

Arabic Language Understanding Tools

Mustafa Jarrar
SinaLab, Birzeit University
Palestine

Artificial Intelligence

Opportunity
to invest in
Language
as industry
and
sustainable
R&D

The 4th Revolution

Opportunity
for
developing
countries to
compete and
impact



We build
tools and
resources
for NLU

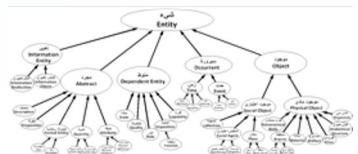
Lexical Resources at SinaLab - Birzeit University

Lexicographic Database



150 lexicons
Largest Arabic lexicographic database

Arabic Ontology/Wordnet



Formal Arabic Wordnet with ontologically clean content

Synonyms 90s%

WSD 84% NER 90s%

Annotated Corpora



Dialects,
NER, WSD, synonyms
Intents, hate
....

Intent 88.4%

NLP library



APIs
Linguistic Data, synonyms, Nested NER, intents, ...

Offensive 88.4%

Big Linguistic Data Graph

<https://ontology.birzeit.edu>



Natural Language Understanding Tools and Datasets



Open Source

<https://sina.birzeit.edu/resources>



Resources

Download and try NLP/NLU datasets, corpora, tools and services



+ Lexicographic Database (150 lexicons)

حوسبة المعاجم

+ Arabic Ontology

الأنطولوجيا العربية

+ Dialect Corpora (Currasat)

كراسات مدونة العاميات

+ Arabic Synonyms

استخراج متزدقات

+ Named Entity Recognition (Wojood)

وجود – استخراج أسماء الاعلام

+ Word Sense Disambiguation (Salma)

سلمى – محلل دلالي

+ ArBanking77 Intent Detection

تحديد المقصود في المساعدات الآلية

+ Offensive Language Detection

خطاب الكراهية بالعبرية

+ Lemmatizer

+ NLP Tools

The Lexicographic Database

- The largest lexicographic Arabic database (150 Lexicons)
 - Contains most lexicon types: glossaries, thesauri, bi/trilingual dictionaries, morph datasets, Arabic Ontology, and more.
 - Covers most domains: science, technology, law, business, art, philosophy, ...
 - Digitized over 10 years.



<https://ontology.birzeit.edu>

Some Statistics

150 Lexicons

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

- **API accessible** for NLP applications.

Natural Language Understanding Tools and Datasets



Open Source

<https://sina.birzeit.edu/resources>



Resources

Download and try NLP/NLU datasets, corpora, tools and services



+ Lexicographic Database (150 lexicons)

حوسبة المعاجم

+ Arabic Ontology

الأنطولوجيا العربية

+ Dialect Corpora (Currasat)

كراسات مدونة العاميات

+ Arabic Synonyms

استخراج متزدقات

+ Named Entity Recognition (Wojood)

وجود - لاستخراج أسماء الاعلام

+ Word Sense Disambiguation (Salma)

سلمى - محلل دلالي

+ ArBanking77 Intent Detection

تحديد المقصود في المساعدات الآلية

+ Offensive Language Detection

خطاب الكراهية بالعبرية

+ Lemmatizer

+ NLP Tools



Arabic Ontology



- Characterization of the intended meaning (i.e., concepts) that the Arabic terms convey.
- Formal **Arabic Wordnet** - with ontologically-clean content.

- Linked with **PWN**, **Wikidata**, **BFO**, **DOLCE**
- Linked with many lexicons

The screenshot shows a search interface with a logo at the top left. Below the search bar are three checkboxes: Translations, Synonyms, and Definitions. The 'Ontology' tab is selected. The main content area displays a list of entity types with their definitions in English and Arabic, examples, and type information.

- Entity | كيان | كائن | شئ | شيء | ذات | ذات فلسفية**
Whatever existed or will exist, and can be realized or imagined
أيما ظُجد أو سيوجد ونستطيع إدراكه أو تخيله
example: كل شيء على ما يرام
293198
- موجود | كائن | قائم | حقيقي | واقع | شئ | ذات**
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 بذاته قياساً أو لذاته اعتباراً يختلف ادراكنا لاي موجود لاختلاف ما يميز أنواعه من المستفات الجوهرية
example: 293200 TypeOf : {Entity}
- متى | متى | متى | متى | متى**
An entity realized by the time of its happening
الشيء الذي يدرك ذاته وأجزاءه بجريانه عبر الزمن
لا يمكن فيه أي حدث بشكل منفصل عن الأطار الزمني له
example: 293202 TypeOf : {Entity}
- متى | متى | متى | متى | متى**
An entity realized by the time of its happening
شيء يعتمد وجوده على وجود أشياء أخرى
طول المبني منوط بوجود المبني والأفلأ طول له
example: 293201 TypeOf : {Entity}
- مجرد | تخييري | غير مادي | نظري**
An entity exists only in mind, cannot be measured or socially realized, and
مجرد | تخييري | غير مادي | نظري

1022977

BIRZEIT UNIVERSITY
Copyright © 2018

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



Arabic Ontology



- Current size so far (but the numbers are dynamic)

1800 fully-done concepts (mostly top levels)

17K partially investigated (ready for NLP applications)

Some branches are elaborated, other not yet.

- English labels are not our target - provided for readability and communication.
- Methodology: Built top-down and bottom-up at the same time.

مترورة | حوث | حصول
occurred
An entity realized by the time of its happening
الشيء الذي تدرك ذاته وأجزائه بجريانه عبر الزمن
لا يمكن فهم أي حدوث بشكل منفصل عن الإطار الزمني له:
example: 293202 TypeOf : {entity}

عملية
process
A cumulative occurred that is composed of a sequence of actions happening respectively in time
حدث تراكمي يتكون من سلسلة من الأفعال المترتبة، التي تحدث بشكل متتابع على خط الزمن
هناك ثلاثة أنواع من الخلايا شاركت في عملية نمو العظم:
example: 293215 TypeOf : {occurred}

عملية بiological process | single-organism process | single organism process | physiological process
biological process
عملية بيولوجية
A biological process represents a specific objective that the organism is genetically programmed to achieve. Biological processes are often described by their outcome or ending state.
See More..
@Collected
455000019 TypeOf : {process} Ontologies

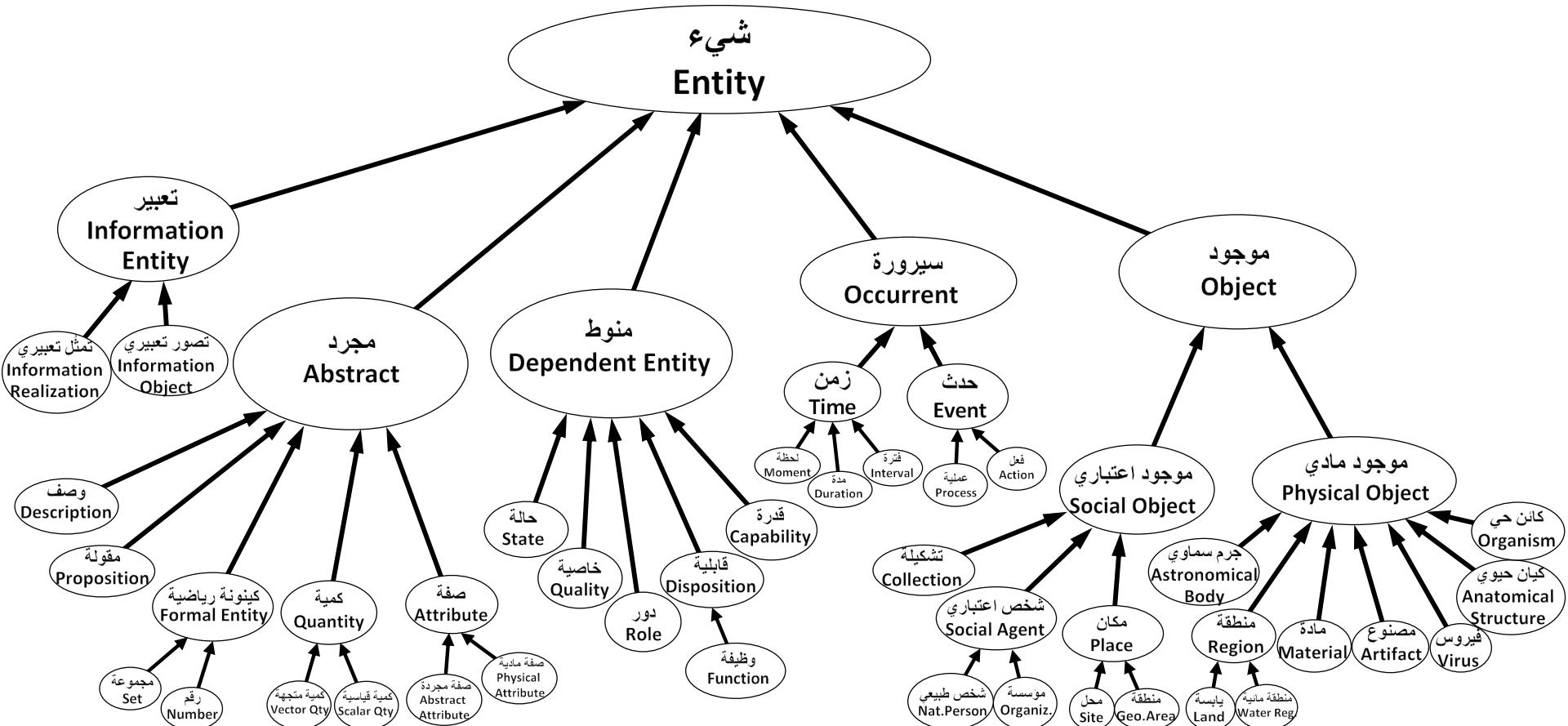
عملية جسمانية
bodily process
عملية في جسم ما
A process in which at least one bodily component of an organism participates. I [OGMS_0000060]
455000574 TypeOf : {biological process} Ontologies

سلوك
behavior
The internally coordinated responses (actions or inactions) of animals (individuals or groups) to internal or external stimuli, via a mechanism that involves nervous system activity. [GO_0007610]
455000020 TypeOf : {bodily process} Ontologies

عملية ادراكية
mental process
A mental process is a bodily process that is of a type



Top Levels of the Arabic Ontology



Natural Language Understanding Tools and Datasets



Open Source

<https://sina.birzeit.edu/resources>



Resources

Download and try NLP/NLU datasets, corpora, tools and services



+ Lexicographic Database (150 lexicons)

حوسبة المعاجم

+ Arabic Ontology

الأنطولوجيا العربية

+ Dialect Corpora (Currasat)

كراسات – مدونة اللهجات العربية

+ Arabic Synonyms

استخراج متزدقات

+ Named Entity Recognition (Wojood)

وجود – لاستخراج أسماء الاعلام

+ Word Sense Disambiguation (Salma)

سلمي – محلل دلالي

+ ArBanking77 Intent Detection

تحديد المقصود في المساعدات الآلية

+ Offensive Language Detection

خطاب الكراهية بالعبرية

+ Lemmatizer

+ NLP Tools

Seven morphologically-annotated Arabic dialect corpora (**1.35 million tokens**)

Curras2 گراس: Palestinian dialect corpus (56K tokens)

Baladi بلدي : Lebanese dialect corpus (10K tokens)

Nabra نبرة : Syrian dialect corpus (60K tokens)

Lisan لسان: Yemeni, Iraqi, Libyan, and Sudanese dialects corpora

Yemeni (1.2 million), Iraqi (46K tokens), Libyan(52K tokens), Sudanese(53K tokens)



- El Haff, K., Jarrar, M., Hammouda, T., Zaraket, F., (2022). **Curras + Baladi: Towards a Levantine Corpus**. In Proceedings of the International Conference on Language Resources and Evaluation (LREC 2022), Marseille, France.
- Amal Nayouf, Tymaa Hammouda, Mustafa Jarrar, Fadi zaraket, Mohamad-Bassam Kurdy: **Nâbra: Syrian Arabic Dialects with Morphological Annotations**. In Proceedings of the 1st Arabic Natural Language Processing Conference (ArabicNLP), Part of the EMNLP 2023. ACL.
- Mustafa Jarrar, Fadi Zaraket, Tymaa Hammouda, Daanish Masood, Martin Waehlisch: **Lisan: Yemeni, Iraqi, Libyan, and Sudanese Arabic Dialect Corpora with Morphological Annotations**. The 20th ACS/IEEE International Conference on Computer Systems and Applications (AICCSA). Pages(-). IEEE. Egypt. 2023

Try and download



هيك

Search

EN | ع

Word Stem Lemma Gloss

Whole Word Substring

Palestinian Lebanese Iraqi Libyan Sudanese Yemeni

About Publications

313 results (4.6 secs)

Gloss	Lemma	POS	Suffix	stem	Prefix	Word	Context
this way, this, thus	هيك / هكذا	ضمير اشارة		هيك		و/عطف	اللي اكلناه ، عزمنا الاسبوع الماضي ، عشان نردهم العزيمة ، وهيك بذلك تحطيلهم ترترسي
this way, this, thus	هيك / هكذا	ضمير اشارة		هيك		و/عطف	وهي قاطعا بحفلة قمة ، صارت عجيبة خلت الفحات يطول حشيشن وهيك قدرت تخبي بيناتن ، ومن هون مصربنا نعمل الفححة بيايامنا .
this way, this, thus	هيك / هكذا	ضمير اشارة		هيك		و/عطف	بيغنو : " بلاطة فوق بلاطة صاحبة البيت ضراطة ! " وهيك منكون وصلنا لآخر جولتنا بعالم عبد البربارا .
inon+like_that;thus	هيك	تعجب		هيك		ب/حرف جر	يمشون وياه الناس ولكن بيهم شي اشنون اندر هم يمتحنون بهيج مكان ياوزير ياعار
inon+like_that;thus	هيك	تعجب		هيج		ب/حرف جر	يجوز ياخذ بريد يخرع بيج جهاله
this way, like this, this, thus	هيك	ضمير اشارة		هيك		هيك	وين يعرف يسولف ويدبر هذا الحجي اكيد واحد كليله احجي هيج
inon+like_that;thus	هيك	تعجب		هيج		ب/حرف جر	و ش حيرتو العالم ولعيتو بيه طوبه كل بيج اينزل كلام اليعجبه ويشر



Rights Reserved © 2022

In cooperation with:



Rights Reserved © 2022

<https://sina.birzeit.edu/currasat/>

Try and download



write

Search

EN ع

Word Stem Lemma Gloss
 Whole Word Substring
 Palestinian Lebanese Iraqi Libyan Sudanese Yemeni

About Publications

313 results (4.6 secs)

Gloss	Lemma	POS	Suffix	Stem	Prefix	Word	Context
be written/be fated/be destined write	كتب	فعل مضارع مفرد، منذكر، متكلّم		كتب	و/عطف+المضارع المتكلّم المفرد	واكتب	وراح تكون هية واكتب وصيّة
be written/be fated/be destined write	كتب	فعل مضارع مفرد، منذكر، مخاطب		كتب	و/عطف+إادة مضارع+ات+المضارع المخاطب المنذكر المفرد	وبنكتب	ونتعد على ليس وبنكتب
be written/be fated/be destined write	كتب	فعل مضارع مفرد، منذكر، مخاطب	ت+للماضي: فاعله مخاطب منكر مفرد	كتب	و/عطف	وكتب	وكتب بوصيتك هو السبب
be written/be fated/be destined write	كتب	فعل مضارع مفرد، منذكر، غائب	ل/حرف جر+ها/اضمير متصل للغائب	كتب	و/عطف+ي/المضارع الغائب المنذكر الغير	ويكتبها	بس ثوفت انو قبل كم يوم كان يحطّلها أغاني ويحكى عنها ويكتبها ...
be written/be fated/be destined write	كتب	فعل مضارع مفرد، منذكر، غائب		كتب	ي/المضارع الغائب المنذكر المفرد	يكتب	الاين : (يكتب ساتوس ع ليس) لا تتحدثوا عن الواقع اكتر الاشياء وجما
be written/be fated/be destined write	كتب	فعل مضارع مجهول مفرد، مؤنث، غائب		كتب	ت/المضارع الغائب المنذكر المفرد	تنكتب	، وفجيئت كيف هالله تظورت من أتحق لوجه ما كان ينسوا تكتب للغة عالميه ما ينسوا ما تكتب .
be written/be fated/be destined write	كتب	فعل مضارع مجهول مفرد، مؤنث، غائب		كتب	ت/المضارع الغائب المنذكر المفرد	تنكتب	أنجح لهجه ما كان ينسوا تكتب للغه عالميه ما ينسوا ما تكتب .
scatter/sprinkle/write in prose	تشر	فعل مضارع مفرد، منذكر، غائب		تشر	ب/إادة مضارعة	بنثر	وبنثر رماد
they (people) + write + اسماء معرفة	كتبوا	فعل مضارع	ـوا/المضارع: فاعله منذكر جمع	كتبوا	ي/المضارع الغائب المنذكر المفرد	يكتبوا	و الكل يتكلّم عن السنن الجرف الأجنبية وكل الصيادين يبغوا يكتبوا



Rights Reserved © 2022

In cooperation with:

 AMERICAN UNIVERSITY OF BEIRUT United Nations
 Rights Reserved © 2022

<https://sina.birzeit.edu/currasat/>

Annotation Tools

SinaLab

Tawseem Portal

بوابة سينا لتوسيم المدونات

Stat Log

Lemmas Annotations

Gloss	MSA Lemma	DA Lemma	Person	Gender	Number	POS	Suffix	Stem	Prefix	Token	Context	Dialect
how	كيف 1	شلون				أداة استفهام		شلون		شلون	شلون دا تدخلني تسلمي	حلبية
will/shall	سُوفَ 1	رَجَّا				أداة استقبال		دا		دا	شلون دا تدخلني تسلمي	حلبية
enter	دَخَلٌ 1	دَخَلٌ	مفرد	مؤنث	مفرد	فعل مضارع	ي/لل مضارع: فاعله مخاطب مؤنث مفرد	دخل	ت/لل مضارع: المخاطب المؤنث المفرد	دخلني	شلون دا تدخلني تسلمي	حلبية
hand over/surren.	سلم 1	سلم	مفرد	مؤنث	مفرد	فعل مضارع	ي/لل مضارع: فاعله مخاطب مؤنث مفرد	سلم	ت/لل مضارع: المخاطب المؤنث المفرد	سلمي	شلون دا تدخلني تسلمي	حلبية
in childbed	نَفْسَاء 1	نَفْسَاء	مفرد	مؤنث		اسم		نفسا	ع/حرف جر+ال/اداة تعريف	نفسا	شلون دا تدخلني تسلمي	حلبية
..				علامة ترقيم		شلون دا تدخلني تسلمي	حلبية
not	ما 2	ما				أداة نفي		ما		ما	شلون دا تدخلني تسلمي	حلبية
become/begin to	صار 2	صار	غائب	ذكر	مفرد	فعل مضارع		صير	ب/أداة مضارعة	صير	شلون دا تدخلني تسلمي	حلبية
to/for + Allah/God +	الله	الله	ـ	ـ	ـ	اسم		الله	ي/أداة نداء	يالله	شلون دا تدخلني تسلمي	حلبية
sisters sister/count	أخت 1	أخت	-	مؤنث		اسم		خيتو		خيتو	شلون دا تدخلني تسلمي	حلبية

Natural Language Understanding Tools and Datasets



Open Source

<https://sina.birzeit.edu/resources>

Resources

Download and try NLP/NLU datasets, corpora, tools and services



+ Lexicographic Database (150 lexicons)

حوسبة المعاجم

+ Arabic Ontology

الأنطولوجيا العربية

+ Dialect Corpora (Currasat)

كراسات مدونة العاميات

+ Arabic Synonyms

استخراج متزدقات

+ Named Entity Recognition (Wojood)

وجود - لاستخراج أسماء الاعلام

+ Word Sense Disambiguation (Salma)

سلمى - محلل دلالي

+ ArBanking77 Intent Detection

تحديد المقصود في المساعدات الآلية

+ Offensive Language Detection

خطاب الكراهية بالعبرية

+ Lemmatizer

+ NLP Tools

Extracting Synonyms from lexicographic graphs

A Benchmark and Scoring Algorithm for Enriching Arabic Synonyms

Sana Ghanem¹, Mustafa Jarrar¹, Radi Jarrar¹, Ibrahim Bounhas^{2,3}

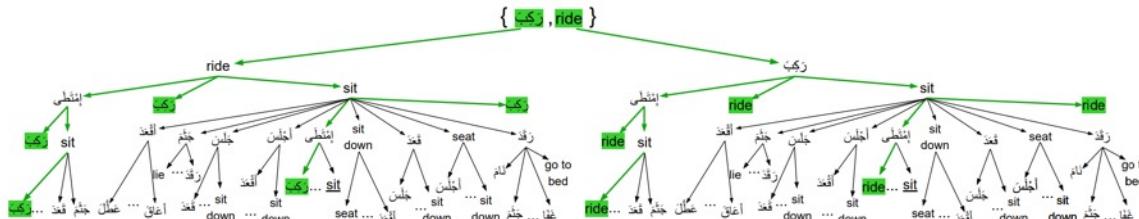
¹Department of Computer Science, Birzeit University, Palestine

²LISI Laboratory of Computer Science for Industrial System, INSAT, Carthage University, Tunisia

³JARIR: Joint group for Artificial Reasoning and Information Retrieval, Tunisia

{swghanem, mjarrar, rjarrar}@birzeit.edu

ibrahim.bounhas@isd.uma.tn



A Benchmark and Scoring Algorithm for Enriching Arabic Synonyms

Dataset Construction

- Four Linguists (500 synsets, 3K synonyms in total).
- Each candidate synonyms was annotated with a **fuzzy value** by four different linguists.

مُحَالَّة		▼	نفس الدلالة، الأسلوب ضعيف ، غير شائعة	60		
إِتْبَالُ		▼	نفس الدلالة، الأسلوب صحيح ، شائعة الى حد قليل	80		
إِتْخَاد		▼	نفس الدلالة والاسلوب والشيوخ	100		
جَامِعَةٌ		▼	نفس الدلالة، الأسلوب ضعيف ، غير شائعة	60		

Score	Meaning
100	Same semantics, style, use
90	Same semantics, style, less used
80	Same semantics, style, rarely used
70	Same semantics, style, not used
60	Close semantics, weak style, uncommon
50	Close semantics, not exact purpose
40	Semantically related
30	Semantically related (somehow)
20	Semantically different
10	Semantically very different
0	Semantically unrelated

Demo





Home > Resources > Synonyms Generator

Synonyms Generator

A dataset and source code for Synonyms Generator

Version: 1.0 (updated on 15/12/2022)

An algorithm to extend a set of synonyms with more synonyms. Given a set of synonyms, the algorithm builds a graph (using many dictionaries) and returns a set of candidate synonyms, each with a fuzzy value to indicate how much it is likely to be a synonym. The more synonyms in the input, the more accurate the candidate synonyms. We trained the fuzzy model using a dataset we built manually (500 synsets, with 3K candidate synonyms by four linguists). Please read this article for the details. Try the service (type synonym separated by | or , or ;):

street | road | شارع | طريق

Extend Evaluate

61 شارع , 62 طريق , road 41% , street 41%

+ Settings

+ Dataset And Downloads

+ API

- Publications

Sana Ghanem, Mustafa Jarrar, Radi Jarrar, Ibrahim Bounhas: [A Benchmark and Scoring Algorithm for Enriching Arabic Synonyms](#). In Proceedings of the Global WordNet Conference (gwc2023), Donostia, January, 2023

[PDF](#) - [Slides](#) - [Video](#)

<https://sina.birzeit.edu/synonyms/>

Download: [Github.Synonyms](#)

A Benchmark and Scoring Algorithm for Enriching Arabic Synonyms

Extend a synset with more synonyms above a given fuzzy value



route 31% , highway 31% , way 15% , via 15% , track 15% ,
thoroughfare 15% , roadway 15% , road way 15% , ride
15% , path 15% , boulevard 15% , avenue 15% , Roadway
15% , Lane 15% , Boulevard 15% , Avenue 15%

سَيِّلٌ %46 , سَيِّلٌ %31 , مَمْرَأٌ %46 , مَمْرَأٌ %31 , مَسْلَكٌ %31 , مَسْلَكٌ %31 , مَوْرٌ %31 , طَرَيْقٌ %31 , طَرَيْقٌ %31 , صِرَاطٌ %31 , صِرَاطٌ %31 , شَارِعٌ %31 , شَارِعٌ %31 , سَنَنٌ %31 , سَنَنٌ %31 , سَيِّلٌ %31 , سَيِّلٌ %31 , رُفَاقٌ %31 , رُفَاقٌ %31 , دَرْبٌ %31 , دَرْبٌ %31 , وَجْهٌ %15 , تَهْجٌ %15 , تَهْجٌ %15 , تَهْجٌ %15

A Benchmark and Scoring Algorithm for Enriching Arabic Synonyms

Evaluate/rank synonyms in a given synset (wordnets, BERT's output, etc.)

street | road | شارع | طريق

Evaluate

61 طريق %62 شارع , road 41% , street 41%

Natural Language Understanding Tools and Datasets



Open Source

<https://sina.birzeit.edu/resources>

Resources

Download and try NLP/NLU datasets, corpora, tools and services



+ Lexicographic Database (150 lexicons)

حوسبة المعاجم

+ Arabic Ontology

الأنطولوجيا العربية

+ Dialect Corpora (Currasat)

كراسات مدونة العاميات

+ Arabic Synonyms

استخراج متزدقات

+ Named Entity Recognition (Wojood)

وجود - لاستخراج أسماء الاعلام

+ Word Sense Disambiguation (Salma)

سلمى - محلل دلالي

+ ArBanking77 Intent Detection

تحديد المقصود في المساعدات الآلية

+ Offensive Language Detection

خطاب الكراهية بالعبرية

+ Lemmatizer

+ NLP Tools

Wojood

Nested Arabic Named Entity Corpus

Wojood

Nested Arabic Named Entity Corpus

Flat NER

Sami	works	at	Jimmy	Carter	Center	✓ mature
B-PERS	O	O	B-ORG	I-ORG	I-ORG	

Nested NER

Sami	works	at	Jimmy	Carter	Center	!
B-PERS	O	O	B-PERS	I-PERS	I-ORG	

Jarrar, M., Khalilia, M., Ghanem, S. (2022). Wojood: Nested Arabic Named Entity Corpus and Recognition using BERT. In Proceedings of the International Conference on Language Resources and Evaluation (LREC 2022), Marseille, France.

Wojood

❖ NER Corpus

Corpus	Nested?	size (tokens)	No. of entities	Entity classes	Which Arabic	Domain
Ontonotes5	No	300k	28k	18	MSA	News
ANERCorp	No	150k	11k	4	MSA	News
Canercorpus	No	258k	72k	14	Classic	Religion
AQMAR	No	74k	-	open	MSA	4 domains
Wojood Corpus	YES	550K	75K	21	MSA & Dialect	Multi domains Media, History, Culture, Health, Finance, ICT, Law, Elections, Politics, Migration, Terrorism

❖ NER BERT

- Multi-task learning (nested entities)
- 88.4% F1-score

Annotation Guidelines

21 entity classes

PERS

People names

فيروز، عادل إمام، ابن احمد، الملك عبدالله، النبي محمد

NORP

Group of people

العرب، المسيحيين، سكان القدس، وزراء الخارجية العرب

OCC

Occupation or professional title

رئيس جامعة بيرزيت، مدير بنك فلسطين، قائد الجيش

ORG

Legal/social body

بنك القاهرة، ريال مدريد، داعش، الجيش المصري،

GPE

Geopolitical: country, city, state

ليبيا، مدينة القدس، الجمهورية اللبنانية، روسيا الاتحادية

LOC

Geographical location (non-GPE)

البحر الميت، قنات السويس، آسيا، الوطن العربي

FAC

Places: landmark, road, building..

مطار صنعاء، سجن ابو غريب، المسجد الأقصى

EVENT

Events of general interest

حرب 1973، القمة العربية 2005، عيد الفطر ، يوم الأرض

DATE

Specific/relative date (>day)

13 يونيو، 2019-2020، الفترة العثمانية

TIME

Specific/relative time (<day)

الساعة ١٢، من الساعة ٥ حتى ٧ مساء، خلال ساعتين

LANGUAGE

Human language or dialect

اللغة العربية، الفصحي، الدارجة المغربية، اللغة الفرنسية

WEBSITE

Website or specific URL

موقع فيسبوك، يوتوب، schema.org

LAW

Geographical location (non-GPE)

قانون الاستثمار ، المادة 114 من قانون العقوبات 2005

PRODUCT

Vehicle, weapon, food, ...

مرسيديس سي.إيه ١٨، آيفون ١٣ ، دبابات مركابا، تروفين

CARDINAL

Numerals in digits/words

١.٥ ، ٣٠ ، ١٥٠ صفر، اثنان ، أربعة وعشرون ، مليون

ORDINAL

does not refer to a quantity

٣ كيلومتر، مئة قدم ، ٣ طن ، ٥٠ غرام ، ٢٥ سم مربع

PERCENT

Word/symbol refers to a percent

٥ بالمائة ، ٩ من كل الف

QUANTITY

Value measured by units

٣ كيلومتر، مئة قدم ، ٣ طن ، ٥٠ غرام ، ٢٥ سم مربع

UNIT

Name/symbol of a unit

ميل ، كيلو ، كيلومتر،إنش ، كيلوغرام ، هكتار ، مل

MONEY

Monetary quantity, incl.

مئة وخمسون درهم اماراتي ، اثنان وثلاثون يورو ، ٨ دولار

CURR

Name/symbol of currency

دولار، جنية مصرى ، دينار ، فرنك ، ريال سعودي ،

Corpus Collection

Source - Topics	Sentences	Tokens
Web Articles ¹ (MSA) Health, Finance, ICT, Law, Elections, Politics, Migration and Terrorism	9,053	258,102
Archive ² (MSA) History and Culture	12,271	227,020
Social Media ³ (Dialect) General topics	5,653	65,342
Total	26,977	550,464

¹ un.org, hrw.org, msf.org, who.org, mipa.institute, elections.ps, sa.usembassy.gov, diplo- matie.ma, quora.com

² Awraq, Birzeit University Digital Palestinian Archive

³ Palestinian and Lebanese dialect corpora

Downloads and Demo

Birzeit University, in cooperation with the Edward Said Foundation, is organizing a folklore festival. The festival will start at 4:00 pm, on May 16, 2016, with the sponsorship of Bank of Palestine for an amount of five thousand dollars.

جامعة بيرزيت وبالتعاون مع مؤسسة ادوارد سعيد تنظم مهرجان للفن الشعبي، سيبدأ المهرجان الساعة الرابعة عصراً، بتاريخ 16/5/2016، وذلك برعاية من بنك فلسطين بمبلغ خمسة الاف دولار.

output format :

جامعة بيرزيت GPE وبالتعاون مع ORG مؤسسة PERS ادوارد سعيد ORG تنظم EVENT مهرجان للفن الشعبي،

سيبدأ المهرجان الساعة الرابعة عصراً، TIME بتاريخ DATE 16/5/2022، وذلك برعاية من بنك GPE فلسطين

MONEY بمبلغ CURRENCY خمسة الاف ORG دولار



<https://ontology.birzeit.edu/wojood>



ArabicNER

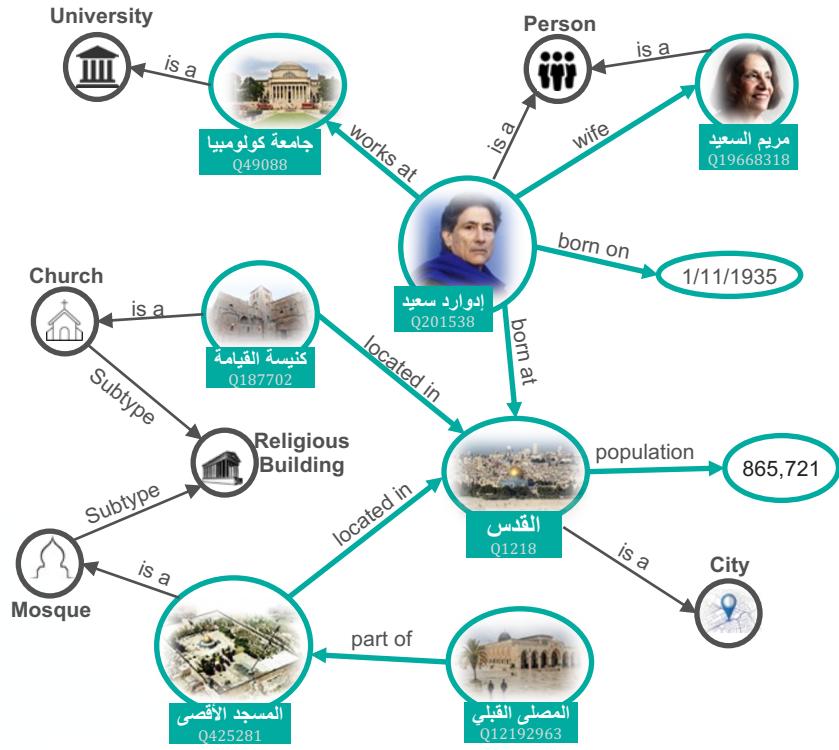


ArabicNER-Wojoood

Text to Intelligence

Text

Knowledge Graph



What Is Named Entity Recognition (NER) Used For?

Customer Support

Recognize and route per products, models, places, phones,

Content Recommendation

understand user profiles and history

Resumes

Recognize name, address, phone, etc.

Chatbots and Question Answering

Recognize names, amounts, places, etc.

Machine Translation

Transliterate entities....

Event Detection

Identify event as they occur (Finance/business events, social media,)

Knowledge Management

knowledge graphs, Digital libraries, content filtering

Classifying content for news providers

Efficient Search Algorithms

Finance

extract figures, loans, credit risk.

News providers

Natural Language Understanding Tools and Datasets



Open Source

<https://sina.birzeit.edu/resources>



Resources

Download and try NLP/NLU datasets, corpora, tools and services

+ Lexicographic Database (150 lexicons)

حوسبة المعاجم

+ Arabic Ontology

الأنطولوجيا العربية

+ Dialect Corpora (Currasat)

كراسات مدونة العاميات

+ Arabic Synonyms

استخراج متزدقات

+ Named Entity Recognition (Wojood)

وجود – استخراج أسماء الاعلام

+ Word Sense Disambiguation (Salma)

سلمى – محلل دلالي

+ ArBanking77 Intent Detection

تحديد المقصود في المساعدات الآلية

+ Offensive Language Detection

خطاب الكراهية بالعبرية

+ Lemmatizer

+ NLP Tools

ArabGlossBERT: Fine-Tuning BERT on Context-Gloss Pairs for WSD

Moustafa Al-Hajj

Lebanese University
Beirut, Lebanon

moustafa.alhajj@ul.edu.lb

Mustafa Jarrar*

Birzeit University
Birzeit, Palestine

mjarrar@birzeit.edu

Context-Gloss Augmentation for Improving Arabic Target Sense Verification

Sanad Malaysha, Mustafa Jarrar, Mohammed Khalilia

Birzeit University, Palestine

SALMA: Arabic Sense-Annotated Corpus and WSD Benchmarks

Mustafa Jarrar, Sanad Malaysha, Tymaa Hammouda, Mohammed Khalilia

Birzeit University, Palestine

{mjarrar, smalaysha, thammouda, mkhalilia}@birzeit.edu

SALMA

ArabGlossBERT for Arabic Word-Sense Disambiguation

The Word Sense Disambiguation (WSD) Task

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

You have beautiful eyes

Set of senses/meanings

1. the organ of sight
2. good discernment (either visually or as if visually)
3. attention to what is seen
4. an area that is approximately central within some larger region)
5. a small hole or loop



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 الإِبصَار فِي الإِنْسَان وَالحَيْوَان: لَهُ عَيْنَانِ كَعْيَيْنِ الصَّفَر - أَلَا إِنَّمَا الْعَيْنَانُ لِلْقَابِ رَائِدٌ ...
2. جَاسُوس، "كَانَ عَيْنًا لَدُولَةً أَجْنبِيَّةً . بِثُّ الْعَيْنَونَ: تَجَسَّسَ، رَاقِبٌ - فَلَانُ عَيْنٌ عَلَى فَلَانَ: نَاظِرٌ عَلَيْهِ
3. أَجْودُ كُلِّ شَيْءٍ وَأَحْسَنُهُ وَنَفِيسُهُ: عَيْنُ الْفَنِ.
4. حَارِسٌ: فَلَانُ عَيْنٌ عَلَى الْمَكَانِ.
5. الْحَاضِرُ مِنْ كُلِّ شَيْءٍ أَصْبَحَ أَثْرًا بَعْدَ عَيْنٍ ...
6. عَيْنُ الْمَاءِ: يَنْبُوعُهُ، تُحَلِّقُ الطَّيْوَرُ فَوقَ عَيْنَ الْمَاءِ
7. عَيْنُ الشَّيْءِ: نَفْسُهُ، ذَاتُهُ (تَسْتَعْمِلُ لِلتَّوْكِيدِ): جَاءَ الْقَوْمُ أَعْيُنَهُمْ - كَتَّا فِي الْمَكَانِ عَيْنَهُ.
8. عَيْنُ الْعُقْلِ: قَدْرَةُ ذَهْنِيَّةٍ مَوْرُوثَةٍ عَلَى التَّخْيِيلِ وَتَذَكُّرِ الْأَحْدَاثِ.
- 9

SALMA

ArabGlossBERT for Arabic Word-Sense Disambiguation

❖ 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]	أجود 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]	أجود 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

- ❖ Fine-tuned three Arabic pre-trained BERT models
The WSD task is converted into **binary sequence-pair classification task**

[CLS] قصيدة من عيون الشعر [SEP]	أجود كل شيء وأحسنـه ونفيـه	[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]	أجود كل شيء وأحسنـه ونفيـه	[SEP]	False

❖ Results

Model	Precision	True	False	Accuracy
AraBERTv02	Precision	81	85	84
	Recall	66	93	
	F1-score	72	89	
CAMELBERT	Precision	77	83	82
	Recall	60	92	
	F1-score	67	87	
QARiB	Precision	73	82	80
	Recall	58	90	
	F1-score	65	86	

Constructing a dataset of context-gloss pairs

Statistics:

	count
Unique Lemmas (undiacritized)	26169
Avg glosses per Lemmas	1.25
Unique Glosses	32839
Unique Contexts	60272
Avg context per gloss	1.83
True context-gloss pairs	60323
False context-gloss pairs	106884
Total True and False pairs	167207

Training and Test Datasets

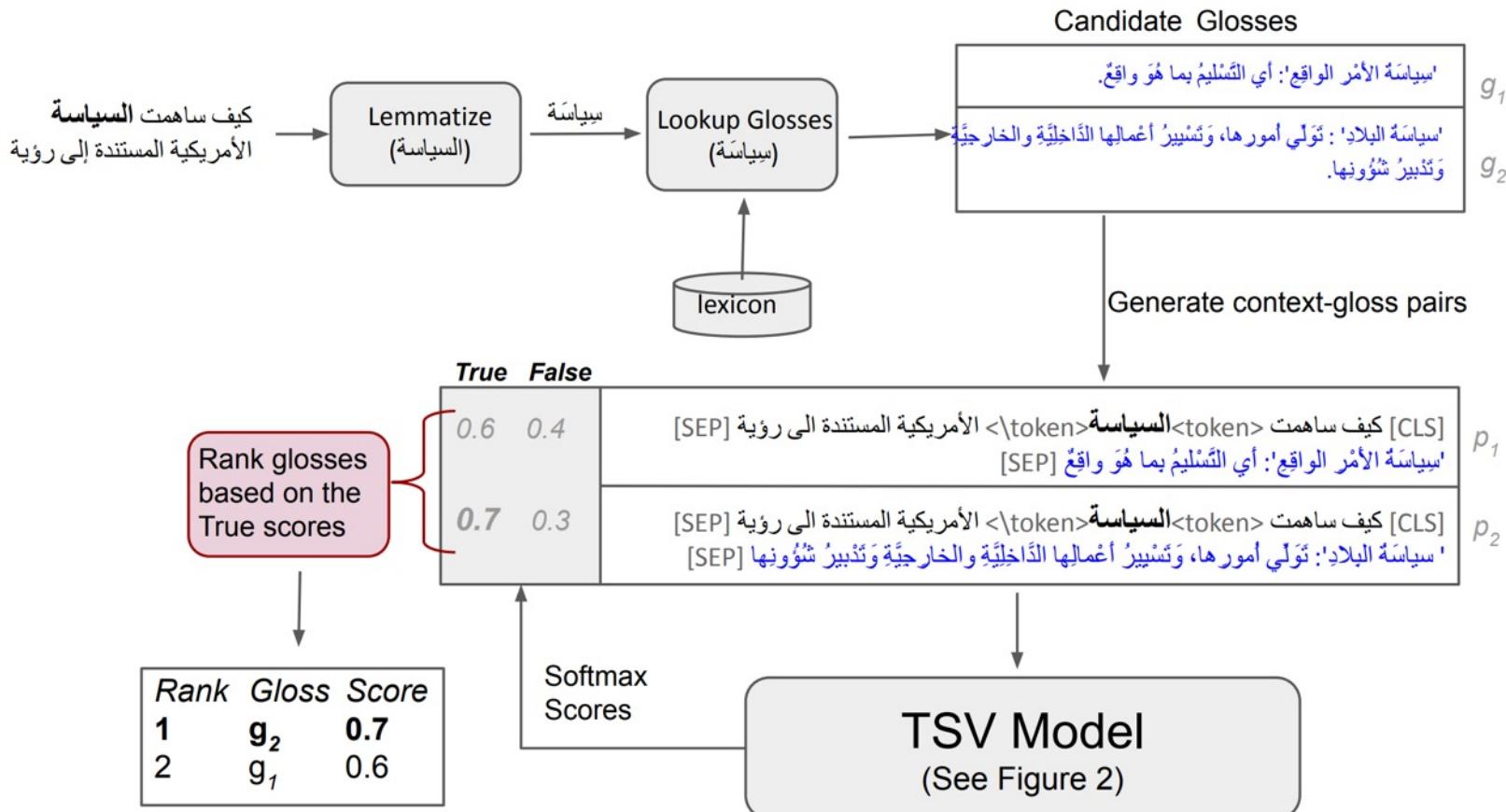
- every context selected in the test set should not be selected in the training set;
- every gloss should be selected in both the training and the test sets.

Datasets	Pairs	Count	Total
Training	True pairs	55,585	152,035
	False pairs	96,450	
Test	True pairs	4,738	15,172
	False pairs	10,434	
	Total	167,207	

Download: <https://ontology.birzeit.edu/downloads>

SALMA

End-to-end WSD system (using TSV)





Try

<https://sina.birzeit.edu/salma/>

سلمى SALMA

A corpus and model for Arabic Word Sense Disambiguation (WSD).

Version: 1.0 (updated on 22/10/2023)

قصيدة من عيون الشعر

WSD

◀ قصيدة (قصيدة 1 [303044571](#)): مجموعة من الأبيات الشعرية متّحدة في الوزن والقافية والرّوّي وهي تتكون من سبعة أبيات فاكثر "قصيدة غزلية".

بيّت القصيدة : البيت المتضمن غاية الشاعر، أو أنفس أبياتها، أو مثل يُضرب في تفضيل بعض الشيء على كله - مطلع القصيدة : أول بيت منها

◀ من (من 1)

◀ عيون (عين 2 [303038477](#)): أجود كل شيء وأحسنه ونفيسه "قصيدة من عيون الشعر - عيون الفن".

◀ الشعر (شعر 1 [303029103](#)): كلام موزون مقفى قصدًا يعتمد على التخييل والتأثير؛ ليوحى بالحساسات مؤثرة وصور خيالية "شعر صافي الدبياجة

- نظم الشعر - ما الشعر إلا شعور المرء يُرسله ... عفو البديهة عن صدق وإيمان - إن من الشعر لحكمة [حديث] - حوما علمناه الشعر وما ينبغي له <يس /

69 ". أنشده الشعر : قرأه عليه - أوايد الشعر : ما لا تماثل جودته أو قوافيه الشاردة - ريبة الشعر : إلهة الشعر عند الوثنيين - شطرا بيت الشعر : الصدر

+ Description

+ Downloads

ArabGlossBERT

Arabic Word-Sense Disambiguation



Arabic Language Understanding Tools

..... Confirm

صورة لعيون جميله

Tokenize Lemmatize NER WSD

• صورة ([صورة 1](#) [303032440](#)): (الطبيعة والفيزياء) ما تراه العين مباشرةً أو من خلال عدسة أو في مرآة أو مرئٍ عنها على سطح ما

• لعيون ([عيون 2](#) [303038475](#)): (التشريح) عُضو الإبصار في الإنسان والحيوان "له عينان كعَيْنَيِّ الصقر - ألا إنما العينان للقلب رائِدٌ ... فما تألفُ العينان فالقلب أَلِفُ - حَفَرَ جَعْنَاكَ إِلَى أَمْكَ كَيْ نَقَرَ عَيْنَهَا > طه/ 40 - حَوْلَتْصَنَعَ عَلَى عَيْنَيِّي > طه/ 39 : لتصنع تحت رعايتي وحفظي وإكرامي ". أَخَذَ بعَيْنَ الاعتبار : قدر، راعى أمرًا ما - أصابته العَيْنُ : حُسِيدٌ - أَعْمَضَ عَيْنَهُ عنِه : تجاهله، تغافله - أَنْتَ عَلَى عَيْنَيِّي : يقال في الإكرام والحفظ جميًعاً - إنسان العَيْنَ :

<https://sina.birzeit.edu/salma/>

BIRZEIT UNIVERSITY

جميع الحقوق محفوظة © 2022 جامعة بيرزيت

43

ArabGlossBERT

Arabic Word-Sense Disambiguation



Arabic Language Understanding Tools

..... Confirm

اعمل بجامعة بيرزيت وأحب زيت الزيتون

Tokenize Lemmatize NER WSD

birzeit.edu

جميع الحقوق محفوظة © 2022 جامعة بيرزيت

44

<https://sina.birzeit.edu/salma/>

Natural Language Understanding Tools and Datasets



Open Source

<https://sina.birzeit.edu/resources>



Resources

Download and try NLP/NLU datasets, corpora, tools and services

+ Lexicographic Database (150 lexicons)

حوسبة المعاجم

+ Arabic Ontology

الأنطولوجيا العربية

+ Dialect Corpora (Currasat)

كراسات مدونة العاميات

+ Arabic Synonyms

استخراج متزدفات

+ Named Entity Recognition (Wojood)

وجود - لاستخراج أسماء الاعلام

+ Word Sense Disambiguation (Salma)

سلمى - محلل دلالي

+ ArBanking77 Intent Detection

تحديد المقصود في المساعدات الآلية

+ Offensive Language Detection

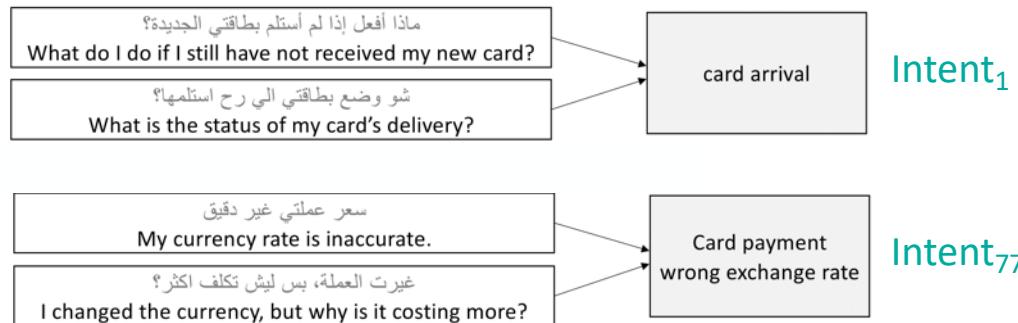
خطاب الكراهية بالعبرية

+ Lemmatizer

+ NLP Tools

Intent Detection

Chatbots need to detect intents before answering
→ build an Arabic intent dataset, and train BERT.



Our Contributions

❖ ArBanking77 dataset

- ArBanking77 dataset consists of **31,404 queries**.
- 2.4x larger than the Banking77 dataset.
- On average, there are 408 queries per intent

C	D	E	F	G	H	I
Qid	Question_en	Question_MSA1	n_MSA2	Question_PAL1	Question_PAL2	Intent_en
1	1 I am still waiting on my card?	ما زلت أنتظر بطاقي؟		بعدني يمتنى في البطاقة؟		card arrival
1	2 What can I do if my card still hasn't arrived after 2 weeks?	شو اساوى اذا ما وصلت بطاقتى ايش بلزم اعمل اذا بطاقتى ما وصلت بعد اسبوعين ماذَا علي افعل اذا لم تصل بطاقتى بعد اسبوعين؟		الى يمتنى اكتر من أسبوعين، كه صارلى اكتر من أسبوعين يمتنى، لسانها البطاقة ما		card arrival
1	3 I have been waiting over a week. Is the card still coming?	انا انتظر منه اكتر من أسبوع. هل ما زالت البطاقة قادمة؟				card arrival
1	4 Can I track my card while it is in the process of delivery?	ابيمكاني تتبع بطاقتى أثناء عملية التسليم؟				card arrival
1	5 How do I know if I will get my card, or if it is lost?	كيف يمكنني ان اعرف اذا اخدي بطاقتى ولا سافت؟	كيف يمكنني ان اعرف اذا ما كنت ساحصل على بطاقتى، او إذا ضاعت	كيف يمكنني ان اعرف اذا اخدي بطاقتى ولا سافت؟		card arrival
1	6 When did you send me my new card?	متى بعثت لي بطاقتى الجديدة؟	متى بعثت لي بطاقتى الجديدة؟	متى بعثت بطاقتى الجديدة؟		card arrival
1	7 Do you have info about the card on delivery?	عنكم معلومات عن بطاقه الدفع عند الاستلام؟	هل هناك معلومات عن البطاقة عند التسليم؟	هل لديك معلومات عن البطاقة عند التسليم؟		card arrival
1	8 What do I do if I still have not received my new card?	ما الذي علي فعله اذا لم استلم بطاقتى الجديدة؟	ما الذي علي فعله اذا لم استلم بطاقتى الجديدة؟	ما الذي علي فعله اذا لم استلم بطاقتى الجديدة؟		card arrival
1	9 Does the package with my card have tracking?	هل هناك تتبع للحفلة يتم تتبع الحزمة يوم يتوجه مع بطاقتى؟	هل هناك تتبع للحفلة يتم تتبع الحزمة يوم يتوجه مع بطاقتى؟	بقدر اثنين الحزمه يوم يتوجه مع بطاقتى؟		card arrival
1	10 I ordered my card but it still isn't here	طلب بطاقتى لكن طلبت بطاقتى لكنها ما زالت غير موجودة.	طلب بطاقتى لكن طلبت بطاقتى لكنها ما زالت غير موجودة.	بطاقتى الي وصيت عليها لسانها طلبت بطاقتى بين لسانها مش موجودة.		card arrival
1	11 Why has my new card still not come?	لماذا بطاقتى الجديدة بعدها؟	لماذا بطاقتى الجديدة بعدها؟	لدي بطاقتى مش موجودة بعدها؟		card arrival
1	12 I still haven't received my card after two weeks, is it broken?	صارلى اسبوعين ما استلمت بطاقتى، معمول ضاعت؟	للان لم استلم بطاقتى ما زلت لم استلم بطاقتى بعد اسبوعين ، هل فضلت؟	صارلى اسبوعين ما استلمت بطاقتى الى اسيوين ما استلمت بطاقتى، معمول ضاعت؟		card arrival
1	13 Can you track my card for me?	ابيمكانتك تتبع بطاقتى هل يمكنك تتبع بطاقتى من احلي؟	ابيمكانتك تتبع بطاقتى هل يمكنك تتبع بطاقتى من احلي؟	ممكن تتبعلي بطاقتى؟		card arrival
1	14 Is there a way to track the delivery of my card?	هل هناك طريقة تتبع تسليم بطاقتى؟	هل هناك طريقة تتبع تسليم بطاقتى؟	في طريقة تتبع تسليم بطاقتى؟		card arrival
1	15 It's been a week since I ordered my card and it's not arrived yet. Am I supposed to wait a week for my card to arrive?	في إمكانية تتبع البطاقة الى يمتنى بقدر اثنين البطاقة الى ايمتنى؟	امكنتك تتبع البطاقة هل سأتمكن من تتبع البطاقة التي تم ارسلها الي؟	في إمكانية تتبع البطاقة الى يمتنى بقدر اثنين البطاقة الى ايمتنى؟		card arrival
1	16 Will I be able to track the card that was sent to me?	منذ أسبوع ولم املك بطاقتى خلال أسبوع واحد. يجب ان اكون فلقا؟	منذ أسبوع ولم املك بطاقتى خلال أسبوع واحد. يجب ان اكون فلقا؟	في اشي داعي للقلق عشان بطاقت المها أسبوع بطاقتى مش معى، لازم افتق؟		card arrival
1	17 I don't have my card in 1 week. Should I be worried?					card arrival

❖ Intent detection model

- F1-scores: MSA (%92) and PAL (%90)

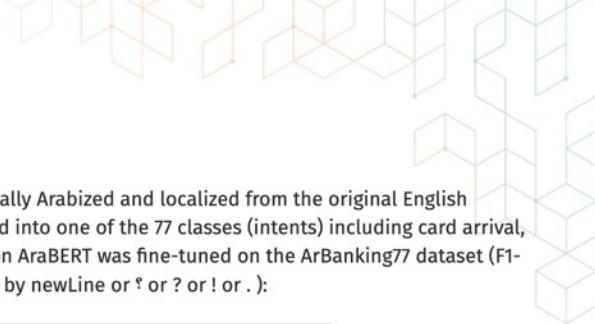
Download





SinaLab

News Team Resources



ArBanking77

A dataset and source-code for ArBanking77
Version: 1.0 (updated on 1/9/2023)

ArBanking77 consists of 31,404 (MSA and Palestinian dialect) that are manually Arabized and localized from the original English Banking77 dataset; which consists of 13,083 queries. Each query is classified into one of the 77 classes (intents) including card arrival, card linking, exchange rate, and automatic top-up. A neural model based on AraBERT was fine-tuned on the ArBanking77 dataset (F1-score 92% for MSA, 90% for PAL). Try the service (type sentences separated by newLine or ! or ? or ! or .):

Detect

- Downloads

ArBanking77 is available to download upon request for academic and commercial use.
[Request to download ArBanking77](#) (whole dataset 31,404 queries, MSA 15,537 queries, Palestinian Dialect 15,867 queries)
[GitHub](#) (download BERT training source code + sample data (~1K queries))
[Hugging Face](#) (download fine-tuned BERT model, ready to use)

+ API

- Publications

Mustafa Jarrar, Ahmet Birim, Mohammed Khalilia, Mustafa Erden, Sana Ghanem: [ArBanking77: Intent Detection Neural Model](#)


Copyright © 2023 Birzeit University

<https://sina.birzeit.edu/arbanking77/>

48

Offensive Language Detection in Hebrew

Offensive content {**Hate** | **Violence** | **Racist** | **Abusive**} in Hebrew is growing.

Type of Offence:

Class	Sub-Class	Definition
Offensive	Abusive	If the tweet contains direct or implicit insults using vulgar or street words.
	Hate	If the tweet contains criticism, attack, or degrade, directly or implicitly, because of race, color, religion, nationality, or gender.
	Pornographic	If the tweet promotes or invites any pornographic or sexual arousal.
	Violence	If the tweet endorses an act that involves physical harm towards any party, regardless of the reason.

Offensive Language Detection in Hebrew

- Collect 16K tweets,
- classify each tweet with different tags,
- then, train a deep learning model – to automatically detect offensive content.

Tweet	English Translation	Annotations
גירוש, הריסת בתים, מאסרי עולם ללא אפשרות חנינה, אחרת יהה עוד יותר גורע, להתייחס אליהם כמו אל מחבלים.	Deportation, demolition of houses, life sentences without the possibility of amnesty, otherwise it will be even worse. Treat them like terrorists.	Class: Violence, Hate Target: Palestinians Topic: punish Palestinians Phrase: Deportation, Demolition of houses, life sentences, terrorists
אין כבר הרתעה. לא מפחדים מהמשטראה. אני חשב שהגע הזמן על פि מראות ההפגנות וההתפרעות בימים האחרונים, כמו שאמר פעם רבין בתחילת האינתיפאדה: „לשבור להם את העצמות“. פה יש כבר אינטיפאדה של התפרויות.	There is no more deterrence. We are not afraid of the police. I think the time has come to face the demonstrations and riots, as Rabin once said at the beginning of the intifada: "to break their bones." There is already an intifada of riots here.	Class: Violence, Hate Target: Palestinians Topic: Demonstrations Phrase: Break their bones
ימח שמה זכרה של אילנה דין. העיתונאית הכי מנולת ושקנית שאני מכיר. ממש מרושעת.	May the name and memory of Ilana Dayan be remembered. The most depraved and lying journalist I know. Really sinister.	Class: Abusive Target: Ilana Dayan Topic: Journalism Phrase: Sinister, Depraved, Lying
פרצופו האמתי של @judash0 אבי ביטון נחשף לעיני כל. מדובר בשמאלי, אנטי ציוני, עוכר ישראל, בוגד שמו מומן ע"י הקרכן החדשה להפיל את שלטונו חימין ולהעלות את המפלגות הערביות לשולטן כדי להרabil למדינת כל אזרחיה.	Avi Beaton's true face clear now. This is a leftist, anti-Zionist, oppressor of Israel, a traitor who is financed to overthrow the right-wing government and bring the Arab parties to power in order to lead to a state for all its citizens.	Class: Hate, Abusive Target: Avi Bitton, Arab Parties Topic: politics Phrase: Traitor, Anti-Zionist
ה... וא לא אם אתה דרוזי, סורי, אנטי ציוני ומגנעל שכמך.	@rabea_bader is irrelevant if you are Druze, Syrian, anti-Zionist and disgusting like you.	Class: Hate, Abusive Target: Rabea Bader, Druze, Syrian Topic: Racism Phrase: disgusting, anti-Zionist
אתה לפחות לא משקר - הייתה ונשارة לאומן ערבי שרוצה בחורבן ישראל כמדינה יהודית.	@Ahmad_tibi At least you're not lying - you were and remain an Arab nationalist who wants the destruction of Israel as a Jewish state.	Class: Hate Target: Ahmad Tibi Topic: Political views Phrase:

Offensive Language Detection in Hebrew

Data set	# of training examples		Validation data	Test data	Accuracy	
	Our data	D_{OLaH}			HeBERT	AlephBERT
D_1	1,750	0	250 (ours)	500 (ours)	63%	68%
D_2	1,750	1,013	250 (ours)	500 (ours)	58%	63%
D_3	1,750	2,026	250 (ours)	500 (ours)	61%	63%
D_4	0	1,418	203 (D_{OLaH})	405 (D_{OLaH})	79%	86%
D_5	1,250	1,418	203 (D_{OLaH})	405 (D_{OLaH})	81%	79%
D_6	2,500	1,418	203 (D_{OLaH})	405 (D_{OLaH})	81%	82%
D_7	0	2,026	203 (D_{OLaH})	500 (ours)	60%	57%
D_8	2,500	0	250 (ours)	405 (D_{OLaH})	64%	69%

Offensive Language Detection in Hebrew

Dashboard Prototype

“Support free speech. Condemn hate speech. That is not so hard at all”

The dashboard features two donut charts. The left chart shows 80.39% of 6238 tweets are Offensive, while the right chart shows 19.61% of 1522 tweets are Not Offensive.

A line chart titled "Num among several durations" shows the count of tweets from May 2022 to July 2023. The count remains low until January 2023, then rises sharply to nearly 3000 by June 2023.

Hebrew Hate Speech in Twitter

Total Tweets : 7760

Category	Count
Offensive	6238
Not Offensive	1522

Most Recent Tweets: Last Updated: 9/5/2023

- Classification : Offensive (Publish):2023-07
- Publisher: : MagiOtsri
- Tweet: @nimrodkadosh קרע ברצף זמן חלל מבהנתי ([Link](https://twitter.com/MagiOtsri/status/1676158077307895808))
- English Translation : @Nimrodkadosh tore a sequence of space for me
- Arabic Translation : @nimrodkadosh ممزق سلسلة من المساحة بالنسبة لي

<< >>

Offensive Language Detection in Hebrew

Download Corpus and models

SinaLab / OffensiveHebrew

Type to search

Code Issues Pull requests Actions Projects Security

OffensiveHebrew Public

Edit Pins Watch 0

main 1 branch 0 tags Go to file Add file Code

naghmaghanim Add files via upload c0fcc0d 2 hours ago 19 commits

Training Add files via upload 3 months ago

data Add files via upload 2 months ago

LICENSE Update LICENSE 3 months ago

README.md Update README.md 2 months ago

tag_vocab.pkl Add files via upload 2 hours ago

README.md

Hebrew Corpus

This corpus contains offensive language in Hebrew manually annotated. The data includes 15,881 tweets, labeled with one or more of five classes (abusive, hate, violence, pornographic, or non-offensive). The corpus is annotated manually by Arabic-Hebrew bilingual speakers.

The Corpus

Hugging Face

Search models, datasets, users...

SinaLab/OffensiveHebrew

License: cc-by-nc-sa-4.0

Use with library

Model card Files Community 1 Settings

Edit model card

Downloads last month 0

Hebrew Corpus

This corpus contains offensive language in Hebrew manually annotated. The data includes 15,881 tweets, labeled with one or more of five classes (abusive, hate, violence, pornographic, or non-offensive). The corpus is annotated manually by Arabic-Hebrew bilingual speakers.

Model Download

Huggingface:

<https://huggingface.co/SinaLab/OffensiveHebrew>

Natural Language Understanding Tools and Datasets



Open Source

<https://sina.birzeit.edu/resources>



Resources

Download and try NLP/NLU datasets, corpora, tools and services

+ Lexicographic Database (150 lexicons)

حوسبة المعاجم

+ Arabic Ontology

الأنطولوجيا العربية

+ Dialect Corpora (Currasat)

كراسات مدونة العاميات

+ Arabic Synonyms

استخراج متزدقات

+ Named Entity Recognition (Wojood)

وجود - لاستخراج أسماء الاعلام

+ Word Sense Disambiguation (Salma)

سلمى - محلل دلالي

+ ArBanking77 Intent Detection

تحديد المقصود في المساعدات الآلية

+ Offensive Language Detection

خطاب الكراهية بالعبرية

+ Lemmatizer

+ NLP Tools

References

1. Mustafa Jarrar, Sanad Malaysha, Tymaa Hammouda, Mohammad Khalilia: SALMA: Arabic Sense-Annotated Corpus and WSD Benchmarks. In Proceedings of the 1st Arabic Natural Language Processing Conference (ArabicNLP), Part of the EMNLP 2023. ACL.
2. Amal Nayouf, Tymaa Hammouda, Mustafa Jarrar, Fadi Zaraket, Mohamad-Bassam Kurdy: Nâbra: Syrian Arabic Dialects with Morphological Annotations. In Proceedings of the 1st Arabic Natural Language Processing Conference (ArabicNLP), Part of the EMNLP 2023. ACL.
3. Mustafa Jarrar, Ahmet Birim, Mohammed Khalilia, Mustafa Erden, and Sana Ghanem: ArBanking77: Intent Detection Neural Model and a New Dataset in Modern and Dialectical Arabic. In Proceedings of the 1st Arabic Natural Language Processing Conference (ArabicNLP), Part of the EMNLP 2023. ACL.
4. Hanene Lqreina, Mustafa Jarrar, Mohammed Khalilia, Ahmed Oumar El-Shangiti, Muhammad AbdulMageed: Arabic Fine-Grained Entity Eecognition. In Proceedings of the 1st Arabic Natural Language Processing Conference (ArabicNLP), Part of the EMNLP 2023. ACL.
5. Mustafa Jarrar, Muhammad Abdul-Mageed, Mohammed Khalilia, Bashar Talafha, AbdelRahim El-madany, Nagham Hamad, Alaa' Omar: WojooodNER 2023: The First Arabic Named Entity Recognition Shared Task. In Proceedings of the 1st Arabic Natural Language Processing Conference (Arabic- NLP), Part of the EMNLP 2023. ACL.
6. Nouran Khallaf, Elin Arfon, Mo El-Haj, Jon Morris, Dawn Knight, Paul Rayson, Tymaa Hammouda, Mustafa Jarrar: Open-source thesaurus development for under-resourced languages: a Welsh case study. The 4th LDK Conference on Language, Data and Knowledge, Vienna, Austria, 12-15 September 2023
7. Nagham Hamad, Mustafa Jarrar, Mohammad Khalilia, Nadim Nashif: Offensive Hebrew Corpus and Detection using BERT. The 20th ACS/IEEE International Conference on Computer Systems and Applications (AICCSA). IEEE. Egypt. 2023
8. Sana Ghanem, Mustafa Jarrar, Radi Jarrar, Ibrahim Bounhas: A Benchmark and Scoring Algorithm for Enriching Arabic Synonyms. The 12th International Global Wordnet Conference (GWC2023), Global Wordnet Association. (pp.). San Sebastian, Spain, 2023
9. Sanad Malaysha, Mustafa Jarrar, Mohammad Khalilia: Context-Gloss Augmentation for Improving Arabic Target Sense Verification. The 12th International Global Wordnet Conference (GWC2023), Global Wordnet Association. (pp.). San Sebastian, Spain, 2023
10. Mustafa Jarrar, Mohammed Khalilia, Sana Ghanem: Wojoood: Nested Arabic Named Entity Corpus and Recognition using BERT. In Proceedings of the International Conference on Language Resources and Evaluation (LREC 2022), Marseille, France. 2022
11. Mustafa Jarrar, Fadi Zaraket, Tymaa Hammouda, Daanish Masood, Martin Waehlisch: Lisan: Yemeni, Iraqi, Libyan, and Sudanese Arabic Dialect Corpora with Morphological Annotations. The 20th ACS/IEEE International Conference on Computer Systems and Applications (AICCSA). Pages(-). IEEE. Egypt. 2023 arXiv, DOI 10.48550/ARXIV.2212.06468. 2023
12. Karim El Haff, Mustafa Jarrar, Tymaa Hammouda, Fadi Zaraket: Curras + Baladi: Towards a Levantine Corpus. In Proceedings of the International Conference on Language Resources and Evaluation (LREC 2022), Marseille, France. 2022
13. Mustafa Jarrar: The Arabic Ontology - An Arabic Wordnet with Ontologically Clean Content. Applied Ontology Journal, 16:1, 1-26. IOS Press. 2021
14. 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
15. Moustafa Al-Hajj, Mustafa Jarrar: LU-BZU at SemEval-2021 Task 2: Word2Vec and Lemma2Vec performance in Arabic Word-in-Context disambiguation. In Proceedings of the Fifteenth Workshop on Semantic Evaluation (SemEval2021) Task 2: Multilingual and Cross-lingual Word-in-Context Disambiguation (MCL-WiC). PP 748–755, Association for Computational Linguistics. 2021
16. Eman Naser-Karajah, Nabil Arman, Mustafa Jarrar: Current Trends and Approaches in Synonyms Extraction: Potential Adaptation to Arabic. In Proceedings of the 2021 International Conference on Information Technology (ICIT). PP 748–755, Association for Computational Linguistics. pp. 428-434, IEEE. 2021
17. Mustafa Jarrar, Eman Karajah, Muhammad Khalifa, Khaled Shaalan: Extracting Synonyms from Bilingual Dictionaries. The 11th International Global Wordnet Conference (GWC2021), Global Wordnet Association. (pp. 215-222). Pretoria, South Africa, 2021
18. Kareem Darwish, Nizar Habash, Mourad Abbas, Hend Al-Khalifa, Huseein T. Al-Natsheh, Houda Bouamor, Karim Bouzoubaa, Violetta Cavalli-Sforza, Samhaa R. El-Beltagy, Wassim El-Hajj, Mustafa Jarrar, Hamdy Mubarak: A Panoramic Survey of Natural Language Processing in the Arab World. Communications of the ACM, April 2021, Vol. 64 No. 4, Pages 72-81
19. Mustafa Jarrar: Digitization of Arabic Lexicons. Arabic Language Status Report. UAE Ministry of Culture and Youth. Pages 214-2017. Dec 2020
20. 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
21. Mustafa Jarrar, Hamzeh Amayreh, John P. McCrae: Representing Arabic Lexicons in Lemon - a Preliminary Study. The 2nd Conference on Language, Data and Knowledge (LDK 2019). Pages(29-33). CEUR, Volume 2402. ISSN:1613-0073. Leipzig, Germany. 2019
22. Diana Alhafi, Anton Deik, Mustafa Jarrar: Usability Evaluation of Lexicographic e-Services. The 16th IEEE/ACS International Conference on Computer Systems and Applications (AICCSA). Pages(1-7). IEEE. Abu Dhabi, UAE. 2019
23. 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, ISSN:2375-4699. December, 2018
24. Mustafa Jarrar: Search Engine for Arabic Lexicons. The 5th Conference on Translation and the Problematics of Cross-cultural Understanding. The Forum for Arab and International Relations. Doha, Qatar. December, 2018
25. Diab Abuaiadah, Dileep Rajendran, Mustafa Jarrar: Clustering Arabic Tweets for Sentiment Analysis. The 2017 IEEE/ACS 14th International Conference on Computer Systems and Applications. Pages(499-506). IEEE Computer Society. ISBN:9781538635810. (doi.10.1109/AICCSA.2017.162). Hammamet, Tunisia. 2017
26. Mustafa Jarrar, Nizar Habash, Faeg Alrimawi, Divam Akra, Nasser Zalmout: Curras: An Annotated Corpus for the Palestinian Arabic Dialect. Journal Language Resources and Evaluation. Pages(745-775). Volume(51).