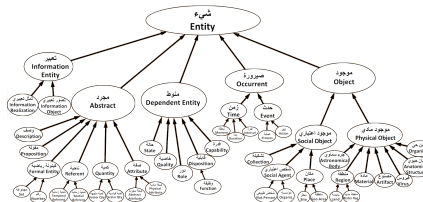➤ **Make available through APIs for NLP applications**

# 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
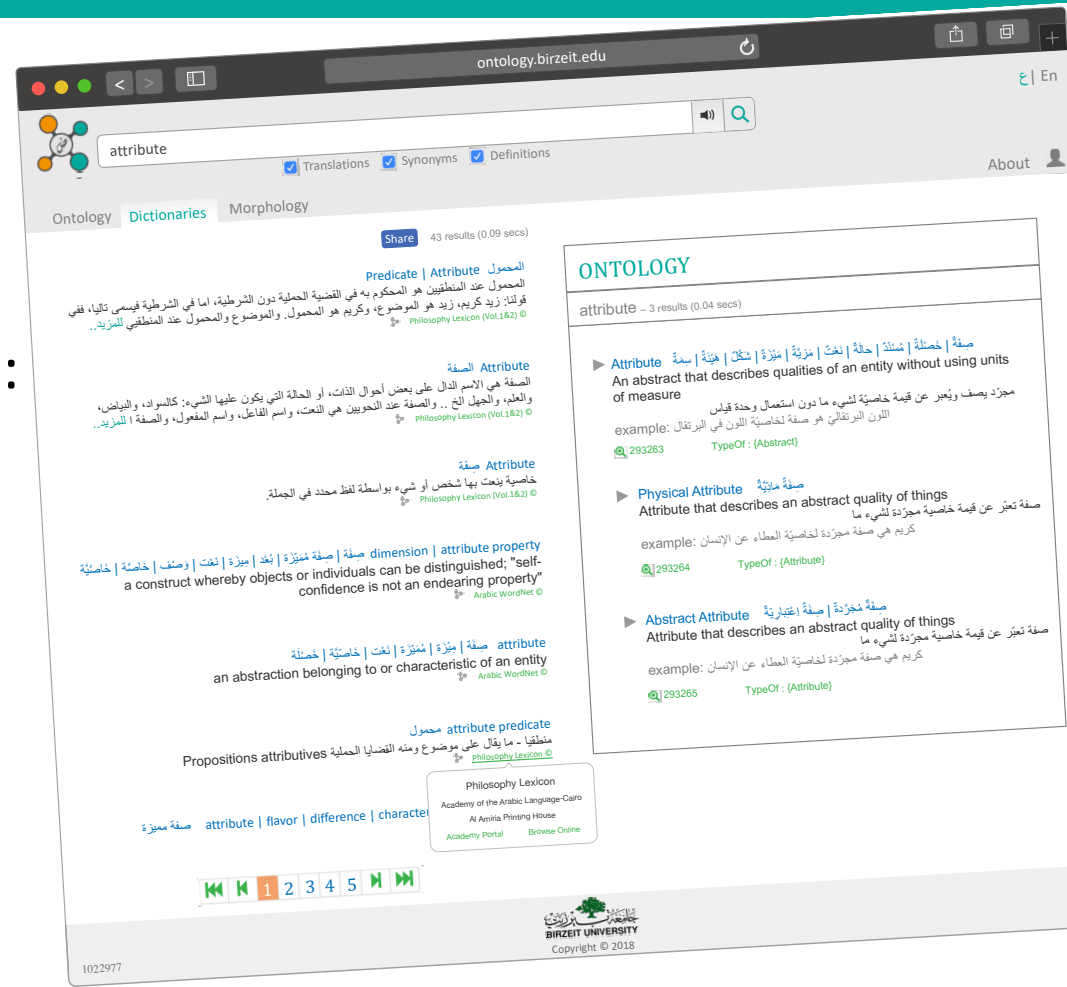
## Dialect Corpus



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

**Linguistic Big Data Graph** via a **Lexicographic Search Engine**

# The Lexicographic

# Database

# The Lexicographic Database

- **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, ...
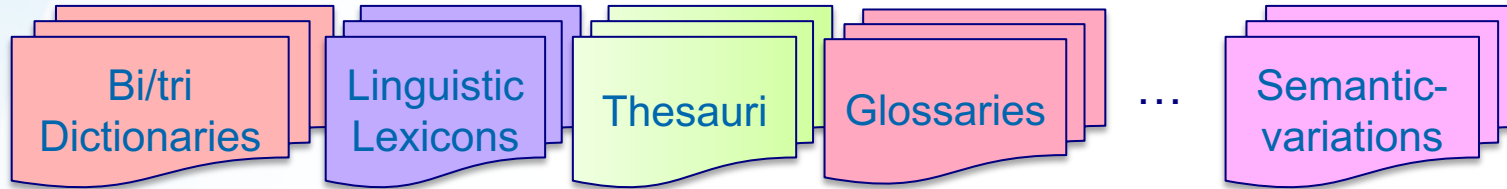


https://ontology.birzeit.edu

# Some Statistics

Currently!

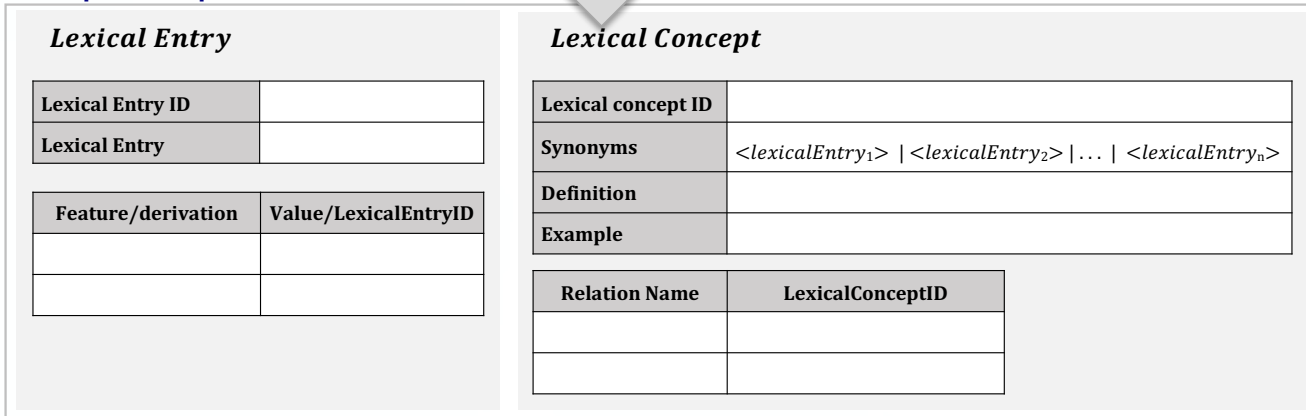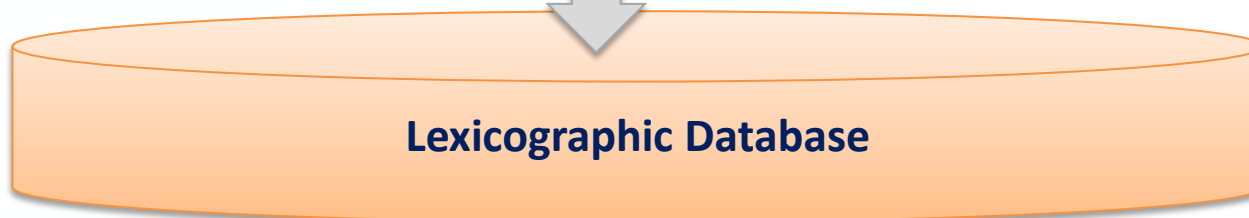| Category | Lexical Concepts | Lexical entries | Synsets | Translations pairs | Glosses | Semantic relations |
|---|---|---|---|---|---|---|
| **Total (Millions)** | 1.1 M | 2.4 M | 1.8 M | 1.5 M | 0.7 M | 0.5 M |
| **Sub Counts** | | 1,100K Arabic<br>1,100K English<br>200K French<br>3K Others<br>1,300K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

# Constructing the Database (9 years)

Bi/tri Dictionaries

Linguistic Lexicons

Thesauri

Glossaries

…

Semantic-variations

manual digitization 150 lexicons

multilingual, semantics, morphology, features…

## Temp Templates

### Lexical Entry

| Lexical Entry ID | |
|---|---|
| Lexical Entry | |

| Feature/derivation | Value/LexicalEntryID |
|---|---|
| | |
| | |

### Lexical Concept

| Lexical concept ID | |
|---|---|
| Synonyms | $<lexicalEntry_1> | <lexicalEntry_2> | \ldots | <lexicalEntry_n>$ |
| Definition | |
| Example | |

| Relation Name | LexicalConceptID |
|---|---|
| | |
| | |

**Lexicographic Database**

Semi-automatic cleaning and normalization

# Cleaning and Normalization

- Lexicons are typically designed to be printed and used as hard copies.

- Big challenges when converting them into a machine processable format.

## Examples of challenges:

*Challenges induced by ordering:*

"accelerator (linear...)", "affinity (chemical)", "drawing (final)", "earth (the)", and "crush (to)", "tube (pipe)", "academy (of art)"

*Subterm synonymy:*

"liquid drier, drier", "calomel electrode, calomel", "kelvin's scale, kelvin's absolute scale"

*Long multiword lexical entries:*

"buildings or other structures recurrent taxes on land". Such cases of long and "poor"

*Character set*: Same characters and symbols have different encodings across different languages (e.g., the dash, quotations, punctuations, and whitespaces),

➤ **See 30 parsers at:**

Hamzeh Amayreh, Mohammad Dwaikat, and Mustafa Jarrar. **Lexicon Digitization -A Framework for Structuring, Normalizing and Cleaning Lexical Entries**, Technical Report, Birzeit University 2018.

# Obtaining Copyrights

- Obtained permission from each lexicons owner (individually contacted).

- Most accepted!

- Show lexicon name and © copyright symbol beside each result.

- Promote lexicons (click to see lexicon info)

# Lexicographic
# Search Engine

# Lexicographic Search Engine

- **Free access to people**: students, translators, researchers, Arabic learners …

- **API accessible** for NLP applications.



## https://ontology.birzeit.edu

Reference:

Mustafa Jarrar, Hamzeh Amayreh: **An Arabic-Multilingual Database with a Lexicographic Search Engine.** NLDB 2019. Pages(234--246), LNCS 11608, Springer. 2019.

# Lexicographic Search Engine

- **Search 150 lexicons** for definitions, synonyms, specialized translations, morphology, ontology...

- **Accurate**! compared with machine translation.

- **The first of its kind**! e.g., there are no similar search engines for English lexicons!



https://ontology.birzeit.edu

# Search Taps

**Ontology tab**: results in this tab are ontology concepts retrieved only from the Arabic ontology. The tab also allows expanding and exploring the ontology tree.

# Search Taps

**Dictionaries tab**: results in this tab are lexical concepts retrieved from the lexicons.

# Search Taps

**Morphology tab:** results are linguistic features, lemma(s), inflections, and derivations of the searched term (partially implemented!).

# Search Engine Architecture

# Conformance with W3C Standards

✓ **W3C's RDF Lemon Model**

Represent (lexical entries, concepts, synsets, …) using the Lemon RDF model

To interlink it with the Linguistic Linked Open Data Cloud

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 "...تحريك التربة أثناء إعداد الأرض للري للوصول إلى سطح مستو أو سطح"@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.

# Conformance with W3C Standards

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

**URLs Schema:**

- Each term is given a URL: `http://{domain}/term/{term}`

  http://ontology.birzeit.edu/term/virus

- Each lexical concept is given a URL:

  `http://{domain}/lexicalconcept/{lexicalConceptID}`

  https://ontology.birzeit.edu/lexicalconcept/304000682

- Each concept in the Arabic Ontology has a URL:

  `http://{domain}/concept/{ConceptID | Term}`

  https://ontology.birzeit.edu/concept/293262

- Each Semantic relation is given a URL:

  `http://{domain}/concept/{RelationName}/{ConceptID}`

  https://ontology.birzeit.edu/concept/instances/293121

- The W3C Lemon representation of each lexical concept is given a

  URL: `http://{domain}/lemon/lexicalconcept/{lexicalConceptID}`

  https://ontology.birzeit.edu/lemon/lexicalconcept/304000682

# API Access

RESTful web services

Ask us for an API Key!



## LexAPI v1.0

LexAPI 1.0 is a set of RESTful webservices 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,"synsetFrequnecy":null,"dataSourceCacheAr":"شبكة المفردات","dataSourceCacheEn":"Arabic WordNet","arabicWordsCache":"| وقت | رابع بُعْد | زمن","englishWordsCache":"fourth dimension | time","superId":1520039870,"superOrder":0,"superTypeCasheAr":"بُعْد","superTypeCasheEn":"dimension","categoryId":null,"area":null,"era":null,"rank":null,"status":null,"subTypesCount":0,"partOfCount":0,"instancesCount":0,"instanceOfID":null,"undiacritizedArabicWordsCache":"| وقت | رابع بعد | زمن","normalizedEnglishWordsCache":"| fourth dimension | time |","exactWord":null}

# Ranking Metrics

We developed three strategies:

❖ **Citation strategy ($R_{cit}$)** frequency of the lexical concept terms:

$$R = \sum_{n=1}^{|A|} \sum_{m=1}^{k} F_{a_{nm}}$$

$$R_{cit} = \frac{R - R_{min}}{R_{max} - R_{min}}$$

❖ **Lexicon renown ranking strategy ($R_{ren}$):** experts assigned each lexicon a rank based on its renown.

❖ **Hybrid ranking strategy ($R_{hyb}$)** is a combination metric:

$$R_{hyb} = R_{ren} + R_{cit}$$

# The 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, WikiData, BFO, DOLCE

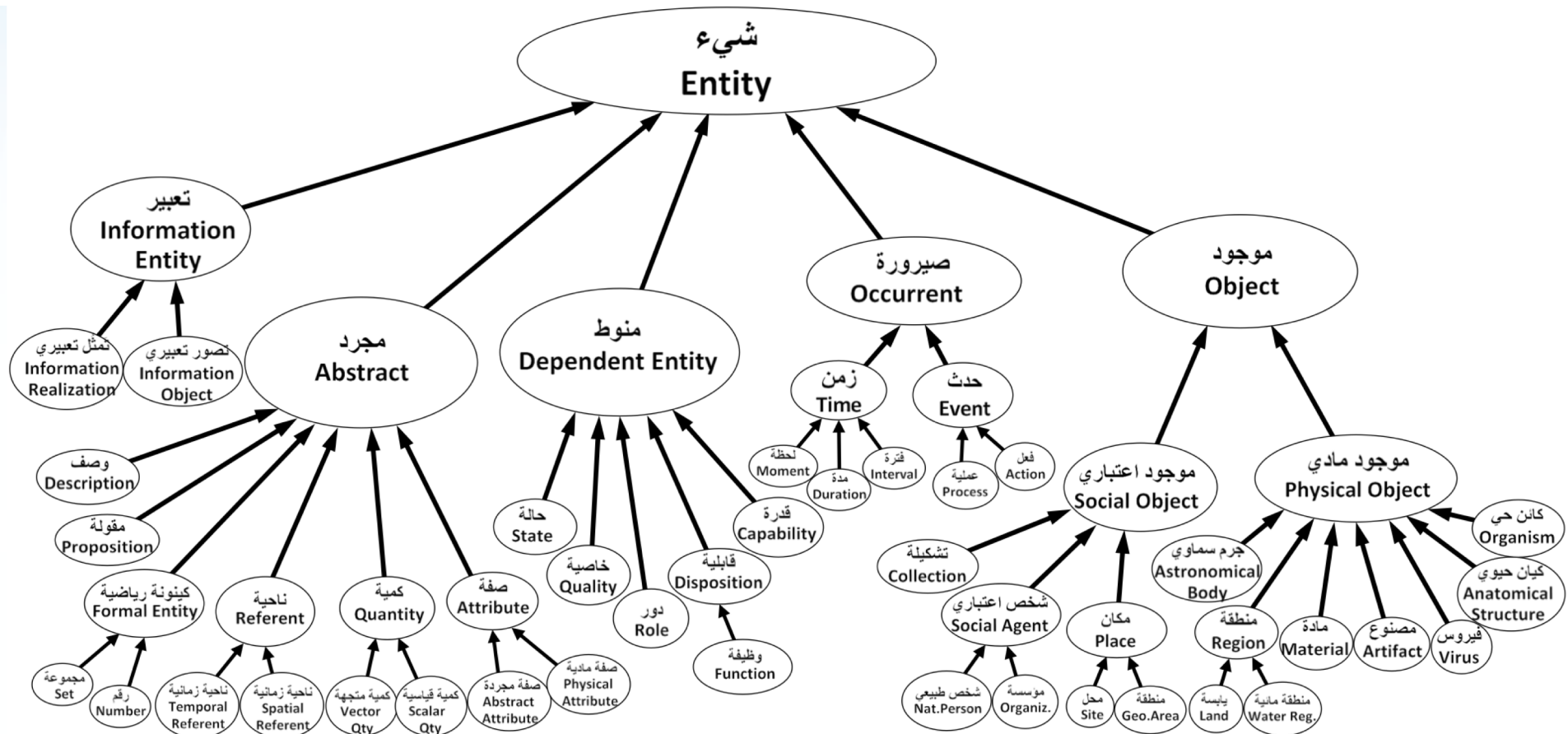- Benchmarked to scientific advances rather than to speakers' naïve beliefs as wordnets do.



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

# Top Levels of the Arabic Ontology



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

# Linking Lexicons with the Arabic Ontology

❖ Lexical concepts (in lexicons) are interlinked with the entities in the ontology.

❖ Given two entities $e_1$ and $e_2$, a *mapping correspondence* between them is defined as the following:

$$< e_1, e_2, R, P, C>$$

❖ Progress so far:

| Relation | Number of Mappings |
|---|---:|
| SameAs | 10500 |
| SubClassOf/SuperClassOf | 950 |
| PartOf/HasPart | 100 |
| InstanceOf/Type | 770 |
| Similar | 80 |
| Total | 12400 |

➢ In this way, lexical concepts across all lexicons would be semantically linked

# Work in Progress

- Lemmatize each lexical entry in every lexicon

- 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, 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). Pages(234--246), LNCS 11608, Springer. 2019.

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

4. Hamzeh Amayreh, Mohammad Dwaikat, and Mustafa Jarrar. Lexicon Digitization -A Framework for Structuring, Normalizing and Cleaning Lexical Entries, Technical Report, Birzeit University 2018.

5. 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

6. 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.

7. 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.

8. 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

9. 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