

نحو بناء قواعد البيانات معجمية للفصحي والعربية

مصطفى جرار
جامعة بيرزيت - فلسطين

الذكاء الاصطناعي

فرصة
الاستثمار باللغة
العربية
صناعة وبحث
علمي مستدام

يمثل قلب
الثورة الصناعية الرابعة

فرصة للدول
الصغيرة
للتنافس
والبروز

100110
1-10-1
0011



A close-up photograph of a woman with long, dark, wavy hair. She is looking upwards and slightly to the left with her eyes closed and her mouth open, as if she is singing or yawning. Her head is tilted back. The background is dark and out of focus. The entire image is enclosed within a white, decorative scalloped border.

A collage of Arabic text snippets on a dark background. The snippets include:

- يرجى ضبط المنبه لheure 05:00
- تشغيل الموسيقى لـ 10:00
- كيف سيكون الطقس اليوم؟
- أرني الموعايد للأسبوع المقبل
- انتقل إلى العنزل زل

The background features a grid pattern with binary code (01010101...) and a partial view of a circuit board.

سوق ضخمة جدا

جهات وأسواق
أخرى
لدواعي اقتصادية
وأمنية

٥,١ مليار
مستخدم للغة
العربية
(كلغة ثانية/دينية)

٣٠٠ مليون
مستخدم في
الوطن العربي
(لغة أولى)

أثر عدم دعم التطبيقات الحديثة للعربية

على التعليم والثقافة والاقتصاد والسياسة وذوي الاحتياجات الخاصة....



❖ ازدادت أهمية المصادر اللغوية (معاجم، مكازن، مسارد، أنطولوجيات..) لبناء التطبيقات الحاسوبية، مثل:

- الترجمة الآلية
- استرجاع البيانات بأكثر من لغة
- فهم وتحليل النصوص المكتوبة والمنطقية
- التحدث مع الآلة

❖ شح في المصادر اللغوية العربية المحosome

تخيلوا أنه لا يوجد قائمة بجميع المدخلات العربية!!!

❖ لا يوجد مصادر لغوية حديثة

مثلاً، لا يوجد مكازن مترادافات تشمل أغلب الكلمات، أو مكازن مصممة بشكل جيد

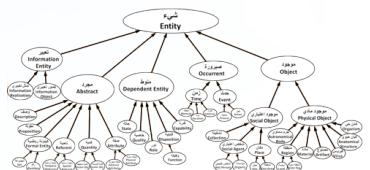
➢ لهذا، يوجد ضعف في دعم التطبيقات للغة العربية

جمع ورقمنة وتنقية ودمج وتوحيد المصادر اللغوية

- واحتها للجمهور العربي والباحثين و المتعلمي اللغة العربية
- واحتها لمطوري التطبيقات عبر واجهات برمجية (APIs)

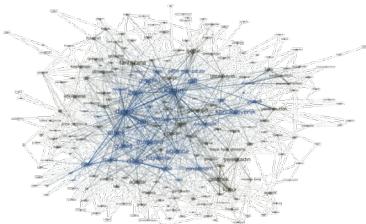
المصادر اللغوية في جامعة بيرزيت

الأنطولوجيا العربية



تصنيف مفاهيم (معاني) الكلمات العربية
Formal Arabic Wordnet
عمل فلسفي حاسوبي لغوی

قاعدة بيانات معجمية



أضخم قاعدة بيانات معجمية
في تاريخ العربية
(تصريفات، اشتتقاقات، دلالة)

مدونة للعامية



مجموعة ضخمة من النصوص
العامية، تم تصريف ووسم كل
كلمة فيها بحوالي 16 سمة.

Linguistic Big Data Graph
شبكة بيانات لغوية ضخمة
تجمع المستويات التصريفية والاشتقاقية والدلالية للفصحي والعامية

<https://ontology.birzeit.edu>



قاعدة بيانات معجمية

Lexicographic Database

محرك بحث معجمي

- البحث في 150 معجماً واسترجاع تعريفات، ومترادفات، وفروق لغوية، وترجمات متخصصة.
- دقيق! مقارنة بالمتجممات الآلية.
- الأول عالمياً، فمثلاً، لا يوجد محرك بحث للإنجليزية، يحتوي هذا الكم من المعاجم في قاعدة بيانات واحدة.



<https://ontology.birzeit.edu>

قاعدة بيانات معجمية

• عربية ومتعددة اللغات

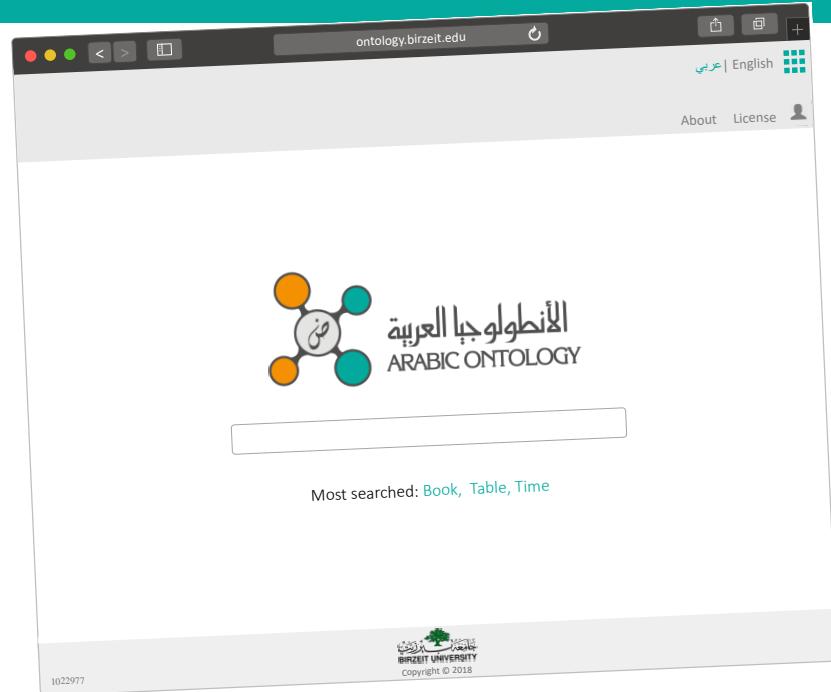
- تحتوي على أنواع كثيرة من المصادر اللغوية:
معاجم لغوية، ومعاجم ثنائية/ثلاثية اللغة،
قواعد بيانات تصريفية، ومسارد، ومكازن،
وفروق لغوية، والأنطولوجيا العربية، وغيرها.
- تغطي معظم المجالات: علوم، وهندسة،
وصحة وطب، وفلسفة، وإنسانيات، وفنون، ...

The screenshot shows a web browser window for ontology.birzeit.edu. The search bar contains the word "attribute". Below the search bar are three checkboxes: "Translations", "Synonyms", and "Definitions", all of which are checked. The main content area is titled "ONTOMY" and displays search results for "attribute". It lists 43 results found in 0.09 seconds. The first result is "Predicate | Attribute" with a definition in Arabic: "المحمول عند المدققين هو المكون به في القضية الجملة دون الترتيبية دون الترتيبية، أما في الترتيبية فيensi تالية، ففي قوله: زرير كريم، زرير هو الموضع، وكريم هو المحمول، والموضع والمحمول عند النقطة المزدوجة". The second result is "Attribute" with a definition in Arabic: "صفة المعرفة عن قيمة خاصية لميزة ما دون استعمال وحدة قياس". The third result is "Physical Attribute" with a definition in Arabic: "صفة مادية". The fourth result is "Abstract Attribute" with a definition in Arabic: "صفة انتقالية". At the bottom of the page, there is a navigation bar with icons for back, forward, and search, and a footer containing the university logo and copyright information.

<https://ontology.birzeit.edu>

محرك بحث معجمي

- متابعة للعامة: طلبة، مترجمين، باحثين، متعلمي اللغة ...
- متاح لمطوري التطبيقات: عبر واجهات برمجية (APIs)



<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

محرك بحث معجمي

تم إنجازه عبر عشرة سنوات

- طباعة يدوية، وإعادة صياغة لبنيّة المعلومات، وتخزينها في قاعدة بيانات محوسبة.
- عمل طوعي، ولم يتم تمويله من أية جهة داعمة.
- عمل غير ربحي.

The screenshot shows a web browser window for ontology.birzeit.edu. The search bar at the top contains the word "attribute". Below the search bar are three tabs: "Ontology", "Dictionaries", and "Morphology", with "Dictionaries" being the active tab. A "Translations", "Synonyms", and "Definitions" checkbox group is also present. The main content area displays search results for "attribute".
The first result is "Predicate | Attribute" with the definition: "المُحَوَّل عَدَ المُتَابِقِينَ فِي الْمُكَوَّبِ بِهِ فِي الْقُصْبَةِ الْجَاهِلِيَّةِ تَوْنُ التَّرْطِيبِ، أَمَّا فِي الشَّرْطَةِ فَيُسَمِّيُ الْمُتَابِقِينَ فِي قَوْلِ زَيْدٍ كَرِيمٍ فِي الْمَوْضِعِ، وَكَرِيمٌ هُوَ الْمُحَمَّلُ، وَالْمُحَمَّلُ وَالْمُحَوَّلُ عَدَ الْمُتَنَقِّبِينَ لِلْمُزَدِّ".
The second result is "Attribute" with the definition: "صَفَّةٌ | خَصْلَةٌ | مُنْتَدٌ | حَاجَةٌ | لَفْظٌ | مُزَدِّةٌ | لَفْظٌ | حَيْنَةٌ | سَمَّةٌ" and the note: "An abstract that describes qualities of an entity without using units of measure".
The third result is "Physical Attribute" with the definition: "صَفَّةٌ مُنْتَدٌ" and the note: "صفة تتعذر عن قيمة خاصية مجزدة ثالثة، ما دون انتصار وحدة قابل".
The fourth result is "Abstract Attribute" with the definition: "صَفَّةٌ مُخَرَّجَةٌ | صَفَّةٌ مُنْتَدَةٌ" and the note: "صفة تعذر عن قيمة خاصية مجزدة ثالثة، ما دون انتصار وحدة قابل".
The footer of the page includes the "Philosophy Lexicon (vol.1&2) ©" logo, "Arabic WordNet ©", "Academy of the Arabic Language-Cairo", "Al Amira Printing House", "Academy Portal", "Browse Online", and the copyright notice "Copyright © 2018 BIRZEIT UNIVERSITY".

<https://ontology.birzeit.edu>

حقوق الملكية

تم التواصل مع أصحاب المعاجم وأخذ إذنهم.

ترويج للمعاجم، وإظهار معلومات المعجم بجانب كل تعريف (المؤلف، الناشر، الصفحة، ...)

The screenshot shows a web browser window for ontology.birzeit.edu. The search bar contains the word "attribute". Below the search bar, there are tabs for Ontology, Dictionaries, and Morphology. The results are displayed under the "Ontology" tab. There are three main sections of results:

- Predicate | Attribute**:
المحمول عند المعنقوبي هو المكون به في القضية الحالية دون الشرطية، بما في التريلفية ليسى ثابلا، ففي قوله زيد كريم، زيد هو الموضع، وكريم هو المحمول. والموضع والمحمول عند المعنقوبي [المزيد...](#)
Philosophy Lexicon (Vol.18.2) ©
- الصلة | Attribute**:
الصلة هي الاسم الذي على ضمن أحوال الأذى، أو الحالة التي يكون عليها الشيء: كالشاد، والوحش، والعلم، والجهل، الخ. والصلة عند المخربين هي التعت، واسم الفاعل، وأسم المفعول، والصلة [المزيد...](#)
Philosophy Lexicon (Vol.18.2) ©
- صلة مميزة | attribute**:
صلة مميزة أو متفردة | افتراضية | خاصية | مميزة | نسخ | وصفت | خصيّة
الصلة مميزة هي صلة تتميز عن قيمة خاصية مميزة شبيه بها
أمثلة على صلة مميزة أو متفردة في الجملة:
 - فكرة زيد كريم هي صلة مميزة لخاصية العطاء عن الإنسان.
 - اللون البرتقالي هو صلة مميزة لخاصية اللون في البرقان.

At the bottom of the page, there is a footer with the text "Propositions attributives" and "محمول ما يقال على موضوعه أنه يتضمنا الحالية". Below this, there is a navigation menu with links to "attribute predicate", "attribute", "flavor", "difference", and "character". A red circle highlights the "attribute" link. The footer also includes the "Philosophy Lexicon" logo and copyright information: "Copyright © 2018 BIRZEIT UNIVERSITY".

<https://ontology.birzeit.edu>

محرك بحث معجمي

LexAPI

RESTful web services to directly retrieve linguistic data from the lexicographic database (ontology, Morph solutions, synonyms, translations, definitions, etc.).

- ## ➤ Request API Token

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 Ontology 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": "نقطة المعلومات العربية", "dataSourceCacheEn": "Arabic WordNet", "arabicWordsCache": "fourth dimension | زن", "englishWordsCache": "fourth dimension | time", "superOrder":0, "superTypeCacheAr": "نقطة", "superTypeCacheEn": "dimension", "categoryID":null, "area":null, "era":null, "rank":null, "status":null, "subtypesCount":0, "instanceOfID":null, "partOfCount":0, "instancesCount":0, "instanceOfID":null, "undiacritizedArabicWordsCache": "نقطة رابع | وقت زن", "normalizedEnglishWordsCache": "fourth dimension | time |", "exactWord":null}
}
```

الأنطولوجيا العربية

Arabic Ontology

(Formal Arabic Wordnet)

الأنطولوجيا العربية

- تستعمل كـ (formal Arabic wordnet)
- عمل حاسوبي وفلسي ولغوی.
- الانطولوجيا العربية: وصف معانی الكلمات اللغة العربية، وعلاقات بين هذه المعانی.
- تصنیف مفاهیم (معانی) الكلمات العربية دلالیاً، وتمثیلها بلغة المنطق ليستطيع الحاسوب فهمها.
- تختلف عن المعاجم كون الحكم على صحة المعلومات فيها يعتمد على ما توصلت إليه العلوم، وليس على ما شاع بين الناس.
- مرتبطة بشبکة المفردات الإنجليزية (وردن特) والويكي بيانات.

The screenshot shows a web interface for an Arabic ontology. At the top, there's a search bar with a magnifying glass icon and three tabs: 'Ontology' (selected), 'Dictionaries', and 'Morphology'. Below the tabs, there are three checkboxes: 'Translations', 'Synonyms', and 'Definitions'. A navigation bar on the right includes links for 'About', 'Contact', and 'Logout'.

Entity | كائن | كيانة | كانن

Whatever existed or will exist, and can be realized or imagined

أيما ظُجد أو سيُوجَد ونستطيع إدراكه أو تخيله

كل شيء على ما يرام

example: [كل شيء على ما يرام](#)

293198

مُؤْخَذٌ | كائن | قائم | حقيقي | واقعٌ | شيء | ذات | قيُّوم

An entity that is wholly and independently present in time, and is realized either for its concrete or social existence

شيء له ذات مستقلة بنفسه، وحاضر كلياً في الزمن، ويدرك ذاته قياساً أو لذاته اعتباراً

يختلف إدراكتنا لأي موجود لاختلاف ما يميز أنواعه من الصفات الجوهرية

example: [إدراكتنا لأي موجود](#)

293200 TypeOf : {Entity}

صَيْرُورَةٌ | حدث | حادث | وقْع | أمر

An entity realized by the time of its happening

شيء يدرك ذاته بزمن حدوثه

لا يمكن فهم أي حدث بشكل منفصل عن الإطار الزمني له

example: [حدث بزمن حدوثه](#)

293202 TypeOf : {Entity}

مُنْوَطٌ | مُنْتَدٌ | مُشْتَقٌ | مُشْرُوطٌ

An entity realized by the time of its happening

شيء يعتمد وجوده على وجود أشياء أخرى

طول المبني منوط بوجود المبني وإلا فلا طول له

example: [طول المبني منوط بوجود المبني وإلا فلا طول له](#)

293201 TypeOf : {Entity}

مُجَدَّدٌ | تجربة | غير مادي | نظري

An entity exists only in mind, cannot be measured or socially realized, and does not have a location

An entity exists only in mind, cannot be measured or socially realized, and does not have a location

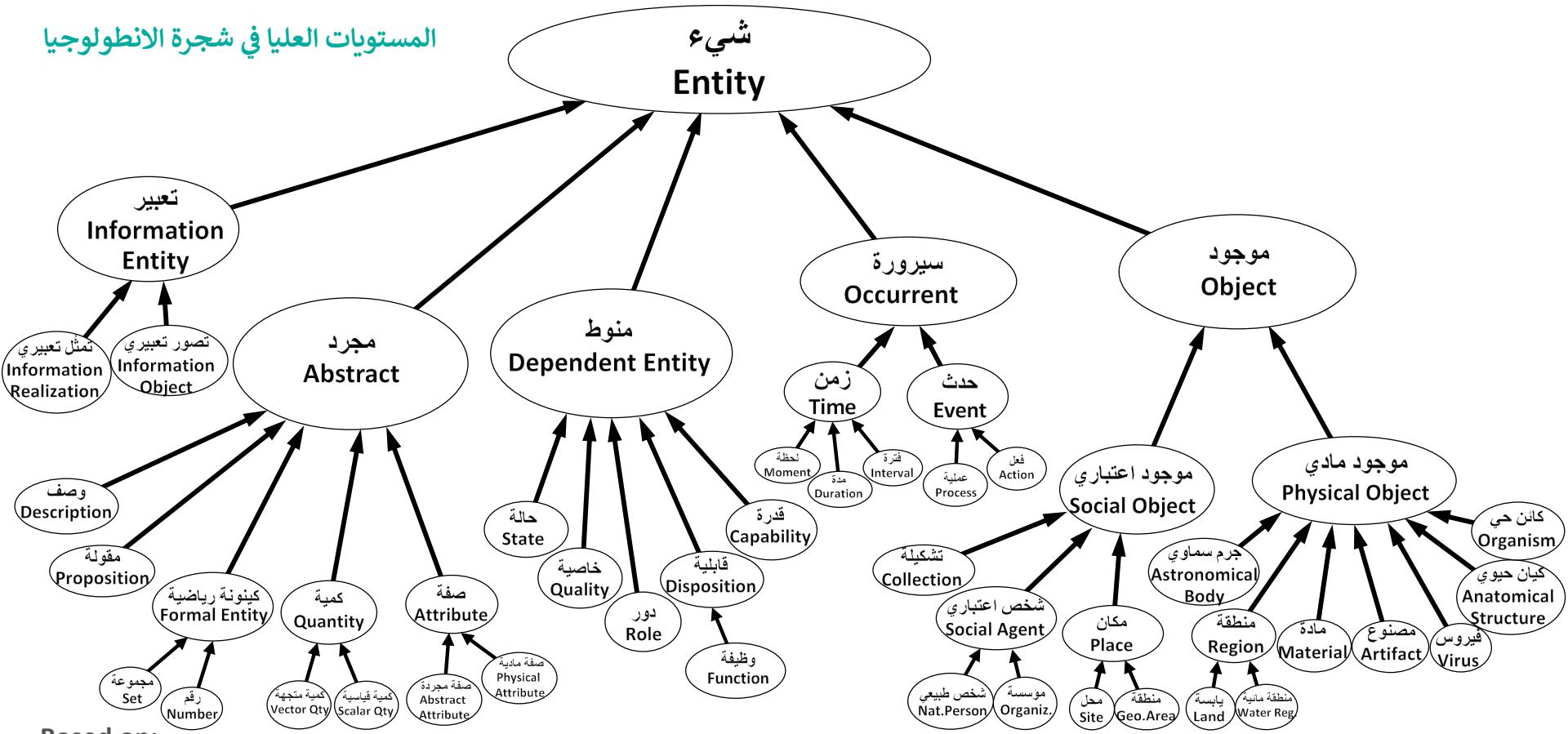
1022977

BIRZEIT UNIVERSITY
Copyright © 2018

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

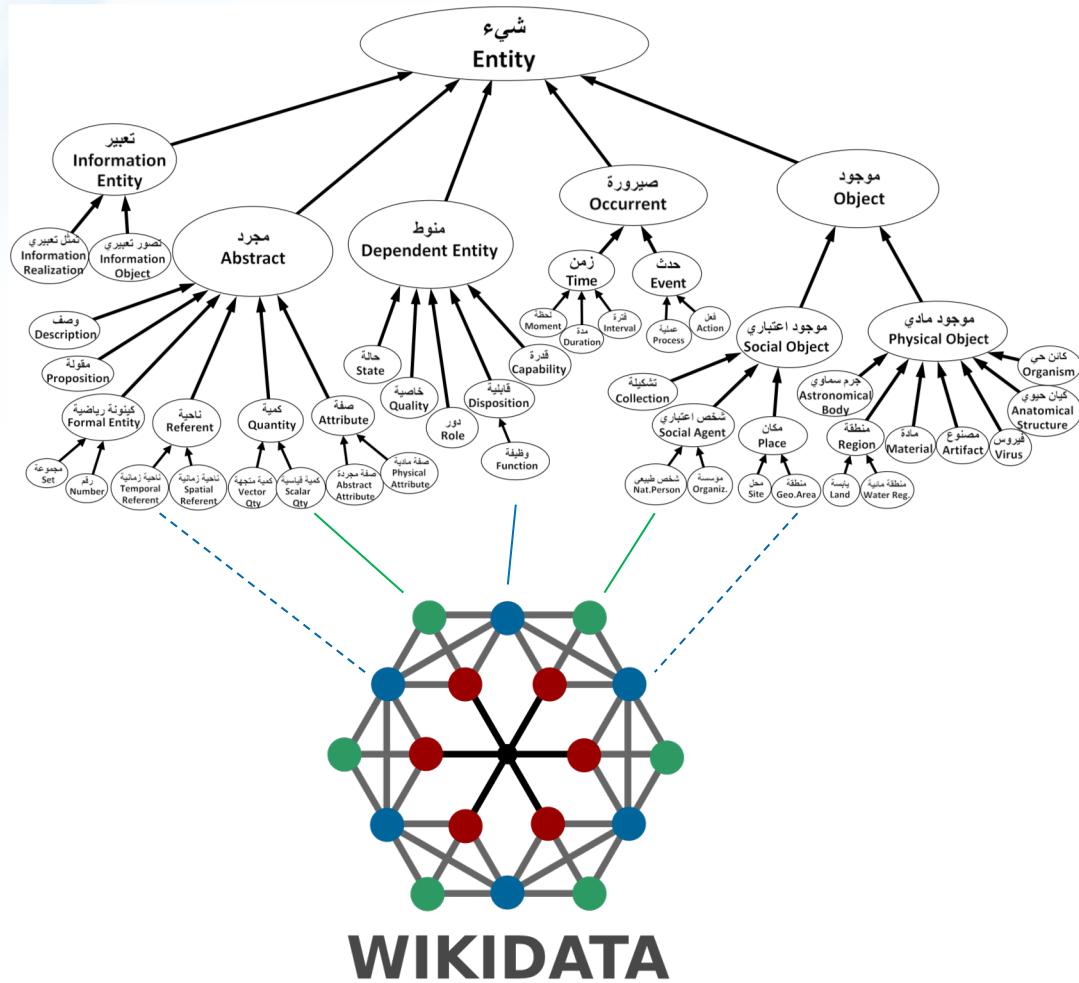
الأنطولوجيا العربية

المستويات العليا في شجرة الأنطولوجيا



Mustafa Jarar: **The Arabic Ontology - An Arabic Wordnet with Ontologically Clean Content.** Applied Ontology Journal, IOS Press.2021.

ربط الأنطولوجيا بالويكيبيديا



- تم ربط كل مفهوم بالأنطولوجيا بمقابله في شبكة بيانات ويكيبيديا (ويكي بيانات).
- وبالتالي، إنشاء شبكة معرفية (Knowledge Graph) مؤصلة بالمفاهيم العربية.
- تطبيقات ذلك مثل، إجابة عن أسئلة معرفية عامة بالعربية.

مدونة اللهجة العامية

Dialectal Corpus

مدونة اللهجة العامية

- عدد ضخم من النصوص المكتوبة باللهجة العامية.
- تصريف كل كلمة بالمدونة ووسمها بحوالي 16 سمة (سوابق، لواحق، جذع الكلمة، فرع الكلمة، المقابل بالفصحي، المقابل بالإنجليزية، ...).
- تحتوي 60 ألف كلمة.



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

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

مدونة اللهجة العامية

تصريف كل كلمة بالمدونة ووسمها بحوالي 16 صفة (سوابق، لواحق، جذع الكلمة، فرع الكلمة، المقابل بالفصحي، المقابل بالإنجليزية، ...).

شو القصة ليش مازمة؟

60k ...

Lemma: قصة

POS: ال/DET + قص/NOUN + ة/NSUFF_FEM_SG

Person: N/A

Gender: female

Gloss: story;stories_

....

Lemma: شو

MSA: ماذا

POS: شو/INTERROG_PRON

Person: N/A

Gender: male

Gloss: what;which

....

مدونة اللهجة العامية

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

Corpus for Palestinian Arabic

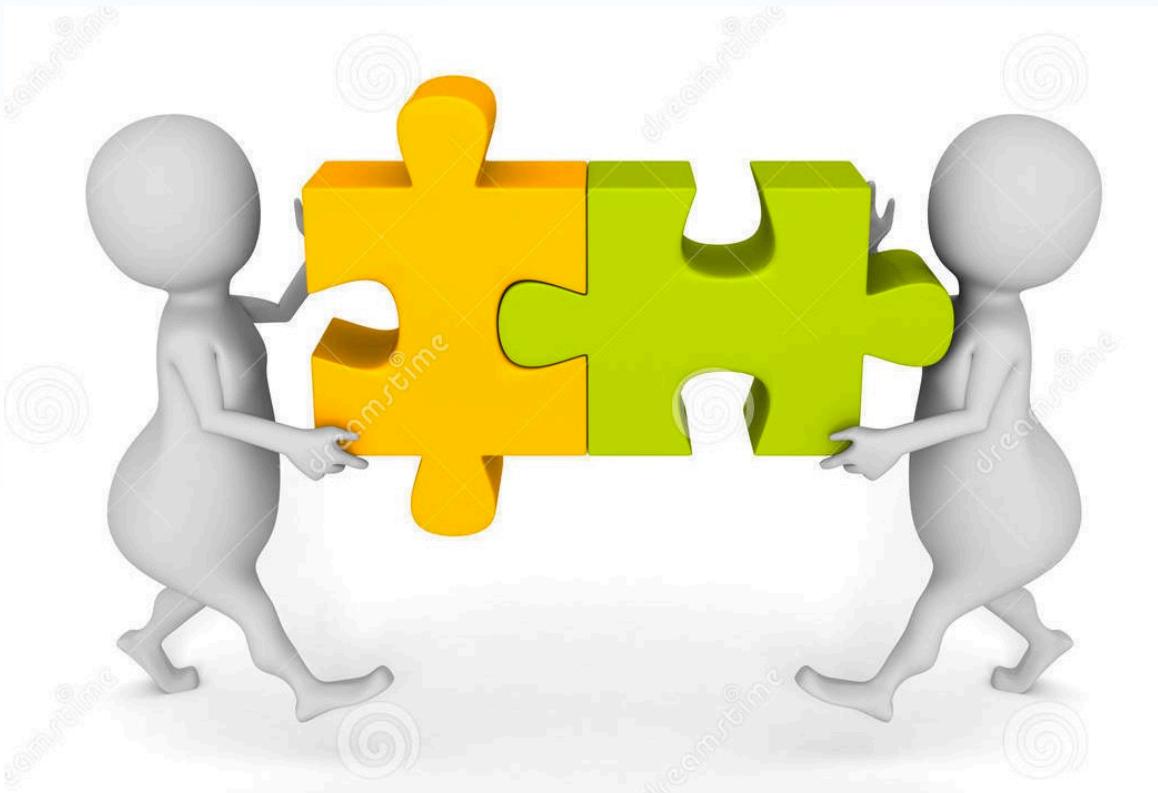
مدونة اللهجة العامية الفلسطينية

بقولك

Word Stem MSA Surface Gloss
 Whole Word Substring

Search

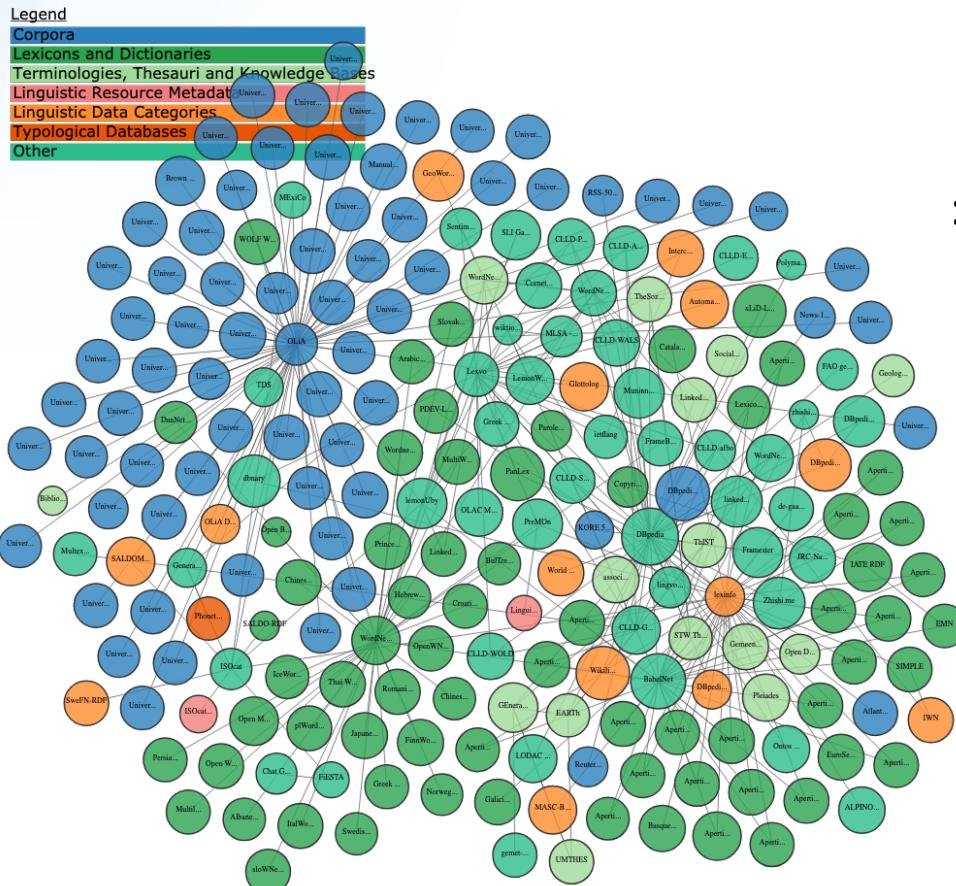
BIRZEIT UNIVERSITY



ربط بين المصادر اللغوية

شبكة عالمية من البيانات المعجمية المترابطة

The Linguistic Linked Open Data Cloud (LLODC)



التوجه العالمي الحديث للتعامل مع المصادر اللغوية:

- جهود بحثية مشتركة لإنشاء شبكة عالمية من البيانات المعجمية المتراقبة.
 - يتم تمثيل كل مدخلة ومفهوم وتعريف... في كل مصدر لغوي بطريقة (lemon) وربطها مع مقابلياتها في المصادر واللغات الأخرى.

قيد البحث والتطوير حالياً

بناء شبكة بيانات لغوية ضخمة للغة العربية تحتوي جميع المستويات التصريفية والاشتقاقية والدلالية والعامية.

Arabic Linked Data Cloud

ربط المعاجم من خلال ربط كل فرع (lemma) وتصريفاتها (inflections) لكل مدخلة في الـ150 معجماً، ثم ربطها جميعاً بشبكات المصادر اللغوية الأجنبية.

Arabic Lemma Index

قائمة بجميع المدخلات العربية القديمة والحديثة والفصحي والعامية، وتحديد سماتها الصرفية ومعانيها وعلاقتها الاشتقاقية.

Example Case Study

Word Sense Disambiguation

Moustafa Al-Hajj, Mustafa Jarrar: [ArabGlossBERT: Fine-Tuning BERT on Context-Gloss Pairs for WSD](#). In Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2021). PP 40--48, 2021

The problem

The Word Sense Disambiguation (WSD) Task

Given a word in a context, which sense (i.e. meaning) this word denotes?

قصيدة من عيون الشعر

Set of senses

1. عُضو الإبصار في الإنسان والحيوان: له عينان كعَيْنَيُ الصقر - ألا إنما العينان للقلب رائد ...
2. جاسوس، "كان عيناً لدولة أجنبية . بِتُّ العيون : تجسس، راقب - فلان عين على فلان : ناظر عليه
3. أجود كل شيء وأحسن ونفيسه: عيون الفن.
4. حارس: فلان عين على المكان.
5. الحاضر من كل شيء أصبح أثراً بعد عين ...
6. عَيْنُ الماء:- ينبع منه، تُحْلِق الطيور فوق عيون الماء
7. عَيْنُ الشَّيْء:- نفسه، ذاته (تستعمل للتوكيد): جاء القوم أعينهم - كُنَّا في المكان عينه.
8. عَيْنُ العقل:- قدرة ذهنية موروثة على التخييل وتذكر الأحداث.
- 9

WSD has been a challenging task for many years but has gained recent attention due to the advances in contextualized word embedding models such as BERT.

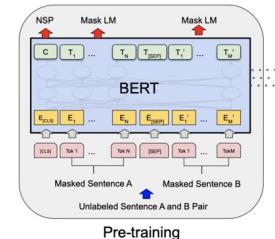
❖ Arabic context-gloss pairs Dataset (167k)

- Extracted from Birzeit University's Lexicographic database
- Annotated target words in context;

Gloss	Context	Label
[CLS] قصيدة من عيون الشعر [SEP] أجدو كل شيء وأحسنه ونفيسيه [SEP]		True
[CLS] قصيدة من عيون الشعر [SEP] عين الشيء : نفسه ، ذاته (تستعمل للتوكيد) [SEP]		False
[CLS] جاء القوم أعينهم [SEP] عين الشيء : نفسه ، ذاته (تستعمل للتوكيد) [SEP]		True
[CLS] جاء القوم أعينهم [SEP] أجدو كل شيء وأحسنه ونفيسيه [SEP]		False

❖ Three Fine-tuned BERT Models

- WSD into **binary sequence-pair classification task**
- **Accuracy 84%**
- 4 types of signals to emphasize target words in context



References

1. Moustafa Al-Hajj, Mustafa Jarrar: [ArabGlossBERT: Fine-Tuning BERT on Context-Gloss Pairs for WSD](#). In Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2021). PP 40–48, 2021
2. Mustafa Jarrar. [**The Arabic Ontology - An Arabic Wordnet with Ontologically Clean Content**](#). Applied Ontology Journal, IOS Press, 2021.
3. Mustafa Jarrar, 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: [**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**](#). Technical Report, Birzeit University. 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

شكرا لاهتمامكم

مصطفى جرار - جامعة بيرزيت

mjarrar@birzeit.edu