

حوسبة المعاجم العربية

دراسة حالة



بقلم: د. مصطفى جرّار
جامعة بيرزيت، فلسطين

مقدّمة

تعرض هذه الورقة ملخصاً لتجربة بناء أضخم قاعدة بيانات لغويّة للغة العربيّة، وكذلك محرّك بحث للمعاجم العربيّة يحتوي على 150 معجمًا تمّت حوسبتها في جامعة بيرزيت بفلسطين.¹⁴² ومحرّك البحث هذا متاح للعامة ولجمهور الطلبة والباحثين والمترجمين ومتعلّمي اللغة وغيرهم، وهو يتيح للمستخدم البحث عن الكلمة في 150 معجمًا عربيًا ومتعدّد اللغات واسترجاع تعريفاتها ومرادفاتها وكذلك ترجماتها المتخصّصة. وهذه الترجمات تعتمد على المعاجم المدقّقة والمنقّحة، وهي أدقّ بكثير من جميع المترجمات الآليّة المتاحة (مثل جوجل للترجمة)، وتجدر الإشارة إلى أنّ محرّك البحث هو الأوّل عالميًا، حيث لا يوجد محرّك بحث لمعاجم أيّ لغة أخرى، حتّى الإنجليزيّة. ويمكن لمطوّري البرمجيات كذلك استعمال محرّك البحث عبر واجهة خاصّة (API) للوصول إلى الترجمات والمترادفات واستخدامها داخل برمجياتهم وتطبيقاتهم. وقد تمّت حوسبة المعاجم بطريقة يدويّة على مدى تسع سنوات، وبعد ذلك قمنا بتوحيد هذه المعاجم في قاعدة بيانات واحدة تشتمل على معاجم لغويّة تقليديّة وحديثة، ومسارد، ومكانز، ومعاجم ثنائيّة وثلاثيّة اللغة، وقواعد بيانات تصريفية واشتقاقية. وتعدّ الأنطولوجيا العربيّة أهمّ المصادر اللغويّة في محرّك البحث، وهي شجرة أو تصنيف لمعاني الكلمات العربيّة وتمثّل لها بلغة المنطق، بحيث يستطيع الحاسوب فهم المعاني ومعالجتها. وقد تمّ بناء هذه الشجرة كمشروع بحثي منفصل في جامعة بيرزيت،¹⁴³⁻¹⁴⁴ ونعمل حاليًا على ربط جميع مدخلات المعاجم

وصولًا إلى شبكة لغويّة محوسبة للغة العربيّة تُسمّى Big Linguistic Data Graph، ومن ثمّ ربط هذه الشبكة باللغات الأخرى.¹⁴⁵

وبأني هذا العمل ضمن مشروع بحثي غير ربحي وطويل الأمد لخدمة اللغة العربيّة، ولإغناء الإنترنت بمحتوى عربيّ نوعي، إذ سيتمكّن المستخدمون من إيجاد المصطلح العربيّ الذي يناسب متطلّباتهم وترجمات متعدّدة، خاصّة وأنّ قاعدة البيانات تحتوي على عدد ضخم من معاجم المصطلحات المعاصرة وفي شتى العلوم والمجالات العلميّة والهندسيّة والتجاريّة والأدبيّة وغيرها. إضافةً إلى ذلك، فإنّ محرّك البحث متاح للعامة مجانًا من خلال موقع جامعة بيرزيت الإلكترونيّ (<https://ontology.birzeit.edu>) الذي يقدّم تعريفًا بالمحرّك وبالمشاريع الأخرى ذات العلاقة على هذا الرابط (<https://ontology.birzeit.edu/about-ar>)، بالإضافة إلى روابط للأوراق العلميّة التي بُني محرّك البحث على أساسها. وتجدر الإشارة إلى أنّ محرّك البحث تمّ تطويره وفق المعايير والمقاييس الصادرة عن منظمّة الشبكة العالميّة (W3C)، وخاصّة المعايير المتعلّقة بأسلوب نشر البيانات على الإنترنت، وكذلك المعايير المتعلّقة بتمثيل وتبادل البيانات اللغويّة.

وفيما يلي يقدّم الشكل التالي مثالاً على طريقة استرجاع ترجمات كلمة «attribute» من عدّة معاجم ومن الأنطولوجيا، ويستطيع الباحث هنا التحكّم فيما إذا كان يرغب في استرجاع ترجمات، و/أو مترادفات، و/أو تعريفات فقط. وتظهر نتائج المعاجم في الجهة اليسرى من الشكل، ونتائج الأنطولوجيا في الجهة اليمنى:



attribute


 Translations Synonyms Definitions

Ontology

Dictionaries

Morphology

About License


 34 results (0.50 secs)

attribute pridecate محمول
Propositions منطقياً - ما يقال على موضوع ومنه القضايا الحملية
attributives

attribute pridecate محمول
قال ابن سينا: "المحمول هو المحكوم به أنه موجود أو غير موجود لشيء آخر"

attribute pridecate محمول
ميتافيزيقياً: ويسمى صفة، وهي خاصة ذاتية للجوهر، ويطلق بوجه خاص على صفات الباري. ومشكلة الإسلامي، وحولها دار معظم نشاط المدارس الكلامية.

attribute صفة
The Unified Dictionary of Social Sciences Terms

attribute صفة مُمَيَّزَة | سَمَة
صفة تميز ملفاً أو سجلاً أو حقلاً وتبين نوعه أو كيفية التعامل معه، أو تبين أي معلومات عن البيانات المخزونة فيه.

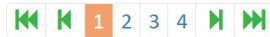
attribute صفة | خاصية
Statistical Terms

inspection by attribute الفحص بالصفات
Statistical Terms

sampling for attribute المعاينة للخاصية | المعاينة للصفة
Statistical Terms

quality attribute معيار الجودة
إحدى الخواص التي تساهم في جودة المادة.
The Unified Dictionary of Nutrition Technologies Terms

attribute samples عيّنات وصفية
The Unified Dictionary for terminologies of General and Nuclear Physics Terms



ONTOLOGY

attribute - 1 results (0.28 secs)

attribute | صِفَةٌ | خَصْلَةٌ | مُسَنَدٌ | حَالَةٌ | نَعَتْ | مَرِيَّةٌ | مَبْرَةٌ | سَكَلٌ | هَيْبَةٌ | سِمَةٌ

An abstract that describes qualities of an entity without using units of measure

مجرد يصف ويُعبّر عن قيمة خاصية لشيء ما دون استعمال وحدة قياس
example: اللون البرتقالي هو صفة لخاصية اللون في البرتقال

293263

Physical Attribute صِفَةٌ مَادِّيَّةٌ
An attribute that describes the value of a physical quality

صفة للتعبير عن قيمة خاصية مادية
الأخضر هي صفة مادية لخاصية لون عيون بعض الأشخاص

293265

Abstract Attribute صِفَةٌ مَجْرَدَةٌ | صِفَةٌ إِعْتِبَارِيَّةٌ
Attribute that describes an abstract quality of things

صفة تعبّر عن قيمة خاصية مجردة لشيء ما
كريم هي صفة مجردة لخاصية العطاء عند الإنسان

293264

الرّسْم التّام

ما يتركب من الجنس القريب والخاصة كتعريف الإنسان بالحيوان الضاحك

41816 ©Al-Jirjani Definitions

الرّسْم النّاقص

ما يكون بالخاصة وحدها أو بها وبالجنس البعيد كتعريف الإنسان بالضحك أو بالجسم الضاحك أو بعرضيات تختصّ جملتها بحقيقة واحدة كقولنا في تعريف الانسان أنه ماشي على قدميه عريض الأظفار للمزيد...

41817 ©Al-Jirjani Definitions

الاستخدام الورقي، بينما يحتاج الاستخدام الحاسوبي إلى تراكيب بنويّة مختلفة تسهّل عمليّات الإضافة والاسترجاع والبحث.

وتشتمل قاعدة البيانات اللغويّة الحاليّة على أنواع مختلفة من المعاجم، مثل المعاجم اللغويّة التقليديّة القديمة والجديدة، وكذلك المسارد (glossaries) التي تشرح المصطلحات، والمكانز (thesauri) التي تحتوي على مترادفات، ومعاجم ثنائيّة وثلاثيّة اللغة، ومعاجم الفروق اللغويّة، وقواعد بيانات تصنيفيّة واشتقاقية. كما تغطّي قاعدة البيانات شتى المجالات مثل العلوم الطبيعيّة، والهندسة، والطب، والاقتصاد، والأدب، والإنسانيّات، والفلسفة،

حوسبة المعاجم

قمنا بحوسبة حوالي مائة وخمسين معجمًا عربيًا ومتعدّد اللغات، واضطررنا لطباعة غالبيتها يدويًا، خاصّةً وأنّه لا توجد نسخ إلكترونيّة للغالبية المعاجم حتّى في حوزة مؤلّفيها. وفي المراحل التالية، عملنا على تطوير برمجيات وخوارزميات لإعادة تشكيل بنويّة (re-structuring) لمدخلات المعاجم¹⁴⁶ بحيث يتمّ الفصل بين المدخلة المعجميّة والتعريف والمشتقات والتصريفات والترجمات والمترادفات، انتهاءً بحفظها في قاعدة بيانات واحدة.¹⁴⁷ وجاءت الخطوة الأخيرة ضروريّةً لكون صناعة المعاجم التقليديّة تركّز في تصميمها على

بحدوده وبصفاته الجوهرية المُميّزة، وإعطاؤه رقمًا فريدًا، ثم يتمّ تصنيف المفهوم إلى أجناسه الأدنى، وهكذا، بحيث يكون التصنيف تصنيفًا مانعًا. وقد تمّ استخدام المنطق الوصفيّ الحديث (Description Logic) كلغة لتمثيل الأنطولوجيا العربية ولتمكين الحاسوب من فهمها والاستنتاج منها. وبالرغم من أنّ الأنطولوجيا العربية مصمّمة لاستخدامها في التطبيقات الحاسوبية، فإنّه يمكن استخدامها كمعجم؛ ولكن تجدر الإشارة إلى أنّه، إضافةً لتصنيفها للمعاني، فإنّ ما يميّز الأنطولوجيا عن المعاجم أيضًا، أنّه يتمّ التحقّق من تعريفات وتصنيفات الأنطولوجيا بالاعتماد على ما وصلت إليه العلوم، خلافًا للمعاجم في اعتمادها على شيوع الدلالة المستخدمة بين المتحدّثين. هذا بالإضافة إلى أنّ المفاهيم في الأنطولوجيا هي كليات/أصناف لأفراد.

ومشروع الأنطولوجيا مشروع طويل الأمد، وقد تمكّننا من إنجاز ما يزيد عن ألف مفهوم، وهي المفاهيم الأعلى والأكثر تجريديًا في اللغة العربية، بالإضافة إلى عشرة آلاف مفهوم آخر تمّ تعريفها وتصنيفها جزئيًا. ويمكن الوصول إلى شجرة الأنطولوجيا كاملةً، واستكشاف المفاهيم الأدنى عبر محرّك البحث الذي يمكن الوصول إليه عن طريق الرابط:

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

The screenshot shows a search interface for the term 'entity'. The search bar contains 'entity' and shows 5 results in 0.33 seconds. The results are listed as follows:

- entity** | كَيْتُونَةٌ | كائِنٌ
 - Whatever existed or will exist, and can be realized or imagined
 - أَيّما وُجِد أو سيوجد ونستطيع إدراكه أو تخيله
 - example: كل شيء على ما يرام 293198
- dependent entity** | مُتَعَلِّقٌ | مُتَعَلِّقٌ | مُتَعَلِّقٌ | مُتَعَلِّقٌ
 - specifically dependent entity
 - An entity whose existence is dependent on the existence of other entities
 - شيء يعتمد وجوده على وجود أشياء أخرى
 - example: طول المبنى منوط بوجود المبنى وآلا فلا طول له 293201
- physical object** | مَوْجُودٌ مَادِّيٌّ | مَحْسُوسٌ | مَلْمُوسٌ | كَائِنٌ | حَقِيقِيٌّ | مَجَسَّدٌ
 - material entity
 - An object that occupies space, and is realized by senses or measuring tools
 - موجود يشغل حيزًا مكانيًا، ويدرك بذاته بالحواس أو بأدوات القياس
 - example: لكل مَوْجُود مَادِّيٍّ حجم يمكن قياسه أو حسابه مهما صغر أو كبر 293254
- social object** | immaterial entity | كيان اجتماعي
 - An object that is realized for its social existence, and can be represented by physical objects
 - موجود يدرك لذاته اعتبارًا، ويمكن أن يمثله موجود مادي
 - example: جامعة بيرزيت كمؤسسة هي موجود اعتباري ولها ممثل قانوني يمثلها أمام الآخرين 293255
- formal entity** | حَدٌّ رِياضِيٌّ | حَدٌّ هَنْدَسِيٌّ | كَيْتُونَةٌ رِياضِيَّةٌ
 - An abstract entity that is formally defined either mathematically or logically, and does not need to be proved
 - شيء مجرد اصطلح عليه رياضيًا أو منطقيًا، وليس بحاجة إلى برهان
 - example: باي هو حد رياضي 293282

والفنّ، والاقتصاد. وتجدر الإشارة إلى أنّ محرّك البحث في هذه المرحلة يُظهر الترجمات من وإلى الإنجليزية، إلّا أنّ قاعدة البيانات تحتوي العديد من اللغات الأخرى، خاصّةً الفرنسية، والتي سنعمل لاحقًا على توفيرها ضمن محرّك البحث.

حقوق الملكية

قمنا بالتواصل الفرديّ مع أصحاب المعاجم وحقوق الملكية الفكرية للحصول على تراخيص باستخدام معاجمهم ضمن محرّك البحث، مشيرين إلى أنّ محرّك البحث، في تصميمه، يعرض اسم المعجم ورمز حقوق الملكية بجانب البيانات التي يتمّ استرجاعها. وعند الضغط على اسم المعجم يظهر اسم المؤلف والناشر وروابط لصفحاتهم الإلكترونية ولكيفية شراء النسخة الورقية من المعجم، ما شجّع أغلبهم على منحنا الترخيص باستعمال معاجمهم. إضافةً إلى ذلك، قمنا بتوقيع مذكرة تفاهم مع منظمة الألكسو، للتعاون في مجال حوسبة اللغة العربية وهندسة المعاجم، حيث منحونا الإذن باستخدام حوالي خمسين معجمًا ثلاثي اللغة صدرت عن مركز تنسيق التعريب بالرباط لحوسبتها وإضافتها إلى قاعدة البيانات الخاصّة بمحرّك البحث.

التطبيقات العلمية والعملية

يتيح محرّك البحث المعجمي، جنبًا إلى جنب مع الأنطولوجيا، خدمات لغوية عديدة لمجموعتين رئيسيتين تشمل الأولى الباحثين واللغويين والمترجمين والطلبة والجمهور العربيّ عمومًا، إضافةً إلى متعلّمي اللغة العربية، وذلك لاسترجاع مترادفات وترجمات وتعريفات.¹⁴⁸ وتشمل الثانية الشركات التي تحتاج لبيانات لغوية من أجل تطوير تطبيقات حاسوبية، مثل تطبيقات الترجمة الآلية، أو البحث الدلاليّ، أو التحليل اللغويّ، أو التدقيق الإملائيّ، أو استخراج البيانات، وغيرها. ويمكن للعاملين في هذه الشركات الاستفادة من محرّك البحث عبر واجهة خاصّة لمطوري البرمجيات (API) تستطيع الشركات من خلالها استرجاع مترادفات وترجمات واستخدامها داخل برمجياتها وتطبيقاتها مباشرة. وفي كلتا الحالتين، فإنّ ما يقدمه المشروع يُسهّم في تقديم ترجمات وتعريفات دقيقة ومتخصّصة، تقلّل من هامش الخطأ - بنسبة كبيرة - في البحث المعرفيّ والبحث، أو في التطبيقات التكنولوجية ذات الصلة.

تصميم محرّك البحث

من أهمّ سمات محرّك البحث اعتماده على المعايير والمقاييس الصادرة حديثًا عن منظمة الشبكة العالمية (W3C) والتزامه بها، خاصّةً تلك المتعلقة بأسلوب نشر المعلومات والمعروفة بـ (Best Practices for Publishing) (Cool URLs, Simple, Linked Data). وبالتالي، فقد تمّ تصميم الروابط لتكون (Stable, Manageable, Linkable) على الإنترنت.¹⁴⁹ إضافةً إلى ذلك، تمّ الالتزام بالمعايير والمقاييس المتعلقة بتمثيل البيانات اللغوية والمعروفة بـ (W3C Lemon model)،¹⁵⁰ حيث يتيح محرّك البحث الوصول إلى البيانات المسترجعة ممثلةً باستخدام (Lemon model).¹⁵¹ انظر، مثلًا، تمثيل معنى كلمة "موضع" في المعجم الفلسفي: <https://ontology.birzeit.edu/lemon/lexicalconcept/300000117>

الأنطولوجيا العربية

الأنطولوجيا هي فرع من فروع الفلسفة وتعني علم الوجود، وهي تعلقو الإبستمولوجيا في نظرية المعرفة، ولكن هذا المصطلح أصبح رائجًا مؤخرًا في علم الحاسوب وهندسة المعرفة، ويعني وصف المعرفة مفاهيميًا ودلاليًا، وتمثيلها بلغة المنطق.¹⁵²⁻¹⁵³ والأنطولوجيا العربية هي مشروع آخر ومستقل عن محرّك البحث المعجمي، نعمل عليه في دائرة علم الحاسوب بجامعة بيرزيت،¹⁵⁴⁻¹⁵⁵ ويهدف إلى وصف وتصنيف مفاهيم الكلمات العربية. أي أنّ الأنطولوجيا العربية وبمعنى آخر، هي شجرة مفاهيم ومعاني الكلمات العربية، وليس الكلمات نفسها، حيث يتمّ وصف كلّ مفهوم/معنى لكل كلمة عربية

1- عدم توفر مصادر لغوية محوسبة

بالرغم من وجود كمّ كبير من المعاجم العربيّة المطبوعة ورقياً، فإنّ عدد المصادر اللغويّة المحوسبة والمتاحة للبحث العلميّ ولمطوّري البرمجيات ضئيل، ممّا يحدّ من دعم اللغة العربيّة في معظم التطبيقات الحاسوبية. وعلى الرغم من أنّ سوق التقنيّات في العالم العربيّ كبير، فإنّ مطوّري البرمجيات يواجهون تحديّات كبيرة في إيجاد مصادر لغوية يمكن استعمالها بسهولة وبتكلفة معقولة.

2- ضعف الصناعة المعجمية وتوانيتها عن مواكبة النظريّات اللسانيّة الحديثة

ثمّة نقص في أنواع جديدة من المعاجم، مثل: مكانز المترادفات، ومعاجم الفروق الدلاليّة، والشبكات الدلاليّة، ومعاجم السلوك التصريفيّ للأفعال، ومعاجم المصطلحات والترجمات الحديثة والكلمات العربيّة المستجدة. وقد أصبحت الحاجة إلى مثل هذه الأنواع من المعاجم، والتي لم تكن موجودة أو منتشرة قديماً، ملحة مع تطوّر احتياجات المجتمع، ووسائل تعليم اللغة، وظهور برمجيات ذكيّة أصبحت اللغة ركنًا أساسياً فيها.

3- عدم وجود منهجيات موحّدة ومعيارية في الصياغة المعجمية العربيّة

فمثلاً، لا يوجد اتّفاق على كيفية صياغة ووسم المدخلات المعجمية، ويظهر ذلك جليّاً في التباين الحاصل بين صياغة المعاجم التي صدرت في السنوات الأخيرة.

4- عدم توفر قواعد بيانات لغوية شاملة

فعلى سبيل المثال، لا توجد قائمة شاملة وموحّدة بالجزور، أو بالأوزان الصرفيّة والاشتقاقية، أو بفروع الكلمات (lemmas)، أو بالاشتقاقات، أو غيرها. ويشكّل هذا النقص تحديّاً لكلّ من يريد المساهمة في تأليف معجم، أو إنشاء قاعدة بيانات لغوية، أو بناء تطبيق حاسوبيّ.

5- عدم توفر ميزات البحث العلميّ الرصين والموجّه نحو حوسبة اللغة العربيّة لدعمها في التطبيقات الحاسوبية.

تتركز أبحاثنا في جامعة بيرزيت حالياً في ثلاثة مسارات:

الأول: إغناء المحرّك المعجميّ بمزيد من التعريفات والشروحات اللازمة، وتوسيع الأنطولوجيا العربيّة بإضافة المزيد من المفاهيم في المستويات الدنيا من الشجرة الدلاليّة.

الثاني: ربط مدخلات المعاجم العربيّة مع شبكة المفردات الإنجليزيّة (Wordnet)، وهي شبكة مفردات تمّ تطويرها في جامعة برنستون وربطها بغالبية اللغات في العالم،¹⁵⁶ وبالتالي، فإنّ ربطها بالمدخلات المعجمية العربيّة يؤدي تلقائياً إلى ربط المدخلات العربيّة بسائر اللغات العالميّة الأخرى.

الثالث: وهو المشروع المستقبليّ الذي شرعنا بأولى خطواته، ويقوم على ربط المعاجم وجميع البيانات اللغوية المتوفرة لدينا ببعضها، وصولاً إلى بناء شبكة بيانات لغوية شاملة من حيث الكمّ ومن حيث المستويات اللغوية الدلاليّة والتصريفية والاشتقاقية، والمقابل في اللغات الأخرى. وبدوره، فإنّ هذا المسار يتضمّن ثلاثة مستويات مترابطة:

- تصريفياً. نعمل على وصف كلّ مدخلة معجمية بجذعها (lemma)، وبالتالي تصبح كلّ مدخلات المعاجم مترابطة تصريفياً، بغضّ النظر عن السوابق واللواحق في المدخلات، أو الصيغ المختلفة من نفس المدخلة في المعاجم الأخرى.
- اشتقاقياً. نقوم بالربط بين المدخلات التي تختلف عن بعضها اشتقاقياً (derivation)، إن كانت فعلاً، أو اسم فاعل، أو اسم مفعول، أو مصدرًا، أو اسم تفضيل، ... إلخ.
- دلائلياً. نستعمل الأنطولوجيا العربيّة لتفعيد وربط المعاجم مفاهيمياً، حيث نقوم بربط كل دلالة لغوية (lexical concept) لكلّ مدخلة معجمية في كلّ معجم بمفهومٍ مقابل في الأنطولوجيا (ontological concept)، وبذلك، تصبح مدخلات المعاجم مترابطة دلائلياً فيما بينها.

وبالرغم من إدراكنا للجهود الكبيرة التي يتطلّبها بناء شبكة لغوية شاملة ومحوسبة للغة العربيّة، فإنّ من شأنها - حين تتحقّق - أن تُحدث نقلة جذريّة في المعرفة اللغوية والدلاليّة وتؤسّس لمنهجيات جديدة للبحث النظريّ والتطبيقيّ.