

The Design and the Construction of the Traditional Arabic Lexicons Corpus (The TAL-Corpus)

Majdi Sawalha¹

¹ Department of Computer Information Systems, King Abdullah II School for Information Technology, The University of Jordan, Amman, Jordan. E-mail: sawalha.majdi@ju.edu.jo

Correspondence: Majdi Sawalha, Department of Computer Information Systems, King Abdullah II School for Information Technology, The University of Jordan, Amman, 11942, Jordan. E-mail: sawalha.majdi@ju.edu.jo

Received: Oct. 20, 2018

Accepted: Nov. 28, 2018

Online Published: January 9, 2019

doi:10.5539/mas.v13n2p95

URL: <https://doi.org/10.5539/mas.v13n2p95>

Abstract

Arabic lexicography is a well-established and deep-rooted art of Arabic literature. Computational lexicography, invests computational and storage powers of modern computers, to accelerate long-term efforts in lexicographic projects. A collection of 23 machine-readable dictionaries, which are freely available on the web, were used to build the Corpus of Traditional Arabic lexicons (the TAL-Corpus). The purpose for constructing the TAL-Corpus is to collect and organize well-established and long traditions of traditional Arabic lexicons which can also be used to create new corpus-based Arabic dictionaries.

The compilation of the TAL-Corpus followed standard design and development criteria that informed major decisions for corpus creation. The corpus building process involved extracting information from disparate formats and merging traditional Arabic lexicons. As a result, the TAL-Corpus contains more than 14 million words and over 2 million word types (different words).

The TAL-Corpus was applied to create useful morphological database. This database was automatically constructed using a new algorithm which is informed by Arabic linguistics theory. The newly developed algorithm processed the text of the TAL-Corpus and as result it extracted 2 781 796 entries. These entries were stored in the morphological database where each represents a word-root pair (*i.e.* an Arabic word and its root).

A comparative evaluation of the TAL-Corpus and other three Arabic corpora showed that the lexical diversity of its vocabulary scored higher. Moreover, its coverage was computed by comparing words and lemmas against their equivalents of other corpora where it scored about 67% when comparing words and 82% when comparing lemmas.

Keywords: lexicography, traditional Arabic lexicon, corpora, dictionary building, the TAL-Corpus

1. Introduction

Lexicography is the applied part of lexicology. It is concerned with the design and construction of lexicons. Lexicography defines the process of collating, ordering of entries, derivations and their meaning, depending on the aim of the lexicon to be constructed and its size. Lexicography is defined as “...the branch of applied linguistics concerned with the design and construction of lexica for practical use.” (Eynde & Gibbon 2000). Moreover, lexicology is also defined as “...the branch of descriptive linguistics concerned with the linguistic theory and methodology for describing lexical information, often focusing specifically on issues of meaning.” (Eynde & Gibbon 2000). Long-term efforts in lexicographic projects have been greatly accelerated since the advent and use of computers: which is known as computational lexicography. However, constructing a large-scale broad-coverage lexicon involves time-consuming development of specifications, design, collection of lexical data, information structuring, and user-oriented presentation formatting (Eynde & Gibbon 2000).

Corpora have been used to construct dictionaries since the release of the Collins-Birmingham University International Database COBUILD. Computer technology was used in the four stages of constructing COBUILD: data-collection, entry-selection, entry construction and entry-arrangement (Ooi 1998). Similarly, the Oxford English Corpus was created to construct the Oxford English Dictionaries. It consists of 2.5 billion words of 21st century English which provides evidences of language use and development. It also draws an accurate picture of the language today. It contains text from literary, novels, specialist journals, everyday newspapers, magazines, blogs, emails, and Internet message boards. These texts were collected from all parts of the world, namely; the UK,

the United States, Ireland, Australia, New Zealand, the Caribbean, Canada, India, Singapore, and South Africa (The Oxford Dictionaries, 2018).

Arabic corpora started to appear in the late 1980s. They differ in size, type, purpose of development, and the materials used to develop them (Al-Sulaiti & Atwell 2006). Freely available Arabic corpora were surveyed by Zaghouni (2014). The survey categorized 66 freely available Arabic corpora into; (i) Raw Text Corpora such as monolingual corpora, multilingual corpora, dialectal corpora, and web-based corpora. (ii) Annotated Corpora such as named entities, error annotation, part-of-speech, syntax, semantic, and anaphora. (iii) Lexicons such as lexical databases and words lists. (iv) Speech Corpora such as audio recordings and transcribed data. (v) Handwriting Recognition Corpora such as scanned and annotated documents. (vi) Miscellaneous Corpora such as questions/answers, comparable corpora, plagiarism detection and summaries.

The first freely available Arabic corpus was the Corpus of Contemporary Arabic (Al-Sulaiti & Atwell 2006). It contains 1 million words collected from newspapers and magazines text. Most monolingual raw text of Arabic corpora cover the news domain. Examples of such corpora are; OSAC: Open Source Arabic Corpora (Saad & Ashour 2010); Khaleej-2004 corpus (Abbas & Smaili 2005); Watan-2004 corpus (Abbas et al, 2011); KACST Arabic Newspaper Corpus (Al-Thubaity et al. 2013); Arabic Words Corpus which is a frequency list of 1.5 million words collected by Al-Saadi. The International Corpus of Arabic (ICA) project by Bibliotheca Alexandrina (BA) was planned to construct a corpus that contains 100 million from Press, Net articles, books and academic text sources. The ICA was planned to cover the Arabic language as being used all over the Arab world (Alansary & Nagi, 2014). Alrabiah et al, (2013) developed the King Saud University Corpus of Classical Arabic (KSUCCA) which consists of around 50 million words. It was collected from authenticated Classical Arabic texts. It was constructed to study the language of the Qur'an using distributional lexical semantics.

Ismail, et al, (2014) developed a set of computational tools and corpus resources that would facilitate research in historical semantics and etymological lexicography. They constructed Historical Arabic Corpus (HAC) with around 45 million tokens in the first phase of development. They analyzed the collected text automatically and they annotated the text with linguistic information such as; part-of-speech, root, and morphological pattern. This corpus was collected from about 500 sources that represent 1 600 years of continuous language use which represent the entire recorded history of the Arabic language. Additionally, they developed a corpus builder that integrates a stemmer with a tagger to process and annotate documents, and then compile them into an XML-formatted corpus. They also created an indexer, a search engine, a concordancer and a dictionary editor that together facilitate searching and extraction of linguistic knowledge from HAC. Also, these tools facilitate the compilation of dictionary entries in a hypothetical dictionary on historical principles.

Since the release of COBUILD, corpora proved to be excellent resources for developing new dictionaries. In addition to lexical information, corpora can provide more useful information that would enrich dictionaries such as idioms, phrases, collocates, word sketches and thesauri of words. Arabic corpora were not yet used to construct most existing Arabic dictionaries (Al-Sulaiti & Atwell, 2006). Ghazali & Barham (2001) criticized existing Arabic dictionaries for literal translations of their lexical items and the lack of idioms; phrasal verbs; collocations; and new words entering the language. Based on a corpus of 1.5 million words, they investigated the different meaning of the verb أَتَى 'ahada 'he took' on both the corpus and Al-Waseet Arabic dictionary. In addition to its literal meaning, they discovered two additional meanings using the corpus which were not mentioned in Al-Waseet Arabic dictionary.

Recently, corpora are used to build bilingual dictionaries for Modern Standard Arabic (van Mol, 2000; van Mol & Paulussen, 2001; Hoogland, 1996; Zemank, 2001). A corpus of 3 million words was constructed to develop a Dutch–Arabic/Arabic–Dutch dictionary (van Mol, 2000). Older version of the Dutch–Arabic/Arabic–Dutch dictionaries were constructed using a 2 million-word corpus (Hoogland, 1996). Likewise, a corpus of 50 million words was used to develop an Arabic–Czech dictionary (Zemank, 2001).

Moreover, The Oxford Arabic Corpus (OAC) consists of 880 million words. It is used to construct the Oxford Arabic Dictionary (OAD) (Arts & McNeil 2013, Arts 2014). They used the Oxford Arabic Corpus with Sketch Engine to eliminate unnatural examples, to add appropriate examples showing natural usage, to identify modern senses of old words, and to include new vocabulary in the constructed dictionary.

Furthermore, a jellyfish dictionary for Arabic was developed using a large-scale Arabic corpus. A jellyfish dictionary is defined as a self-updating and automatically monitoring language change. Three motivations for developing a jellyfish dictionary for Arabic are (i) discovering new words; (ii) flagging obsolete words; and (iii) recognizing new senses. The large-scale Arabic corpus used to develop the jellyfish dictionary is consisted of 1 billion words (Attia & van Genabith, 2013).

In addition, many Arabic lexical databases were constructed. The morphological Analyzer for Arabic (BAMA) (Buckwalter 2002; 2004) contains Arabic-English lexicon files. One of them is contains 82 185 stems which was reused in many Arabic NLP tools such as morphological analyzers and spell checkers. Similarly, AyaSpell a spell checker for Arabic depends on a lexicon which was built by analyzing 5 traditional Arabic lexicons. It contains more than 50 000 entries distributed on more than 10 000 verbs and more than 40 000 nouns, particles and residuals (Zarrouki & Kebdani 2009; Zarrouki & Balla 2009). A third example is the Arabic WordNet (AWN). It is a lexical resource for MSA which is based on the design and the contents of the Princeton WordNet (PWN) for English. The semantic background for the AWN were encoded in a large ontology that contains around 1 000 terms and 4 000 definition statements (Elkateb & Black 2001; Black & El-Kateb 2004; Elkateb, Black et al. 2004; Rodríguez, Farwell et al. 2008). Likewise, Arabic Verbnet is a large lexicon for Arabic verbs. It contains verb entries where each entry is a third person masculine singular perfect verb. It has 173 classes which contain 4 392 verbs and 498 frames (Mousser 2010). Aralex is a lexical database which was developed to study the cognitive processing of Arabic on relation of precise frequency counts. Aralex was built depending on a 40-million word MSA corpus which was collected from online newspapers. It provides information about orthographic forms, stems, roots and patterns and their frequencies (Boudelaa & Wilson 2010). Quranic Arabic WordNet (QAWN) is a word net for the Qur'an and consists of 6 918 synsets that were constructed from about 8 400 unique word senses, on average of 5 senses for each word (Al-maayah et al 2015). These lexical databases are designed and built for a specific purpose and for specific Arabic NLP application. They are small in size and they are designed for MSA only (Sawalha 2011).

This paper describes an important lexical resource that is constructed to improve Arabic lexicography and Arabic NLP tools. The Corpus of Traditional Arabic Lexicons (the TAL-Corpus) is constructed from the text of 23 traditional Arabic lexicons spanning the period of over 1 200 years. The TAL-Corpus will be used as part of a large lexicographic corpus of Arabic to build new modern Arabic dictionaries. The TAL-Corpus can also be used to study the evolution of Arabic vocabulary system. The TAL-Corpus is accessible via an online interface which allows users to search for lexical entries.

2. Traditional Arabic Lexicons and Lexicography

Arabic lexicography is a well-established and deep-rooted art of Arabic literature. Arabic lexicography was founded by *al-farāhīdī* (died in 791) who constructed the first Arabic lexicon *kitāb al-‘ayn* ‘al-‘ayn lexicon’. Over the past 1 400 years, many Arabic lexicons were constructed. The lexical entries (*i.e.* roots) appear in Arabic dictionaries and followed by a definition part which may span several pages. The definition part is written as a unit or an encyclopaedic article which defines all the derived words from a certain root. These lexical entries are not arranged or distinguished with special formatting. Figure [1] shows samples of a lexical entry (the root **ك ت ب** k-t-b) with the definition part from the traditional Arabic dictionary “*lisān al-‘arab*”. Figure [2] is the human English translation of the lexical entry sample listed in the first figure. The derived words in both figures are manually underlined and highlighted in blue.

ك ت ب: الكتاب: معروف، والجمع كُتُبٌ وكُتُبٌ. كَتَبَ الشيءَ يَكْتُبُهُ كُتُبًا وكِتَابًا، وكُتِبَ: خَطَّهُ؛ قال أبو النجم: أَقْبَلْتُ من عِنْدِ زِيَادٍ كَالْخَرْفِ، تَخَطُّ رَجُلَانِ بِخَطِّ مُخْتَلَفٍ، كُتِبَانِ في الطَّرِيقِ لَمْ أَلْفِ قَالَ: ورَأَيْتُ في بعض النسخ يَكْتُبَانِ، بكسر التاء، وهي لغة بَهْرَاءَ، يَكْتُبُونَ التاءَ، فيقولون: يَغْلُمُونَ، ثم أُنْبِغَ الكاف كسرة التاء. وَالْكِتَابُ أيضًا: الاسمُ، عن اللحياني. الأزهرى: الكتابُ اسم لما كُتِبَ مَجْمُوعًا، والكتابُ مصدر؛ والكتابةُ لِمَنْ تَكُونُ لَهُ صِنَاعَةٌ، مثل الصَّيَاغَةِ والخِيطَةِ. والْكُتْبَةُ: اِكْتَبْتُكَ كِتَابًا تتسخه. ويقال: اِكْتَبْتُ فلانًا أي سألته أن يَكْتُبَ له كتابًا في حاجة. واِسْتَكْتَبَهُ الشيءَ أي سألَهُ أَنْ يَكْتُبَهُ له. ابن سيده: اِكْتَبْتُهُ كُتْبَةً. وقيل: كُتِبَ خَطَّهُ؛ واِكْتَبْتُهُ: اِسْتَمْلَاهُ، وكذلك اِسْتَكْتَبْتُهُ. واِكْتَبْتُهُ: كُتِبَ، واِكْتَبْتُهُ: كُتِبَ. وفي التنزيل العزيز: اِكْتَبْهَا فهي تُمْلَى عليه بُكْرَةً وَأَصِيلًا؛ أي اِسْتَكْتَبْتُهَا. ويقال: اِكْتَبْتُ الرجلُ إذا كُتِبَ نفسه في ديوان السُّلْطَانِ. وفي الحديث: قال له رجلٌ إنَّ امرأتِي خَرَجَتْ حَاجَةً، وإنِّي اِكْتَبْتُ في غزوة كذا وكذا؛ أي كُتِبْتُ اسمي في جملة الغزاة. وتقول: اِكْتَبْنِي هذه القصيدة أي أَمْلِهَا عَلَيَّ. والكتابُ: ما كُتِبَ فيه. وفي الحديث: مَنْ نَظَرَ في كِتَابٍ أخيه بغير إذنه، فكأنما يَنْظُرُ في النار؛ قال ابن الأثير: هذا تمثيل، أي كما يَحْذَرُ النارَ، فَلْيَحْذَرْ هذا الصنيعَ، قال: وقيل معناه كأنما يَنْظُرُ إلى ما يوجبُ عليه النارَ؛ قال: ويحتمل أنه أرادَ عَقُوبَةَ البَصَرِ لأنَّ الجناية منه، كما يُعاقَبُ السَّمْعُ إذا اسْتَمَعَ إلى قوم، وهم له كارهون؛ قال: وهذا الحديث محمولٌ على الكتابِ الذي فيه سِرٌّ وأمانة، يَحْزَرُه صاحبه أن يُطْلَعَ عليه؛ وقيل: هو عامٌّ في كل كتابٍ.

Figure 1. A sample of text from the traditional Arabic dictionary “*lisān al-‘arab*” for the lexical entry (ك ت ب k-t-b) where the derived words of the root (k-t-b) are underlined and highlighted in blue

Four main classes of ordering lexical entries in lexicons were developed and followed by authors of Arabic lexicons. Three arrangement methodologies depend on the roots of the words as lexical entries for Arabic lexicons. The fourth one groups lexical entries according to their conceptual themes or topical frames. These arrangement methodologies are different than those used in modern English dictionaries. Lexical entries of common English dictionaries, which are words (*i.e.* lexical entries in form of lemmas), are arranged alphabetically followed by the

type (*i.e.* part of speech) and the meaning of that word. On the other hands, Arabic lexicons depend on roots as lexical entries.

The first arrangement methodology of lexical entries of Arabic lexicons is the *al-ḥalīl* methodology. It was developed by *الخليل بن أحمد الفراهيدي al-ḥalīl bin aḥmad al-farāhīdī* (died in 791). The second arrangement methodology is the *al-ḡawharī* methodology which was developed by *'ismā'īl bin ḥammād al-ḡawharī* (died in 1002). The *al-barmakī* methodology is the third arrangement methodology. This arranging method was developed by *أبو المعالي محمد بن تميم البرمكي abū al-ma'ālī moḥammad bin tamīm al-barmakī*, who lived in the same time period as *al-ḡawharī*. *al-barmakī* did not construct a new lexicon; but he alphabetically re-arranged a lexicon called *aṣ-ṣiḥāḥ fī al-luḡa* 'الصحيح في اللغة' 'The Correct Language' by *al-ḡawharī*. For these three ordering methods, roots are considered the lexical entries. The last methodology is the *abū 'ubayd* methodology which was developed by *أبو عبيد القاسم بن سلام abū 'ubayd al-qāsim bin sallām* (died in 838). The following sections discuss the arrangement methodologies for lexical entries of traditional Arabic dictionaries.

k t b: [*al-kitāb*] the book; is well known. The plural forms are [*kutub^{un}*] and [*kutb^{un}*]. [*kataba Aṣ-ṣḥay'*] He wrote something. [*yaktubuhu*] the action of writing something. [*katb^{an}*], [*kitāb^{an}*] and [*kitābat^{an}*] means the art of writing. And [*kattabahu*] writing it means draw it up. Abu Al-Najim said: I returned back from Ziyad's house [after meeting him] and behaved demented, my legs drawn up differently (means walking in a different way). They wrote [*tukattibāni*] on the road the letters of *Lam Alif* (describing how he was walking crazily and in a different way). He said: I saw in a different version, the word "they wrote" [*ukittibāni*] using the short vowel *kasrah* on the first letter [*tā'*], as it is used by *Bahrā'* [Arab tribe] dialect. They say: (*ti'lamuwn*) (you know). Then the short vowel *kasrah* is propagated to the following letter (*kāf*). Moreover, [*al-kitāb*] the book is a noun. *Al-liḥyānī Al-'zharī* definition is: [*al-kitāb*] The book is the name of a collection of what has been written (a collection of written materials or texts). And the book has gerund [*al-kitābatu*] writing (art of writing) for whoever has a profession, similar to drafting and sewing. And [*al-kitābatu*]: is copying a book (copying a book in several copies). It is said: [*iktataba*] someone subscribed another means; he asked to write him a letter in something. [*istaktabahu*] He dictated someone something means to write him something. Ibn Sayyedah: [*iktatabahu*] is similar to [*katabahu*]. It is said: [*katabahu*] write something down means draw up. And [*iktatabahu*] writing something down means dictate someone something, which is the same meaning of [*istaktabahu*]. [*iktatabahu*] registering (masculine), and [*iktatabathu*] registering (feminine). In the Qur'an: [*iktatabaha*] He registered it, he has dictated it every sunrise and sunset, which means dictating it. It is said: [*iktataba ar-rajul*] The man registered, if he registered himself in the Sultan's office. In Hadith: a man said to him (the prophet): my wife is pilgrimaging (to Mecca), and I have registered [*uktutibtu*] in a conquest, which means that I have written my name among the conquerors. And you say: [*'aktibnī*] let me copy this poem, means dictate me the poem. Also, [*al-kitāb*] the book is something which has been written on. And in Hadith: who looks at his brother's book without permission is as looking to hell. Ibn Al-Atheer said: it is a similarity; which means as he avoids hell, he should avoid doing this. He said: the meaning (of the Hadith) is the punishment by hell will be applied if someone looks at a book without permission. He said: it might be the punishment of visual explorers as the crime is done by sight. Hearing explorer is punished if someone intentionally listened to other people who do not like anyone to listen to them. He said: this Hadith is specific for books of secrets and secure books, whose owners hate anybody to look at these books. It is also said: the Hadith is general; applied to any type of books [*kitāb*].

Figure 2. A human translation of the sample of text from the traditional Arabic lexicon "*lisān al-'arab*", the target lexical entries are highlighted in blue and square brackets

2.1 The *al-ḥalīl* Ordering Methodology

The first traditional Arabic lexicon is called *كتاب العين kitāb al-'ayn* "al-'ayn lexicon". It was developed by *الخليل بن أحمد الفراهيدي al-ḥalīl bin aḥmad al-farāhīdī* (died in 791). The *al-ḥalīl* ordering methodology, which was followed in constructing 'The al-'ayn' lexicon, arranges the lexical entries phonologically according to places of articulation of phonemes from the mouth and throat, working forwards from glottal through to labial regions. The *al-'ayn* lexicon was divided into books, where one book was dedicated for each letter. Each book was then divided into 4 sections according to their internal structure: (i) doubled biliteral roots; (ii) intact trilateral roots; (iii) doubly-defective roots; and (iv) quadrilateral and quinquiliteral roots. Many lexicons followed *al-ḥalīl*'s methodology with slight modifications. Table [1] lists some of traditional Arabic Lexicons that followed *al-ḥalīl*'s methodology.

2.2 The *al-ğawharī* Methodology

'ismā'īl bin ḥammād *al-ğawharī* اسماعيل بن حماد الجوهري (died in 1002) constructed a lexicon called *الصحيح في اللغة* *aṣ-ṣiḥāḥ fī al-luġa^h* 'The Correct Language'. Roots are the lexical entries of this lexicon. They were alphabetically ordered according to their last letter, then the first letter. This methodology is called the *al-ğawharī* methodology. The lexicon was organized into chapters where each chapter corresponds to the last letter of the root. Each chapter includes sections corresponding to the first letter of the root, then the second letter of trilateral roots, then the third letter of quadrilateral roots, then the fourth letter in quinquilateral roots. For example, the word *بَسَطَ* *baṣaṭa* "spread" which is derived from the root (b-s-t) is found in chapter ط *t* representing the last letter of the root, and in section ب *b* representing the first letter of the root. Table [1] lists some of traditional Arabic Lexicons that followed this ordering methodology.

2.3 The *al-barmakī* Methodology

The third lexicon ordering methodology is "The *al-barmakī* methodology". It was developed by *abū al-ma'ālī muḥammad bin tamīm al-barmakī* أبو المعالي محمد بن تميم البرمكي (died in 1006). In this methodology, lexical entries (i.e. roots) are alphabetically arranged according to the first letter of the root. *al-barmakī* lived in the same period as *al-ğawharī*. *al-barmakī* did not construct a new Arabic lexicon. Instead, he re-arranged the lexical entries of *الصحيح في اللغة* *aṣ-ṣiḥāḥ fī al-luġa^h*, which was developed by *al-ğawharī*. The *al-barmakī* methodology was followed by *az-zamahṣarī* (died in 1143) in constructing his lexicon *أساس البلاغة* *asās al-balāġa^h* "Fundamentals of Rhetoric". Table [1] lists Arabic lexicon which followed the *al-barmakī* methodology for ordering lexical entries. The *al-barmakī* methodology for ordering lexical entries becomes the most widely used ordering methodology for Arabic lexicons.

2.4 The *abū 'ubayd* Methodology

abū 'ubayd al-qāsim bin sallām أبو غبيد القاسم بن سالم (died in 838) developed the fourth ordering methodology for Arabic lexicons which is called "The *abū 'ubayd* methodology". This methodology arranges and groups together lexical entries according to their semantic fields. This arrangement methodology is similar to arranging lexical entries in modern thesauri. Many lexicons followed this ordering methodology. *الغريب المصنف في اللغة* *al-ġarīb al-muṣannaf fī al-luġa^h* "The Irregular Classified Language" by *abū 'ubayd al-qāsim bin sallām* was the first lexicon that followed this methodology. This lexicon includes many small books that describe similar topics (i.e. group words of similar meanings) such as books describing horses, milk, honey, flies, insects, palms, and human creation. Then, more than thirty small books were collated into one large lexicon. Figure [3] shows a sample from Colours' Book taken from *al-ġarīb al-muṣannaf fī al-luġa^h* lexicon. Table [1] lists traditional Arabic lexicons that followed *abū 'ubayd* methodology.

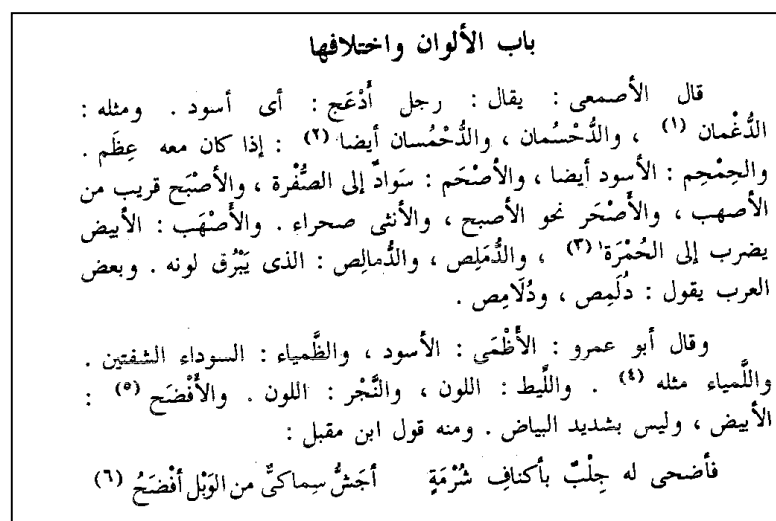


Figure 3. A sample of *الغريب المصنف في اللغة* *al-ġarīb al-muṣannaf fī al-luġa^h* "The Irregular Classified Language" lexicon

Table 1. Examples of Traditional Arabic Lexicons classified according to their Arrangement Methodology

Arrangement Methodology	Traditional Arabic Lexicons following this Arrangement Methodology
1. The <i>al-ḥalīl</i> Methodology	1- كتاب العين <i>kitābu al-‘ayn</i> “al-‘ayn Lexicon” by الخليل ابن أحمد الفراهيدي <i>al-ḥalīl bin aḥmad al-farāhīdī</i> (died in 175H / 791AD).
	2- مُعْجَمُ الْمُحِيطِ فِي اللُّغَةِ <i>mu‘ğam al-muḥīt fī al- luğah</i> “The Comprehensive Language” by الصاحب بن عباد <i>aṣ-ṣāhib bin ‘abbād</i> (died in 385H / 995AD).
	3- المحكم والمحيط الأعظم <i>al-muḥkam wa al-muḥīt al-‘a‘āzam</i> “The Greatest Verified and Comprehensive Lexicon” by ابن سيده أبو الحسن علي بن اسماعيل النحوي الأندلسي <i>‘ibn sayyidah, abū al-ḥasan bin ‘ismā ‘l an-naḥawī al-lağawī al-‘andalusī</i> (died in 458H / 1065AD).
	4- لسان العرب <i>lisān al-‘rab</i> “Arab tongue” by جمال الدين محمد بن منظور <i>ğamāl ad-dīn moḥammed bin manzūr</i> (died in 629H / 1311AD).
	5- معجم تهذيب اللغة <i>mu‘ğam taḥḍīb al-luğah</i> “The Lexicon of Refined Language” by أبو منصور إسماعيل بن حماد <i>abū manṣūr al-‘azharī</i> (died in 1205H / 1790AD).
2. The <i>al-ğawharī</i> Methodology	1- الصحاح في اللغة <i>aṣ-ṣiḥāḥ fī al-luğah</i> “The correct language” by أبو نصر إسماعيل بن حماد <i>abū naṣr ‘ismā‘l bin ḥammād al-ğawharī al-farābī</i> (died in 400H / 1009AD).
	2- العباب الزاخر في اللغة <i>al-‘ibāb az-zāḥir fī al-luğah</i> “The High Flood Water of Language” by الحسن بن محمد الصنعاني <i>al-ḥasan bin muḥammad aṣ-ṣağānī</i> (died in 650H / 1252AD).
	3- القاموس المحيط <i>al-qāmūs al-muḥīt</i> “The Comprehensive Dictionary” by مجد الدين أبو طاهر <i>mağd ad-dīn abū ṭāhir muḥammad bin ya‘qūb al-fayrūz ‘ābādī</i> (died in 817H / 1414AD).
	4- تاج العروس من جواهر القاموس <i>tag al-‘arūs min ġawāhir al-qāmūs</i> “Bridal Crown Jewel of Dictionaries” by الزبيدي <i>az-zubaydī</i> (died in 1205H / 1790AD).
3- The <i>al-barmakī</i> methodology	1- معجم الجيم <i>mu‘ğam al-ğīm</i> “The jīm Lexicon” by أبو عمرو الشيباني <i>abū ‘amr aṣ-ṣībānī</i> (died in 206H / 821AD).
	2- جمهرة اللغة <i>ğamharat al-luğah</i> “The Gathering of the Language” by ابن دُرَيْدٍ <i>‘ibn durayd</i> (died in 256H / 869AD).
	3- معجم مقاييس اللغة <i>mu‘ğam maqāyīs al-luğah</i> “The Lexicon of the Standard Language” by أبي الحسين أحمد بن فارس بن زكريا <i>‘abī al-ḥusayn aḥmad bin fāris bin zakariyyā</i> (died in 395H / 1004AD).
	4- معجم ما استعجم <i>mu‘ğam mā ‘ista‘ğam</i> “A Lexicon of Foreign Words” by البكري الأندلسي <i>al-bakrī al-‘andalusī</i> (died in 487H / 1094AD).
	5- تهذيب الأفعال <i>taḥḍīb al-af‘āl</i> “The Refined Verbs” by أبو القاسم علي بن جعفر السعدي <i>(‘ibn al-qīṭā) abū al-qāsim ‘alī bin ġa‘far as-sa‘dī</i> (died in 515H / 1121AD).
	6- المصباح المنير في غريب الشرح الكبير <i>al-muṣbāḥ al-munīr fī ġarīb aṣ-ṣarḥ al-kabīr</i> “The Illuminating Light on the Irregularity of the Great Explanations” by أحمد بن محمد بن علي <i>aḥmad bin muḥammad ‘alī al-fayyūmī ṭumma al-ḥamawī, abū al-‘abbās</i> (died in 538H / 1143AD).
	7- أساس البلاغة <i>asās al-balāğah</i> “Fundamentals of Rhetoric” by أبو القاسم محمود بن عمرو بن <i>abū al-qāsim maḥmūd bin ‘amr bin aḥmad, az-zamaḥṣarī ġār allāh</i> (died in 538H / 1143 AD).
	8- الْمُغْرِبُ فِي تَرْتِيبِ الْمُغْرِبِ <i>al-muğrib fī tartīb al-mu‘rib</i> “Irregular Declinable Words” by أبو الفتح ناصر الدين المطرزي <i>‘abū al-faṭḥ nāṣir ad-dīn al-muṭrazī</i> (died in 610H / 1213AD).
	9- مختار الصحاح <i>muḥtār aṣ-ṣiḥāḥ</i> “The Selected of the Correct Language” by أبو بكر الرازي <i>abū bakr ar-rāzī</i> (died in 666H / 1267AD).
	10- المعجم الوسيط <i>al-mu‘ğam al-wasīṭ</i> “The Intermediary Lexicon” by إبراهيم مصطفى - أحمد - حامد عبد القادر - محمد النجار <i>ibrāhīm muṣṭafā, aḥmad az-zayyāt, ḥāmid ‘abdul-qādir, muḥammad an-nağğār</i> (published in 1960).
	11- معجم الأفعال المتعدية بحرف <i>mu‘ğam al-af‘āl al-muta‘adyyah bi ḥarf</i> “The Lexicon of Transitive Verbs” by موسى بن محمد بن الملياني الأحمد <i>mūsā bin muḥammad al-malyānī al-‘aḥmadī</i> (published in 1979).

4- The <i>abū 'ubayd</i> Methodology	1- <i>al-ġarīb al-muṣannaf fī al-luġa</i> "The Irregular Classified Language" by أبي غيبه القاسم بن سلام <i>'abi 'ubayd al-qāsim bin sallām</i> (died in 223H / 838AD).
	2- <i>al-munaḡḡad fī al-luġa</i> "The Decorated Language" by علي بن حسن الهنائي <i>ali bin ḥasan al-hunā'ī al-'azdī</i> (died in 310H / 922AD).
	3- <i>al-muḥaṣṣaṣ fī al-luġa</i> "The Specified Language" by أبو الحسن (ابن سيدة) <i>'ibn ṣayyidah, abū al-ḥasan bin 'ismā'īl an-naḥwī al-luġawī al-'andalusī</i> (died in 458H / 1065AD).

3. The Design of the TAL-Corpus

The motivation behind constructing the TAL-Corpus is to collect and organize well-established and long traditions of traditional Arabic lexicon in one freely available resource. The TAL-Corpus will help Arabic lexicographers to design and construct new modern Arabic dictionaries. These dictionaries can have new ordering methodology where derived words can be easily linked with their lexical entries whether they are roots or lemmas. The TAL-Corpus can be used to determine the origin of Arabic vocabulary and can track the development and changes of their meanings. The TAL-Corpus can also be used to extract useful information that supports Arabic NLP applications such as root extraction applications, morphological analyzers, semantic networks of Arabic vocabulary, WordNets, ontologies ... etc.

The following sections show the design criteria followed in constructing the TAL-Corpus. Atkins et al., (1992) proposed general criteria of corpus design. These principal aspects and standards were recommended to be followed to inform major decisions for corpus creation. These criteria were designed to support high-quality and compatible corpora regardless of the corpus language, purpose, and location. Sections 3.1 to 3.5 discuss the design criteria followed to construct the TAL-Corpus.

3.1 Text

The text of the TAL-Corpus was collected from 23 freely available traditional Arabic lexicons. These lexicons are listed in Table 1. Al-Meshkat Islamic Network¹ *شبكة مشكاة الإسلامية* *šabaka' miškā' al-'islāmiyya* provides most of these lexicons freely. These lexicons have been key-boarded (*i.e.* typed) and put online in machine readable formats as MS-Word (.doc) or HTML text files.

The texts of the collected Arabic dictionaries were organized using different ordering methodologies as discussed in Section 2. However, most of these lexicons use roots as their main lexical entries. The definition of a root in each lexicon is written as an encyclopaedic article that contains the derived words from that root, their meanings, and examples of usages. These definitions vary in size from half a page to span several pages. Figure [4] shows a sample of text of a lexical entry taken from a traditional Arabic lexicon; the derived words are underlined and highlighted in blue. The text of the collected lexicons is fully vowelized, partially vowelized or non-vowelized. Texts (*i.e.* definitions) of similar roots from the different traditional Arabic dictionaries were grouped together in the TAL-Corpus. Then, several automatic processing steps and algorithms were applied to extract relevant linguistic information such as derived words and lemmas. Sections 3.4 and 3.5 discusses in detail these processing steps and algorithms.

(ك ت ب):
 (كُتِبَ) كُتِبَ وَكُتِبَ وَقَوْلُهُ وَإِذَا كَانَتْ السَّرْفَةُ صُحُفًا لَيْسَ فِيهَا كِتَابٌ أَيِ مَكْتُوبٌ (وَفِي حَدِيثِ أَنَسٍ) وَأَحْكُمُ بِكِتَابِ اللَّهِ أَيِ بِمَا فَرَضَ اللَّهُ مِنْ كِتَابٍ عَلَيْهِ كَذَا إِذَا أُوجِبَ وَفَرَضَ (وَمِنْهُ) الصَّلَاةُ الْمَكْتُوبَةُ وَأَمَّا قَوْلُهُ - صَلَّى اللَّهُ عَلَيْهِ وَآلِهِ وَسَلَّمَ - [مَا بَالُ أَقْوَامٍ يَشْتَرِطُونَ شَرْطًا لَيْسَتْ فِي كِتَابِ اللَّهِ تَعَالَى] فَقِيلَ الْمُرَادُ قَوْلُهُ تَعَالَى (أَذْعُوهُمْ لِأَبَائِهِمْ) إِلَى أَنْ قَالَ وَمَوَالِيَهُمْ أَنَّهُ نَسَبُهُمْ إِلَى مَوَالِيَهُمْ كَمَا نَسَبَهُمْ إِلَى آبَائِهِمْ فَلَمَّا لَمْ يَجْزِ التَّحُولُ عَنْ الْأَبَاءِ لَمْ يَجْزِ عَنْ الْأَوْلِيَاءِ وَيَجُوزُ أَنْ يُزَادَ بِكِتَابِ اللَّهِ قَضَاؤُهُ وَحُكْمُهُ عَلَى لِسَانِ رَسُولِ اللَّهِ - صَلَّى اللَّهُ عَلَيْهِ وَآلِهِ وَسَلَّمَ - إِنَّ الْوَلَاءَ لِمَنْ أَعْتَقَ (وَأَكْتُبَ الْغُلَامَ وَكُتِبَ) عِلْمُهُ الْكِتَابُ (وَمِنْهُ) سَلَّمَ غُلَامُهُ إِلَى مَكْتُوبٍ أَيِ إِلَى مُعَلِّمِ الْخَطِّ رُويَ بِالتَّخْفِيفِ وَالتَّشْدِيدِ (وَأَمَّا الْمَكْتُوبُ) وَالْكِتَابُ فَمَكَانُ التَّعْلِيمِ وَقِيلَ الْكِتَابُ الصَّبِيحَانِ (وَكُتِبَ) عَزَّدَهُ مُكَاتِبَةً وَكِتَابًا قَالَ لَهُ حَزْرَتُكَ يَدَا فِي الْحَالِ وَرَقَبَةٍ عِنْدَ آدَاءِ الْمَالِ (وَمِنْهُ) قَوْلُهُ تَعَالَى (وَالَّذِينَ يَبْتِغُونَ الْكِتَابَ) وَقَدْ يُسَمَّى بِذَلِكَ الْكِتَابَةُ مُكَاتِبَةً وَأَمَّا الْكِتَابَةُ فِي مَعْنَاهَا فَلَمْ أَجِدْهَا إِلَّا فِي الْأَسَاسِ وَكَذَا تَكَاتِبَ الْعَبْدُ إِذَا صَارَ مُكَاتِبًا وَمَذَارُ التَّزْكِيكِ عَلَى الْجَمْعِ (وَمِنْهُ كِتَابُ النَّعْلِ وَالْقَرْبَةِ) حَزَرَهَا (وَالْكِتَابُ) الْخَزْرُ الْوَاجِدَةُ كُتِبَ (وَمِنْهُ كِتَابُ الْبَغْلَةِ) إِذَا جَمَعَ بَيْنَ شَفَرَتَيْهَا بِحَلْقَةٍ (وَالْكِتَابَةُ) الطَّائِفَةُ مِنَ الْجَيْشِ مُجْتَمِعَةً (وَبِهَا سُمِّيَ) أَحَدُ خُصُوفِ خَيْبَرَ (وَقَوْلُهُمْ) سَمِيَ هَذَا الْعَقْدُ مُكَاتِبَةً لِأَنَّهُ ضَمُّ حَرِيَّةِ الْيَدِ إِلَى حَرِيَّةِ الرَّقَبَةِ أَوْ لِأَنَّهُ جَمَعَ بَيْنَ تَجْمِينِ فَصَاعِدًا ضَعِيفًا جَدًّا وَإِنَّمَا الصَّوَابُ أَنْ كَلَّمَ مِنْهُمَا كِتَابٌ عَلَى نَفْسِهِ أَمْرًا هَذَا الْوَقَاءَ وَهَذَا الْإِدَاءَ.

Figure 4. A sample of text from the traditional Arabic dictionary "*al-muġrib fī tartīb al-mu'rib*", the derived word of the target lexical entry are underlined and highlighted in blue.

For all collected lexicons, common processing steps were applied. These steps include; (i) converting the file

formats from MS Word or HTML web pages into standard text files in Unicode ‘utf-8’ encoding. (ii) A statistical analysis was applied that computed the words frequencies and the vocabulary size for both vowelized and non-vowelized text of the corpus. As a result, the complete TAL-Corpus contains 14 369 570 words, 2 184 315 vowelized word types and 569 412 non-vowelized word types (*i.e.* after removing short vowels (diacritics) from the text). Table [2] shows the summary of the statistical analyses of the lexicon texts used to construct the TAL-Corpus. Figure [5] shows the highest 25 frequent words in the TAL-Corpus of partially vowelized and non-vowelized forms of words.

Table 2. Statistical Analysis of the Lexicons’ Text used to construct the TAL-Corpus

Number of files		247
Size		178.32 MB
Vowelized word	Number of words	14 369 570
	Number of word types	2 184 315
Non-vowelized word	Number of words	14 369 570
	Number of word types	569 412

Partially-vowelized			Non-vowelized		
Word		Frequency	Word		Frequency
في	<i>fī</i> “in”	292 396	من	<i>min</i> “from”	322 239
من	<i>min</i> “from”	269 200	في	<i>fī</i> “in”	301 895
قال	<i>qāl</i> “he said”	172 631	قال	<i>qāl</i> “he said”	190 918
و	<i>wa</i> “and”	120 060	أي	<i>’ayy</i> “which”	132 635
على	<i>’alā</i> “over”	108 252	و	<i>wa</i> “and”	130 809
ما	<i>mā</i> “what”	89 195	على	<i>’alā</i> “over”	119 639
وقال	<i>wa qāl</i> “and he said”	88 233	إذا	<i>’ihā</i> “if”	115 842
عن	<i>’an</i> “about”	82 027	وقال	<i>wa qāl</i> “and he said”	99 601
إذا	<i>’ihā</i> “if”	81 479	ابن	<i>’ibn</i> “son of”	94 980
أي	<i>’ay</i> “which”	78 622	ما	<i>mā</i> “what”	94 530
وهو	<i>wa huwa</i> “and he”	75 149	بن	<i>bin</i> “son of”	92 213
لا	<i>lā</i> “no”	69 737	عن	<i>’an</i> “about”	87 064
ابن	<i>’ibn</i> “son of”	58 334	وهو	<i>wa huwa</i> “and he”	80 375
به	<i>bihi</i> “in it”	53 343	لا	<i>lā</i> “no”	73 066
وفي	<i>wa fī</i> “and in”	53 197	أبو	<i>abū</i> “father”	72 231
وقد	<i>wa qad</i> “and perhaps”	50 648	أن	<i>’an</i> “that”	65 419
أبو	<i>abū</i> “father”	47 915	أو	<i>’aw</i> “or”	62 298
بن	<i>bin</i> “son of”	46 880	الله	<i>allā^h</i> “Allah”	59 511
أي	<i>’ay</i> “which”	46 788	به	<i>bihi</i> “in it”	58 941
هو	<i>huwa</i> “he”	45 916	يقال	<i>yuqāl</i> “it is said”	58 062
يقال	<i>yuqāl</i> “it is said”	45 794	وفي	<i>wa fī</i> “and in”	55 077
عليه	<i>’alayhi</i> “about him”	44 786	وقد	<i>wa qad</i> “and perhaps”	53 992
ولا	<i>wa lā</i> “and not”	42 190	عليه	<i>’alayhi</i> “about him”	50 906
الله	<i>allā^h</i> “Allah”	39 961	هو	<i>huwa</i> “he”	49 785

Figure 5. The first 25 words of the frequency list generated from the TAL-Corpus Corpus

The analysis represented by Tables [3] and [4] and Figure [6] classifies the traditional Arabic lexicons which were include in the TAL-Coprus, according to the time of construction. The time period spans around 14 centuries since the first Arabic lexicon was created (*i.e.* from the second Hijri century to the fifteenth Hijri century). This time span was divided into 14 time frames where each corresponds to 100 years. These time frames were defined by the creation times of the traditional Arabic dictionaries which are indicated by the death date of dictionaries' authors. The first time frame includes one lexicon *kitābu al-‘ayn* which consists of 348 114 words and 141 098 word types which forms 2.42% of the text size and 3.72% of the vocabulary size of the TAL-Corpus. The lexicons from 12th century are the largest. They contain 5 215 917 words and 1 211 432 word types. They form 36.30% of the TAL-Corpus text and 31.90% of its vocabulary size. The lexicons included in this time frame are *tag al-‘arūs min ġawāhir al-qāmūs* and *mu‘ġam tahqīb al-luġa^h* which represent the largest in terms of number of words and vocabulary size.

Table 3. Text and vocabulary size of the Traditional Arabic Dictionaries and their percentage in the TAL-Corpus

	Time Frame	Lexicon Name	Date (Died in)	# Words	# Types	% of Words	% of Types
1	100-199H, 718-814AD	<i>kitābu al-‘ayn</i> كتاب العين	175H (791AD)	348,114	141,098	2.42%	3.72%
2	200-299H, 815-911AD	<i>mu‘ġam al-ġīm</i> معجم الجيم	206H (821AD)	125,676	56,274	0.87%	1.48%
3	300-399H, 912-1008AD	<i>al-ġarīb al-muṣannaf fī al-luġa^h</i> الغريب المصنف في اللغة	223H (838AD)	16,541	7,775	0.12%	0.20%
		<i>ġamharat al-luġa^h</i> جمهرة اللغة	256H (869AD)	396,144	123,576	2.76%	3.25%
		<i>al-munaġġad fī al-luġa^h</i> المنجد في اللغة	310H (922AD)	32,173	16,942	0.22%	0.45%
		<i>mu‘ġam al-muḥīt fī al-luġa^h</i> مغمم المحيط في اللغة	385H (995AD)	392,246	168,870	2.73%	4.45%
		<i>mu‘ġam maqāyīs al-luġa^h</i> معجم مقاييس اللغة	395H (1004AD)	445,126	129,838	3.10%	3.42%
4	400-499H, 1009-1105AD	<i>aṣ-ṣiḥāḥ fī al-luġa^h</i> الصاحح في اللغة	400H (1009AD)	593,654	118,591	4.13%	3.12%
		<i>al-muḥkam wa al-muḥīt al-‘a‘āz</i> المحكم والمحيط الأعظم	458H (1065AD)	1,020,137	279,157	7.10%	7.35%
		<i>al-muḥaṣṣaṣ fī al-luġa^h</i> المخصص في اللغة	458H (1065AD)	902,324	274,780	6.28%	7.24%
		<i>mu‘ġam mā ‘ista‘ġam</i> معجم ما استعجم	487H (1094AD)	278,713	43,289	1.94%	1.14%
5	500-599H, 1106-1202AD	<i>tahqīb al-af‘āl</i> تهذيب الأفعال	515H (1121AD)	132,319	38,102	0.92%	1.00%
6	600-699H, 1203-1299AD	<i>asās al-balāġa^h</i> أساس البلاغة	538H (1143AD)	289,436	95,887	2.01%	2.52%
		<i>al-muġrib fī tartīb al-mu‘rib</i> المغرب في ترتيب المغرب	610H (1213AD)	128,047	39,930	0.89%	1.05%
		<i>al-‘ibāb az-zāhir fī al-luġa^h</i> العباب الزاخر في اللغة	650H (1252AD)	261,658	100,536	1.82%	2.65%
		<i>muḥtār aṣ-ṣiḥāḥ</i> مختار الصحاح	666H (1267AD)	171,487	40,295	1.19%	1.06%
7	700-799H, 1300-1396AD	<i>lisān al-‘rab</i> لسان العرب	711H (1311AD)	2,146,545	507,860	14.94%	13.37%
		<i>al-muṣbāḥ al-munīr fī ġarīb aṣ-ṣarḥ al-kabīr</i> المصباح المنير في غريب الشرح الكبير	770H (1368 AD)	219,276	61,422	1.53%	1.62%
8	800-899H,	<i>al-qāmūs al-muḥīt</i>	817H	563,460	203,600	3.92%	5.36%

	1397-1493AD	القاموس المحيط (1414AD)					
12	1200-1299H, 1785-1881AD	<i>mu'ḡam tahdīb al-luḡa^h</i> معجم تهذيب اللغة <i>taḡ al-'arūs min ḡawāhir al-qāmūs</i> تاج العروس من جواهر القاموس	1205H (1790AD) 1205H (1790AD)	1,351,837 3,864,080	379,928 831,504	9.41% 26.89%	10.00% 21.89%
13	1300-1399H, 1882-1978AD	<i>al-mu'ḡam al-wasīṭ</i> المعجم الوسيط	Modern 1960	615,352	112,164	4.28%	2.95%
14	1400H-Today, 1979AD	<i>mu'ḡam al-'af'āl al-muta'diyah bi ḥarf</i> معجم الأفعال المتعدية بحرف	Modern 1979	75,225	26,299	0.52%	0.69%

Table 4. The 14 Time frames and their percentage of words and vocabulary size in the TAL-Corpus.

Frame	Time frame	# of dictionaries	# words	# types	% of words	% of types
1	100H-199H (718AD-814AD)	1	348,114	141,098	2.42%	3.72%
2	200H-299H (815AD-911AD)	3	538,361	187,625	3.75%	4.94%
3	300H-399H (912AD-1008AD)	3	869,545	315,650	6.05%	8.31%
4	400H-499H (1009AD-1105AD)	4	2,794,828	715,817	19.45%	18.85%
5	500H-599H (1106AD-1202AD)	1	132,319	38,102	0.92%	1.00%
6	600H-699H (1203AD-1299AD)	4	850,628	276,648	5.92%	7.28%
7	700H-799H (1300AD-1396AD)	2	2,365,821	569,282	16.46%	14.99%
8	800H-899H (1397AD-1493AD)	1	563,460	203,600	3.92%	5.36%
9	900H-999H (1494AD-1590AD)	0	-	-	-	-
10	1000H-1099H (1591AD-1687AD)	0	-	-	-	-
11	1100H-1199H, 1688AD-1784AD	0	-	-	-	-
12	1200H-1299H, 1785AD-1881AD	2	5,215,917	1,211,432	36.30%	31.90%
13	1300H-1399H, 1882AD-1978AD	1	615,352	112,164	4.28%	2.95%
14	1400H-Today, 1979AD	1	75,225	26,299	0.52%	0.69%

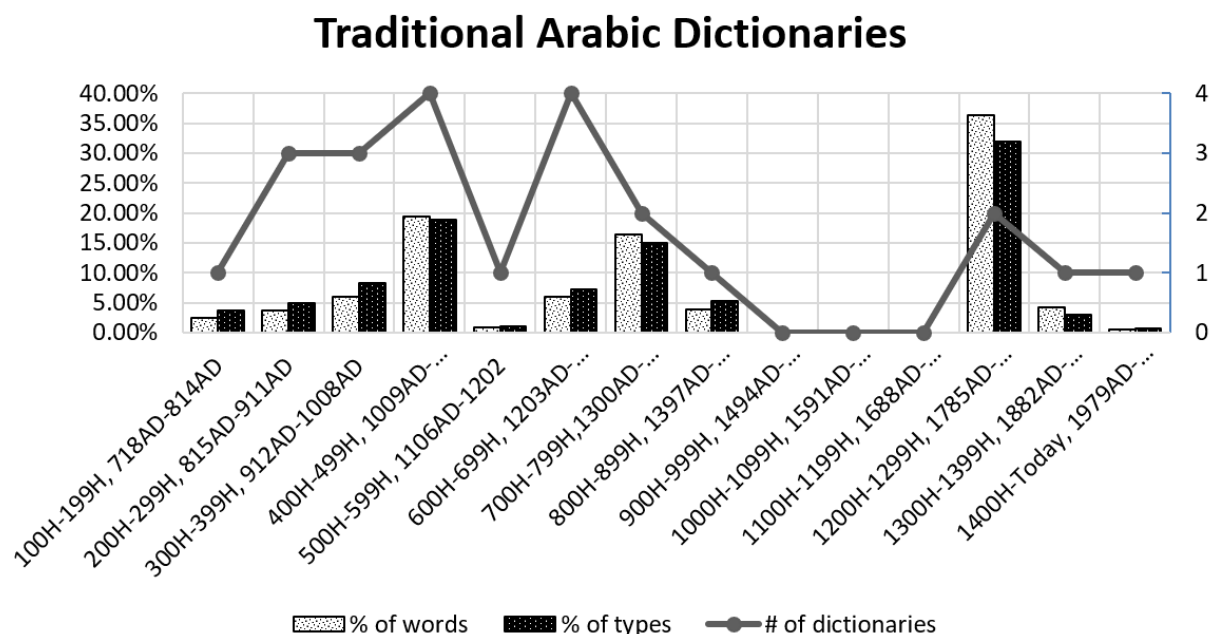


Figure 6. Traditional Arabic dictionaries included in the TAL-Corpus

3.2 Text Handling

After collecting the text of 23 traditional Arabic dictionaries, common pre-processing steps were applied. First, all dictionaries' files were converted into standard text files using Unicode 'utf-8' encoding. Then, the SALMA-Tokenizer and the SALMA-root extractor and Lemmatizer (Sawalha, 2011) were used to tokenize and process Arabic words by stripping diacritics, and extracting the root and the lemma for each word in the TAL-Corpus. Third, frequency lists of both vowelized and non-vowelized word were generated (see Table [1] and Figure [5]).

Special algorithm was developed to extract the derived words of the lexical entries for the dictionaries included in the TAL-Corpus. The purpose of this algorithm is to group together roots and their definition parts and then to extract derived words of roots from their related definition articles. To achieve this goal, a specific treatment were applied to each dictionary text. The 23 collected dictionaries were originally constructed following an ordering methodology of their lexical entries as discussed in Section 2. Most of them use roots as their main head words of lexical entries. These dictionaries were typed into machine-readable files in different formats without using any lexicographic representations that can be recognized by Computers. Therefore, specialized programs were developed for each dictionary to reformat and extract useful information such as roots, definitions and derived words.

The root-definition structure is the common basic structure for most traditional Arabic dictionaries. Each lexical entry consists of the root as a head word and the definition part. The definition part is written as an encyclopaedic article featuring free writing style. These encyclopaedic articles defines the root and its derived words and their linguistic attributes are specified. However, the derived words of a root within the definition part are neither structured nor ordered. This free writing style requires the authors of dictionaries to add affixes and clitics to the derived words within the definition parts. Clitics, such as conjunctions, prepositions and connected pronouns, are used to connect sentences and paragraphs of these definition articles.

For the above mentioned-reasons, the free writing style of the definition part adds extra challenges to extract the derived words and their definitions. Therefore, a dedicated algorithm was developed to extract the roots and their derived words from the dictionaries' texts. The tokenizing module in the program specifies the boundaries of a lexical entry which is normally starts with a root followed by an article that defines that root. For each lexical entry, the algorithm extracts and pairs words from the definition part with the root and stores them in vectors (*i.e.* bag of words). Many of these word-root pairs are not correct matches (*i.e.* the word is not derived from the associated root). A normalization analysis verified these word-root pairs by throwing out pairs where the word is not derived from its associated root. The normalization procedure applies linguistic knowledge that governs the derivation process of words from their roots. These linguistic rules were used to match the consonant letters of words and roots and their order for each word-root pair. The first linguistic rule checks if all consonant letters forming the root appear in the paired word. The second rule examines if all root letters orderly appear in the derived word. Both rules must be applied to every word-root pair for verification. This process is applied to extract the derived words of a root and later to build a morphological lexicon (See Section 3.3.1). Figure [7] shows the process of selecting word-root pairs. Table [5] shows the number of words and the percentage of words extracted from the original text of the dictionaries.

Word-root vector for the root كَتَبَ <i>k-t-b</i>				
(مُخْتَلَفٌ ، كَتَبَ)	(عُدَّ ، كَتَبَ)	(خَطَّه ، كَتَبَ)	(كَتَبَ ، الشَّيْءَ)	(الْكِتَابُ ، كَتَبَ)
(تَكْتَبَانِ ، كَتَبَ)	(زَادَ ، كَتَبَ)	(قَالَ ، كَتَبَ)	(كَتَبَ ، يَكْتُبُهُ)	(مَعْرُوفٌ ، كَتَبَ)
(فِي ، كَتَبَ)	(كَالْخَرْفِ ، كَتَبَ)	(أَبُو ، كَتَبَ)	(كَتَبَ ، كُنْأً)	(وَالْجَمْعُ ، كَتَبَ)
(الطَّرِيقِ ، كَتَبَ)	(تَخَطَّ ، كَتَبَ)	(النَّجْمِ ، كَتَبَ)	(وَكُنْأً ، كَتَبَ)	(كُنْأً ، كَتَبَ)
(لَامٌ ، كَتَبَ)	(رَجُلًا ، كَتَبَ)	(أَقْبَلْتُ ، كَتَبَ)	(وَكُنْأَةً ، كَتَبَ)	(كُنْأً ، كَتَبَ)
(أَلْفٌ ، كَتَبَ)	(بَحَطَّ ، كَتَبَ)	(مَنْ ، كَتَبَ)	(وَكُنْأَةً ، كَتَبَ)	(كُنْأً ، كَتَبَ)

Figure 7. Using linguistic knowledge to select word-root pairs from traditional Arabic lexicons. The selected word-root pairs are underlined and highlighted in blue

Table 5. Words and Roots Extracted from 8 Traditional Arabic lexicons

Lexicon name	Word types	Words extracted	Roots extracted
--------------	------------	-----------------	-----------------

1	<i>tağ al-‘arūs min ġawāhir al-qāmūs</i>	831,504	474,351	57.05%	11,101
2	<i>lisān al-‘rab</i>	507,860	274,305	54.01%	9,355
3	<i>mu‘ğam al-muḥīṭ fī al- luğā^h</i>	168,870	66,763	39.54%	6,411
4	<i>kitābu al-‘ayn</i>	141,098	54,970	38.96%	5,826
5	<i>al-mu‘ğam al-wasīṭ</i>	112,164	45,614	40.67%	6,489
6	<i>al-muṣbāḥ al-munīr fī ġarīb aš-šarḥ al-kabīr</i>	61,422	29,742	48.42%	2,947
7	<i>muḥtār aš-ṣiḥāḥ</i>	40,295	17,636	43.77%	3,420
8	<i>al-muğrab fī tartīb al-mu‘rab</i>	39,930	13,798	34.56%	2,322

3.3 Advanced Text Handling

The TAL-Corpus implements advanced text handling tools which can automatically process linguistic information in a corpus and allow more sophisticated statistical analyses. Lexical database (*i.e.* the SALMA-ABCLexicon) was created using the extracted information from the TAL-Corpus text.

3.3.1 Link to Lexical Database

The TAL-Corpus was used to construct the SALMA-ABCLexicon. The SALMA-ABCLexicon is a lexical database that contains around three million word-root pairs. This lexical database was extracted from the text of the TAL-Corpus following the analyses steps as described in Section 3.2. These steps include (i) manually converting the traditional Arabic dictionaries’ text into a unified format; (ii) a specialized algorithm extracts a bag of words from the definition part text of Arabic dictionaries where word-root pairs are stored; (iii) two linguistic rules were applied to the word-root pairs to verify that words are derived from the associated roots.

Later, a specialized program combines the disparate lexicon information into one large broad-coverage lexical resource the SALMA-ABCLexicon. A lexical information of a large dictionary called *لسان العرب* *lisān al-‘rab* ‘Arab tongue’ was feed to the program as a seed for the SALMA-ABCLexicon. All word-root pairs of the first dictionary were included in the SALMA-ABCLexicon which represent around 48% of the total records. Around 82% of the words and roots of *المحيط في اللغة* *mu‘ğam al-muḥīṭ fī al- luğā^h* dictionary were added which represents around 14% of total records. *تاج العروس من جواهر القاموس* *tağ al-‘arūs min ġawāhir al-qāmūs* dictionary contributes 74% of its records which represents around 22% of the total records. The percentage of added records decreases during the combination process. This decrement indicates the termination of the combination process and which traditional Arabic dictionaries are better to construct the morphological lexicon. Table [6] shows the traditional Arabic dictionaries that were used to construct the SALMA-ABCLexicon. It also shows the number of records and their percentage that contribute to the construction of the SALMA-ABCLexicon.

The SALMA-ABCLexicon contains 2 774 866 word-root pairs that represent 509 506 different words and 261 125 different non-vowelized words. It contains 12 729 roots that are distributed into 12 biliteral roots; 8 585 trilateral roots; 4 038 quadrilateral roots; 63 quinqueliteral roots; and 31 different sexiliteral roots. The 509 506 word types of the lexicon are distributed into; 117 word types derived from biliteral roots; 483 356 word types of trilateral roots; 30 873 word types of quadrilateral roots; 615 word types of quinqueliteral; and 335 word types of sexiliteral roots. Figure [8] shows the first 60 derived words of the root *كتب* *k-t-b* ‘wrote’.

Table 6. Number of records extracted and inserted in the SALMA-ABCLexicon.

#	Lexicon	Word types	Records	Percentage	
		[B]	inserted [A]	(A/B)%	(A/C)%
1	<i>lisān al-‘rab</i>	207,992	207,992	100.00%	47.80%
2	<i>mu‘ğam al-muḥīṭ fī al- luğā^h</i>	74,507	61,113	82.02%	14.04%
3	<i>tağ al-‘arūs min ġawāhir al-qāmūs</i>	128,119	95,415	74.47%	21.93%
4	<i>muḥtār aš-ṣiḥāḥ</i>	19,540	16,573	84.82%	3.81%
5	<i>al-muğrib fī tartīb al-mu‘rib</i>	12,396	9,805	79.10%	2.25%
6	<i>kitāb^u al-‘ayn</i>	30,292	18,878	62.32%	4.34%
7	<i>al-mu‘ğam al-wasīṭ</i>	36,660	25,364	69.19%	5.83%
Totals		509,506	435,140 [C]		

أكتبه	‘aktabahu	الكتاب	al-kitāb	الكتبة	al-kutba ^u
-------	-----------	--------	----------	--------	-----------------------

اَكْتَبَ	'aktaba	الكتابة	al-kitāba'	الْكُتْبَةُ	al-kutba ^{uu}
اَكْتَبْتُ	'aktabtu	الكتابة	al-kitāba ^{ta}	الْكِتَابُ	al-kitāb
اَكْتَبْنِي	'aktibnī	الكتابة	al-kitāba'	الْكِتَابَةُ	al-kitāba ^{uu}
اِكْتَابًا	'iktāb ^{an}	الكتائب	al-katātib	الْكِتَابُ	al-kitāba
اسْتَكْتَبَهُ	'istaktabahu	الكتابة	al-kitba'	الْكِتَابَةُ	al-kitāba ^{uu}
اسْتَكْتَبْتَهُ	'istaktabahu	الكتيبة	al-katība'	الْكِتَابُ	al-kitābu
اسْتَكْتَبْتُهَا	'istaktabahā	وكتيبة	wa katība'	الْكِتَابِ	al-kitābi
اَكْتَتَبَ	'iktataba	الكتائب	al-katā'iba	المكاتب	al-mukātib
اَكْتَتَبْتُ	'iktataba	الكتائب	al-katā'ibu	المكاتبة	al-mukātiba'
اَكْتَتَبَهُ	'iktatabahu	الكتيبة	al-katība ^{ta}	المكتب	al-maktab
اَكْتَتَبْتُهَا	'iktatabahā	الكتائب	al-katā'iba	المكتبة	al-maktaba'
اَكْتُبْ	'uktub	الكتابة	al-kataba'	المكتوبة	al-maktūba'
اَكْتُبْتُ	'uktutibtu	الكتب	al-katbu	الْكِتَابُ	al-kuttābu
اَكْتُبْكَ	'iktītābuk	الكتب	al-katbi	الْكِتَابُ	al-kitāba
اَكْتُبْكَ	'iktītābuka	الكتب	al-kutabu	الْكِتَابَةُ	al-kitāba ^{uu}
الْاِكْتِتَابُ	al-'iktītābu	الكتيبة	al-kutayba ^{uu}	الْكِتَابَةُ	al-kitāba ^{ti}
التكاتب	at-takātubu	الكتائب	al-kuttāba	الْمَكْتُبُ	al-maktabu
الكاتب	al-kātib	الكتائب	al-kuttābi	الْمَكْتُوبَةُ	al-maktūba ^{uu}
الكاتب	al-kātibu	الكتابة	al-kutba'	اسْتَكْتَبَ	'istaktaba

Figure 8. The first 60 lexical entries of the root كتب k-t-b 'wrote' stored in the SALMA-ABCLexicon

4. The TAL-Corpus Markup

Markups are introduced to the TAL-Corpus to indicate its features such as lexicon name, lexical entry, and definitions of lexical entries. The TAL-Corpus is formatted using XML technology where lexicons are reformatted and their lexical entries are alphabetically arranged. All traditional Arabic lexicons that form the TAL-Corpus are stored using XML files. XML is a markup language that facilitates the labelling or tagging of corpus features. The use of XML allows formatting and labelling the features of the TAL-Corpus. Figure [9] shows the XML structure and the labels used to format the corpus files.

```
<Lexicon id = "1" ar_name = "القاموس جواهر من العروس تاج" eng_name = "tağ al-'arūs min ġawāhir al-qāmūs" author_ar = "الزبيدي" author_eng = "az-zubaydī">
...
<lexicon_entry id = "8391">
<root>كتب</root>
<text>
الْبَحْيَانِي عَنْ ، كَالْبَاسِ اسْمٌ هُوَ : وَقِيلَ . الْقِيَّاسُ خِلَافَ عَلَى الْكُسْرِ ( وَكِتَابًا ) ، الْمَقْيُوسُ الْمَصْدَرُ بِالْفَتْحِ ( كُتِبَ ) ، يَكْتُبُ ، ( كُتْبَةٌ ) : كُتِبَ .
مِنْ أَقْبَلْتُ : النَّجْمُ أَبُو قَالَ ، ( خَطُّهُ ) : فِيهِمَا بِالْكَسْرِ ، وَكِتْبَةٌ ، كِتَابَةٌ : وَكَذَا . شَيْخُنَا قَالَ . مَعَانِيهِ مِنْ سِيَائِي فِيمَا اسْتَعْمِلْتُ ، الْمَصْدَرُ أَصْلُهُ : وَقِيلَ ،
النَّاءُ بِكَسْرِ ( يَكْتُبَانِ ) : النَّسَخُ بَعْضُ فِي وَرَأَيْتُ : قَالَ ، الْعَرَبُ لِسَانٍ فِي الْفَتْحِ لَمْ يَطْرُقَ فِي ثَكْنِيَّانِ مُخْتَلِفٌ بِخَطِّ رَجُلَيْنِ خُطُّ كَالْخَرْفِ زِيَادٌ عِنْدَ
أَوْ ) ، كَكُتْبَةٍ ( اَكْتُتَبَهُ ) : سَيِّدَةُ ابْنِ عَن ( وَ ) ، مُضَعَّفًا ( كَكُتْبَةٍ ) ، النَّاءُ كَسْرَةً الْكَافِ أَتْبَعَ ثُمَّ . تَعْلَمُونَ : فَيَقُولُونَ ، النَّاءُ يَكْسِرُونَ ، بَهْرَاءُ لُغَةٌ وَهِيَ
لَهُ يَكْتُبُهُ أَنْ سَأَلَهُ أَيُّ : الشَّيْءِ وَاسْتَكْتَبَهُ . لَهُ يَكْتُبُ أَنْ سَأَلَ أَيُّ : كِتَابًا فَلَانَ وَكَتَبْتُ . ( كَاسْتَكْتَبَهُ ، اسْتَمْلَاهُ ) إِذَا : ( وَكَتَبْتُ ) . ( خَطُّهُ ) إِذَا : ( كُتْبَةٌ
... </text>
</lexicon_entry>
...
<lexicon_entry id = "9657">
<root>نَجَح</root>
<text>
اللَّهُ وَأَنْجَحَهَا ) . لَكَ وَأَنْجَحْتُهَا ( وَأَنْجَحْتُ ، كَمَنْعَ ، الْحَاجَةُ نَجَحَتْ ) وَقَدْ . وَالْفَوْزُ ( بِالشَّيْءِ الطَّغْرِ : بِالضَّمِّ وَالنَّجْحُ ، بِالْفَتْحِ ، النَّجَاحُ ) : نَجَحَ :
خَطْبَةٌ وَفِي . لَهُ قَضِيَّتُهَا إِذَا ، حَاجَتُهُ أَنْجَحْتُ وَقَدْ . ( وَمَنَاجِيحُ مَنَاجِيحُ ) قَوْمٌ ( مِنْ ، مُنَجِّحٌ وَهُوَ . نَجَحَ ذَا صَارَ : زَيْدٌ وَأَنْجَحَ ) . بِإِدْرَاكِهَا أَسْعَفَهُ : ( تَعَالَى
، أَنْفَعَهُ وَبِاللَّهِ : الْأَسَاسُ سَجَعَاتٍ وَمِنْ . هِيَ وَنَجَحْتُ ، ( تَنْجَحُهَا ) إِذَا ، ( وَاسْتَنْجَحَهَا الْحَاجَةُ وَتَنْجَحُ ) . ( أَكْدَيْتُمْ إِذْ وَأَنْجَحَ ) : عَنْهَا اللَّهُ رَضِيَ عَانِشَةُ
يُحَدِّثُ نِقَابَ مَا قَطِ أَخُو جَوَادٍ نَجِيحٌ : أَوْسُ قَالَ ، الْحَاجَاتِ مُنَجِّحٌ أَيُّ ، ( النَّاسُ مِنَ الْمُنَجِّحِ ) : النَّجِيحُ ( وَ الرَّأْيُ مِنَ الصَّوَابِ : وَالنَّجِيحُ ) . اسْتَنْجَحَ وَإِيَّاهُ
( ، وَشَيْكًا أَيُّ ، نَجِيحًا سِيرًا فَلَانَ سَارَ : يَقَالُ ، ( السَّيْرُ مِنَ الشَّيْءِ ) : النَّجِيحُ : الْمَجَازُ مِنْ ( وَ ) . نَجَحَ ذُو : مُنَجِّحُ رَجُلٍ : ( الْأَسَاسُ ) وَفِي بِالْغَائِبِ
... </text>
</lexicon_entry>
...
</Lexicon>
```

Figure 9. XML structure of The Corpus of Traditional Arabic Lexicons



Search for the meaning of any Arabic root

Enter Arabic Root (أدخل الجذر العربي) :

The Root الجذر

كتب

تاج العروس من جواهر القاموس *tag al-'arūs min ḡawāhir al-qāmūs*

كتب : (كَتَبَ) ، يَكْتُبُ ، (كَتَبًا) بالفتح المصدر المقيس ، (وكتباً) بالكسر على خلاف القياس . وقيل : هر اسم كالتباس ، عن اللحياني . وقيل : أصله المصدر ، ثم استعمل فيما سباني من معانيه . قاله شيخنا . وكذا : كتابة ، وكتبته ، بالكسر فيهما : (خطه) ، قال أبو النجم : أقبلت من عند زياد كالمخرف تمشط رجلاي بغط مختلف لكتبان في الطريق لأم الغ وفي لسان العرب ، قال : ورأيت في بعض النسخ : (وكتبان) بكسر التاء ، وهي لغة بهراء بكسرون التاء ، فيقولون : تَعْلَمُونَ . ثم أجمع الكاف كسرة التاء ، (وكتبته) مضمناً ، (و) عن ابن سيده : (وكتبته) ككتبه ، (أو ككتبته) : إذا خطه . (وكتبته) : إذا استناده ، كاستكتبته . وكتب فلان كتاباً : أي سأل أن يكتب له . واستكتبته الشيء : أي سأله أن يكتبه له . وفي التنزيل العزيز : { اكتبها فهي تملى عليه بكرة وأصيلاً } (الفرقان : 5) ، أي : استكتبها . (والكتاب : ما يكتب فيه) ، وفي الحديث : (من نظر إلى كتاب العزيز :

Figure 10. Web interface for searching the traditional Arabic dictionaries

These corpus markups were effectively used when a web interface² for searching the contents of the corpus was developed. The web interface allows users to access the contents of the corpus, to search for a root and to retrieve the definition parts from the traditional Arabic lexicons included in the TAL-Corpus. Figure [10] shows part of the web interface for part of the results after searching for the root “كتب” *k-t-b*.

5. Evaluation

The purpose of constructing the TAL-Corpus is to introduce a new lexicographic corpus that contains the majority of standard Arabic vocabulary. This kind of corpus will not only help in the design and development of Arabic monolingual dictionaries but also it can support constructing Computational Linguistics resources such as; morphological dictionaries, frequency lists, lexical and morphological databases, etc. The SALMA-ABCLexicon is a lexical and morphological dictionary that was constructed using the TAL-Corpus text (see Section 3.3.1). It contains slightly under three million word-root pairs.

There are no mature standard criteria for evaluating newly constructed text corpora (Atkins et al, 1992). Therefore, our criteria for evaluating the TAL-Corpus should meet the goal for construction. We need our corpus to include the majority of standard Arabic vocabulary. Moreover, these vocabularies should be diverse and cover contemporary as well as classical ones. Lexical diversity is defined by McCarthy and Jarvis (2010) as “*the range of different words used in a text, with a greater range indicating a higher diversity*”. Lexical diversity (LD) is computed as the token-type ratio. The lexical diversity of the TAL-Corpus scored 0.152. It was evaluated by comparing it against the LD of rival Arabic corpora. The Arabic Web 2012 (arTenTen) corpus belongs to the TenTen corpora family which was created by harvesting web pages using SpiderLing. It contains around 7.5 billion tokens which represents around 2 million word types (Arts et al, 2014). Its LD scored about 0.000263. Similarly, the Arabic Internet Corpus was developed by harvesting articles from webpages published in Arabic. It contains around 165 million tokens and more than 4 million different tokens. Its LD is computed and scored 0.025965. The third corpus used in this comparative evaluation is the Arabic Wikipedia corpus (wiki-ar)³. It contains around 16 million tokens and slightly less than 1 million types. The LD for this corpus scored 0.057. Table [7] summarizes the LD for the 4 corpora used in the comparative evaluation. It shows that the LD of the TAL-Corpus scored the highest. Although it is similar size compared to the Arabic Wikipedia Corpus, its LD is 2.7 times higher. In comparison with large Arabic corpora namely: the Arabic Internet Corpus and the Arabic Web 2012 Corpus, although these large corpora contains large amounts of texts harvested from webpages, their LD is less in magnitude of times than the LD of the TAL-Corpus.

Table 7. Comparative evaluation of the LD for four Arabic corpora

<i>Corpus</i>	<i># tokens</i>	<i># Types</i>	<i>Lexical Diversity</i>
The TAL-Corpus	14 369 570	2 184 315	0.152009
Arabic Wikipedia Corpus (wiki-ar)	16 425 960	0 933 895	0.056854
Arabic Internet Corpus	165 674 718	4 301 727	0.025965
Arabic Web 2012 (arTenTen12) Corpus	7 464 566 176	1 965 566	0.000263

Another criteria for evaluating the TAL-Corpus is based on the coverage of its vocabulary on different types of text corpora. The evaluation experiments were performed using the SALMA-ABCLexicon and three text corpora: the Qur'an, the Arabic Internet Corpus⁴, and the Corpus of Contemporary Arabic. The SALMA-ABCLexicon was used because it was constructed using the TAL-Corpus and it contains all the vocabulary instances from the TAL-Corpus. The three corpora were selected to represent different types of Arabic text. The Qur'an represents Classical Arabic; the Corpus of Contemporary Corpus represents Modern Standard Arabic; and a snapshot of current Arabic language on the web is represented by the Arabic Internet Corpus.

Two experiments were conducted to compute the coverage of the TAL-Corpus. The first experiment is based on exact matching of the non-vowelized words of the three corpora with the non-vowelized words of the SALMA-ABCLexicon. The results of this experiment scored a coverage of 67.53% for the Qur'an⁵; 65.58% for the Arabic Internet Corpus; and 67.5% for the Corpus of Contemporary Arabic. Table [8] and Figure [11] show the results of the first coverage experiment. Some tokens are not words (*i.e.* Arabic words) but numbers, dates, currency symbols, punctuations, HTML or XML tags and English words. Only Arabic words were selected to compute the coverage of the SALMA-ABCLexicon.

Table 8. The coverage of the lexicon using exact word-match method

Corpus	Tokens	Arabic words	Covered words	Coverage %
Qur'an	77 800	77 799	52 536	67.53%
CCA	684 726	594 664	389 133	65.44%
Internet	1 128 114	833 916	546 880	65.58%

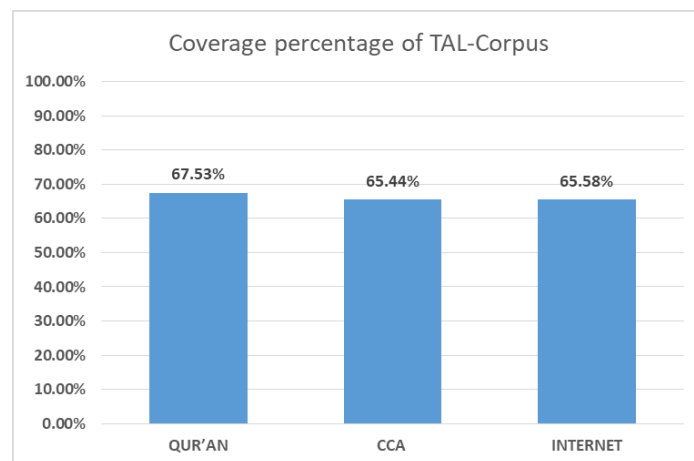


Figure 11. The coverage percentage of the TAL-Corpus using exact match method

Arabic is a morphologically rich language. Therefore, most Arabic words in context are complex words. Clitics and affixes are attached to the words in context which remarkably increase the various forms of words. Clitics make the matching process with lexical entries of the SALMA-ABCLexicon not an easy task. Hence, the coverage percentage would decrease. As an alternative, the coverage of the TAL-Corpus was computed by matching the lemmas of the SALMA-ABCLexicon with the lemmas of the three corpora. The SALMA-Lemmatizer (Sawalha, 2011) was used to lemmatize the three corpora and the lexical entries of the SALMA-ABCLexicon. The SALMA-Lemmatizer also includes a list of function words. The other part of this experiment excludes function words from the coverage calculations. Tables [9] and [10] show the coverage percentage of the TAL-Corpus computed by matching lemmas including and excluding the function words respectively. Figure 12 shows a summary of the coverage of the TAL-Corpus based on matching lemmas.

Table 9. Coverage of lemmas including function words

Corpus	Tokens	Words	Covered words	Coverage %
Qur'an	77 804	77 803	64 065	82.34%
CCA	685 161	595 099	507 943	85.35%
Internet	1 128 624	834 426	708 101	84.86%

Table 10. Coverage lemmas excluding function words

Corpus	Tokens	Words	Covered words	Coverage %
Qur'an	77 804	54 004	42 532	78.76%
CCA	685 161	411 482	338 790	82.33%
Internet	1 128 624	576 407	476 190	82.61%

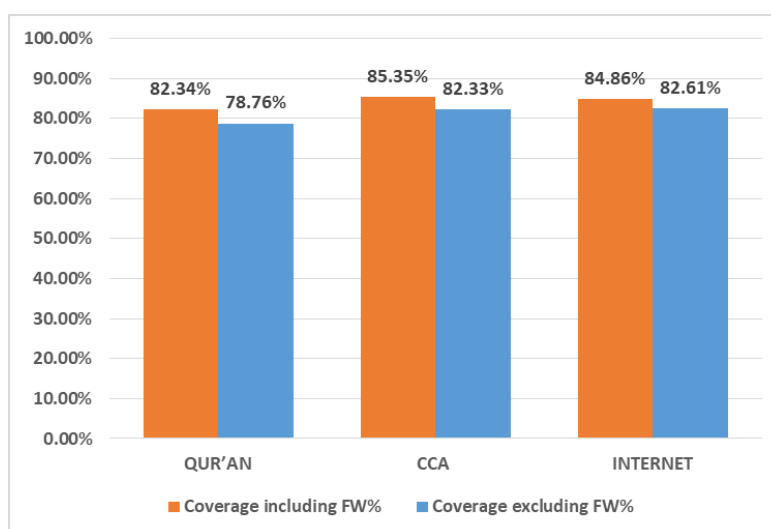


Figure 12. Coverage percentage of the TAL-Corpus using the lemmatizer

The average coverage percentage of the TAL-Corpus is 84.18% when matching the lemmas of the three corpora with the lemmas of the SALMA-ABCLexicon including function words. The coverage of the TAL-Corpus scored highest at 85.35% when computed using the CCA Corpus. The coverage scored 84.86% and 82.34% using the Internet corpus and the Qur'an respectively. The average coverage percentage of the TAL-Corpus is 81.23% after excluding function words. The highest coverage percentage was achieved using the Arabic Internet Corpus at 82.61%. Similar coverage percentage at 82.33% was achieved using the CCA corpus. Finally, 78.76% was the coverage percentage scored when the Qur'an lemmas were matched excluding function words.

The evaluation experiments of the TAL-Corpus by computing its coverage against three Arabic corpora showed that it does not fully cover words that belong to the categories; (i) function words; (ii) new Arabic terms; (iii) relative nouns; and (iv) borrowed words. Function words such as *ذَلِكَ* *dālika* "that"; *وَالِى* *wa-'ilā* "and to"; *إِنَّهُمْ* *'innahum* "they are"; and *الَّتِي* *allatī* "which" were not covered in the TAL-Corpus. These words can be easily added by including traditional Arabic grammar books in the corpus (Diwan 2004). Second, new Arabic terms such as *دردشة* *dardaša* "chat"; *انقر* *'unqur* "click" and *الانتخابات* *al-'intihābāt* "elections" are not covered because these words have appeared recently due to recent technical and social developments. Unfortunately, modern Arabic dictionaries are not available in machine readable format. Therefore, including these dictionaries in the TAL-Corpus requires retyping these dictionaries and reformatting them in a machine readable format. Third, relative nouns *الاسماء المنسوبة* *al-'asmā' al-mansūba^h* are nouns that indicate affiliation of something to these nouns. Relative nouns such as *السياحية* *as-siyāhiyya* "tourism"; *الاجتماعية* *al-iğtimā'iyya* "social"; and *الثقافية* *at-taqāfiyya* "cultural" have become widely used in the media and modern standard Arabic. Annexing this group of words to the TAL-Corpus can be achieved by including modern Arabic dictionaries. Fourth, borrowed words such as *الدكتور* *ad-duktūr* "doctor"; *الإيميل* *al-'imayl* "e-mail"; *التليفون* *at-tilifūn* "telephone"; and *الإنترنت* *al-'intarnit* "Internet" are foreign words transliterated into Arabic by using Arabic letters. Borrowed words are frequently found in newspaper and web pages text because of the lack of standard translations of them. However, Arabic

Language Academies (*i.e.* organizations which are responsible for standardizing Arabic) are producing specialized dictionaries and word lists that translate these technical terms⁶ into Arabic. These specialized dictionaries can be included in the TAL-Corpus to increase its coverage. Figure [13] shows a sample of words which are not covered in the TAL-Corpus.

ذَلِكَ	<i>dālika</i>	That	الاقتصادية	<i>al- 'iqtiṣādiyya'</i>	Economical
السَّمَاوَاتِ	<i>Assamāwāti</i>	Skies	الإنسان	<i>al- 'insān</i>	The human
إِنَّهُمْ	<i>'innahum</i>	They are	الإنترنت	<i>al- 'intarnit</i>	Internet
بِاللَّهِ	<i>Billāhi</i>	Swear to God	التلفون	<i>at-tilfūn</i>	Telephone
عَنْهُمْ	<i>'anhum</i>	After them	الفلسطيني	<i>al-filasṭīnī</i>	Palestinian
بِالْحَقِّ	<i>bilḥaqqi</i>	By the right	دردشة	<i>dardaṣa'</i>	Chat
فَأُولَئِكَ	<i>fa- 'ulā 'ika</i>	And those	انقر	<i>'unqur</i>	Click
فَبِأَيِّ	<i>fabi- 'ayyi</i>	In what	الأمريكية	<i>al- 'amrīkiyya'</i>	American
وَالِى	<i>wa- 'ilā</i>	And to	الداخلية	<i>ad-dāḥiliyya'</i>	Interior
فَسَوْفَ	<i>Fasawfa</i>	It will	الانتخابات	<i>al- 'intiḥābāt</i>	Elections
التي	<i>Allatī</i>	which	الولايات	<i>al-wilāyāt</i>	States
المتحدة	<i>al-muttaḥida'</i>	United	الاجتماعية	<i>al-iḡtimā 'iyya'</i>	Social
الدكتور	<i>ad-duktūr</i>	Doctor	الإنترنت	<i>al- 'intarnit</i>	Internet
السياحية	<i>as-siyāḥiyya'</i>	Tourism	التنمية	<i>at-tanmiya'</i>	Developmental
الغربية	<i>al-ḡarbiyya'</i>	Western	الثقافية	<i>at-ṭaqāfiyya'</i>	Cultural

Figure 13. A sample of common words which are not covered by the TAL-Corpus

6. Potential Users and Uses

The purpose for constructing the TAL-Corpus was to provide a collection of traditional Arabic dictionaries that can be analysed, studied and used to create comprehensive language resources such as; new Arabic dictionaries; frequency lists; collocates; morphological dictionaries, etc. Obviously, the potential users for the TAL-Corpus are lexicographers, Arabic linguists, language learners and computational linguists. The following is a discussion of potential uses of each expected user of this corpus.

- **Lexicographers:** This corpus was constructed as a resource for building new Arabic dictionaries. Therefore, lexicographers could use it to find examples of usage for words from different periods, track the changes in meaning of a certain vocabulary, and mark the origin of words and when they first appeared. The TAL-Corpus represents a bank of citations which are essential for the construction of new Arabic dictionaries. Citations denote objective evidence of language in use (Atkins and Rundell 2008).
- **Arabic linguists:** the TAL-Corpus provides the Arabic linguists with a repository of 23 traditional Arabic dictionaries. Feature labels (*i.e.* annotations (See Section 4) which were added to the corpus) make the search for a word, root, phrase or idiomatic expression easier via the corpus than paper based versions of traditional Arabic dictionaries. Arabic linguists are interested in studying the structures as well as the semantic features of words. The TAL-Corpus is an excellent resource for providing both. Word structures can be studied because roots and their derived words are provided. Semantic features of words such as the senses of the words; the changes to the meaning of the word; or new usage can be investigated and tracked using the TAL-Corpus. In addition, linguists can compare between the traditional Arabic dictionaries in terms of vocabulary size, ordering methodology and definitions of words. They also can conduct a comparison of other criteria such as features included in the dictionaries. These features can be the derived words, the different senses of words, phases, idioms and examples of usage.
- **Language learners:** Arabic language learners of both native and nonnative speakers use Arabic dictionaries mainly to search for words' meanings. Searching traditional Arabic dictionaries, where roots are the lexical entries, is not easy as it requires learners to know the root of the words. The TAL-Corpus provides a collection of 23 traditional Arabic dictionaries which were annotated to facilitate searching for definitions of either a word or a root. Learner can search for a word and retrieve the definition of it in addition to other linguistic information such root, lemma, derived words of the same root or lemma, examples of usage, phrases and idioms.
- **Computational linguists:** Corpora are essentially used by computational linguists to build language models for machine learning algorithms. The TAL-Corpus could be used to build language models for Arabic morphological analysers, stemmers and lemmatizers. As well as, language models for semantic analysis can

be built for Arabic using the TAL-Corpus. Computational linguists can build tracking programs that investigate the development of Arabic vocabulary and the changes of their meanings. The TAL-Corpus includes traditional Arabic dictionaries of a period that span more than 1 200 years which enables tracking the development and changes of meaning for Arabic vocabulary. In conclusion, the TAL-Corpus is an essential resource for extracting useful information that supports a wide variety of Arabic NLP applications such as; root extraction applications, morphological analysers, semantic networks of Arabic vocabulary, WordNets, ontologies ... etc.

7. Discussion of the Results, Limitations and Improvement

The TAL-Corpus is constructed using text from traditional Arabic dictionaries. It is characterized by a wide coverage of Arabic words, word types and roots. The evaluation proved that the TAL-Corpus has a wide coverage of about 85% of the test corpora words. Despite the time span of 13 centuries of the traditional Arabic lexicons from which the TAL-Corpus has been derived, only 15% of the test corpora words were not captured. The latest Arabic dictionary included in the TAL-Corpus is *المعجم الوسيط* *al-mu'ğam al-wasīf* which appeared in 1960s. Hence, new vocabulary items added to Arabic in the past 50 years are not covered in the TAL-Corpus. Moreover, due to the advances in telecommunication and information technology; globalization; and the wide and intensive use of social networks, words of foreign languages have been increasingly used in both spoken and written Arabic. These foreign words do not have a proper translation into Arabic, but are written using Arabic letters (*i.e.* transliterated). Advances in telecommunication and information technology imply new products with their original names have entered Arab countries. These products keep their original names which have been widely used and become part of the contemporary Arabic vocabulary. Moreover, the use of dialectal Arabic has increased in the written and spoken forms due to open systems such as chat rooms, blogs and forums, and social networks which allow people to write text without restrictions.

The TAL-Corpus was used to construct a broad-coverage morphological database the SALMA-ABCLexicon. This database did not involve any manual correction due to the limitations in funding. However, an automatic correction and verification procedure was applied to part of the database. The verification procedure was performed by counting how many times the word-root pairs appear in the analyzed traditional Arabic dictionaries. 976 427 word-root pairs representing 35.19% of the lexicon's word-root pairs scored a count of 2 or more. This means that these word-root pairs appeared in different dictionaries. Therefore, these word-root pairs have a high potential to be valid and correct.

This is the first version of the SALMA-ABCLexicon. It can be extended to include the full morphological analyses of the lexical entries and other useful information that will enhance the performance of NLP applications. Special linguistic lists such as compounds, collocations, idiomatic phrases, phrasal verbs and named entities can be added to extend the lexicon. Moreover, morphological lists such as broken plurals, intransitive and transitive verbs, rational and irrational words and primitive nouns can be another extension to the lexicon. The SALMA-ABCLexicon can also be extended by adding modern and dialect vocabularies from newly constructed Arabic corpora and the web.

8. Conclusions

The Corpus of Traditional Arabic Lexicons (the TAL-Corpus) is a special corpus which is constructed from the text of 23 traditional Arabic dictionaries. These dictionaries are spanning over a period of 1 200 years. The corpus contains 14 369 570 words and 2 184 315 word types. The motivation for building the TAL-Corpus is to collect and organize well-established and long traditions of traditional Arabic lexicons. The TAL-Corpus can also be used to construct new corpus-based Arabic dictionaries. Corpora were not used to construct Arabic dictionaries and lexical databases yet. Therefore, building corpora for the purpose of building new Arabic dictionaries is needed.

Thousands of traditional Arabic dictionaries were constructed in the past 1 200 years. These dictionaries are different size, type and ordering of their lexical entries. The wide variety of traditional Arabic dictionaries represent rich base for building a corpus that can be further used and exploit to construct new corpus-based Arabic dictionary.

The TAL-Corpus followed standard design and development criteria that informed major decisions in corpus creation. The text of the TAL-Corpus is composed from the text of 23 freely available and machine readable traditional Arabic dictionaries. These dictionaries were processed to have a unified format. The unified format is based on arranging the contents of the corpus by roots (*i.e.* the head words for the majority of traditional Arabic dictionaries) and their definitions. Then, the SALMA-root extractor and lemmatizer were used to tokenize, strip diacritics, and extract roots and lemmas for each word in the corpus. Frequency lists of both vowelized and non-vowelized word were also generated.

The SALMA-ABCLexicon is constructed by analysing the TAL-Corpus text. The processing steps in constructing the SALMA-ABCLexicon involve; applying linguistic rules that were encoded in a specialized program to extract the root and the words derived from that root. Second, a combination algorithm merges the information extracted from the previous step into one large broad-coverage lexical database. The SALMA-ABCLexicon contains 2 781 796 vowelized word-root pairs which represent 509 506 different non-vowelized words.

The TAL-Corpus is stored and distributed using XML technology. The corpus XML files contain all markups which indicate the corpus features. The choice of using XML technology is to facilitate the distribution and the use of the corpus. The TAL-Corpus is an open-source resource which is licenced under a Creative Commons Attribution-NonCommercial 4.0 International Licence.

The evaluation of the TAL-Corpus was done by computing its coverage over three Arabic corpora; the Corpus of the Contemporary Arabic; the Qur'an text; and the Arabic Internet Corpus. The coverage was computed by matching the words of the test corpora to the words in the SALMA-ABCLexicon, which scored about 67%. A lemmatizer program was used to compute the coverage by matching the lemmas of the test corpora and the lemmas of the SALMA-ABCLexicon. This method scored a coverage of about 82%.

The potential users for the TAL-Corpus are lexicographers, Arabic linguists, language learners and computational linguists. The potential practices for TAL-Corpus are to provide a collection of traditional Arabic dictionaries that can be analysed, studied and used to create comprehensive language resources such as; new Arabic dictionaries; frequency lists; collocates; morphological dictionaries, etc.

References

- Abbas, M., & Smaili, K. (2005) Comparison of Topic Identification Methods for Arabic Language, RANLP05 : Recent Advances in Natural Language Processing, pp. 14-17, 21-23 September 2005, Borovets, Bulgaria.
- Abbas, M., Smaili, K., & Berkani, D. (2011) Evaluation of Topic Identification Methods on Arabic Corpora, *Journal Of Digital Information Management*, 9(5), 185-192.
- Al-Sulaiti, L., & Atwell, E. (2006). The design of a corpus of contemporary Arabic. *International Journal of Corpus Linguistics*, 11(2), Jan 2006, p. 135-171. <https://doi.org/10.1075/ijcl.11.2.02als>
- Al-Thubaity, A., Khan, M., Al-Mazrua, M., & Al-Mousa, M. (2013). New Language Resources for Arabic: Corpus Containing More Than Two Million Words and a Corpus Processing Tool. *In Proceedings of the Asian Language Processing (IALP) Conference*, pp.67-70. <https://doi.org/10.1109/IALP.2013.21>
- Alansary, S., & Nagi, M. (2014) The International Corpus of Arabic: Compilation, Analysis and Evaluation. *In Proceedings of the EMNLP 2014 Workshop on Arabic Natural Language Processing (ANLP)*, 8-17, October 25, 2014, Doha, Qatar.
- AlMaayah, M., Sawalha, M., & Abushariah, M. (2016) Towards an automatic extraction of synonyms for Quranic Arabic WordNet. *Int J Speech Technol*, 19, 177. <https://doi.org/10.1007/s10772-015-9301-9>
- Alrabiah, M., Al-Salman, A. & Atwell, E. (2013). The design and construction of the 50 million words KSUCCA King Saud University Corpus of Classical Arabic. *In Proceedings of the Second Workshop on Arabic Corpus Linguistics (WACL-2)*, 22 Jul. 2013, Lancaster University, UK.
- Arts, T. (2014). The Making of a Large English-Arabic/Arabic-English Dictionary: the Oxford Arabic Dictionary. *In Proceedings of XVI EURALEX International Congress: The User in Focus*, 15-19 July 2014, Bolzano/Bozen.
- Arts, T., & McNell, K. (2013) Corpus-based lexicography in a language with a long lexicographical tradition: The case of Arabic. *In Proceedings of Second Workshop on Arabic Corpus Linguistics (WACL-2), Workshop in conjunction with the Corpus Linguistics 2013 conference*, Monday 22nd July 2013 – Lancaster University, UK.
- Arts, T., Belinkov, Y., Habash, N., Kilgarrieff, A., & Suchomel, V. (2014). arTenTen: Arabic Corpus and Word Sketches. *Journal of King Saud University-Computer and Information Sciences*, 26(4), 357-371. <https://doi.org/10.1016/j.jksuci.2014.06.009>
- Atkins, S., & Rundell, M. (2008). *The Oxford guide to practical lexicography* Oxford; New York; Oxford University Press.
- Atkins, S., Clear, J., & Ostler, N. (1992) Corpus Design Criteria. *Lit Linguist Computing*, 7(1), 1-16. <https://doi.org/10.1093/lit/7.1.1>
- Attia, M., & van Genabith, J. (2013). *A Jellyfish Dictionary for Arabic*. Retrieved Nov., 1, 2017, from

- https://www.researchgate.net/publication/259494599_A_Jellyfish_Dictionary_for_Arabic
- Black, W. J., & ElKateb, S. (2004). A Prototype English-Arabic Dictionary Based on WordNet. *The Second Global Wordnet Conference 2004* Brno, Czech Republic, January 20-23, 2004, 67-74. Retrieved from <http://www.fi.muni.cz/gwc2004/proc/95.pdf>
- Buckwalter, T. (2004). *Buckwalter Arabic Morphological Analyzer Version 2.0*: Linguistic Data Consortium, catalog number LDC2004L02 and ISBN 1-58563-324-0.
- Boudelaa, S., & Marslen-Wilson, W. D. (2010). Aralex: A lexical database for Modern Standard Arabic. *Behavior Research Methods*, 42(2), 481-487. <https://doi.org/10.3758/BRM.42.2.48>
- Diwan, A. H. (2004). المعجم النحوي لمفردات اللغة العربية *The Syntactic Lexicon of Arabic Words* (First edition). Aleppo, Syria: Fusselat Publishers.
- Elkateb, S., Black, W. J., & Farwell, D. (2006). Arabic WordNet and the Challenges of Arabic. In *Proceedings of The Challenge of Arabic for NLP/MT International Conference at The British Computer Society (BCS)*, London.
- Elkateb, S., & Black, W. J. (2001). Towards the Design of English-Arabic Terminological Knowledge Base. *Paper presented at the Proceedings of ACL 2000*, Toulouse, France:113-118.
- Eynde, V. E., & Gibbon, D. (Eds.). (2000). *Lexicon development for speech and language processing*. Dordrecht, The Netherlands: Kluwer Academic Publishers.
- Ghazali, S., & Braham, A. (2001). Dictionary Definitions and Corpus-Based Evidence in Modern Standard Arabic. *Arabic NLP Workshop at ACL/EACL*. Toulouse, France.
- Hoogland, J. (1996). The use of OCR software for Arabic in order to create a text corpus of Modern Standard Arabic for lexicographic purposes. In A. Ubaydli (Ed.), *Proceedings of the international conference and exhibition on multi-lingual computing* (pp. 2701-2716). Cambridge University.
- Ismail, O., Yagi, S., & Hammo, B. (2014). Corpus Linguistic Tools for Historical Semantics in Arabic. *International Journal of Arabic-English Studies (IJAES)* 15, 135-152.
- Khalil, H. (1998). *Dirasat fi al-lughah wa al-ma'ajim " دراسات في اللغة والمعاجم " Studies of language and lexicons* (First Edition ed.). Beirut, Lebanon: Dar al-nahdhah al-arabiah.
- Mousser, J. (2010). A Large Coverage Verb Taxonomy for Arabic. In *Proceedings of the Seventh conference on Language Resources and Evaluation (LREC'10)*, Valletta, Malta.
- Ooi, Vincent B. Y. (1998). *Computer corpus lexicography* Edinburgh: Edinburgh University Press.
- Oxford Dictionaries. (2017) The Oxford English Corpus, Retrieved Nov., 1, 2017, from <https://en.oxforddictionaries.com/explore/oxford-english-corpus>
- Rodríguez, H., Farwell, D., Farreres, J., Bertran, M., Alkhalifa, M., & Martí, M. (2008). Arabic WordNet: Semi-automatic Extensions using Bayesian Inference. In *proceedings of the 6th Conference on Language Resources and Evaluation LREC2008*, Marrakech (Morocco). Retrieved from <http://www.lsi.upc.edu/~nlp/papers/rodriguez08b.pdf>
- Saad, M., & Ashour, W. (2010), OSAC: Open Source Arabic Corpora, in *'EEECS'10 the 6th International Symposium on Electrical and Electronics Engineering and Computer Science'*, European University of Lefke, Cyprus, 118-123.
- Sawalha, M. (2011). *Open-source Resources and Standards for Arabic Word Structure Analysis*. (PhD), University of Leeds, Leeds.
- Sawalha, M., & Atwell, E. (2010). Constructing and Using Broad-Coverage Lexical Resource for Enhancing Morphological Analysis of Arabic. in: *Proceedings of the Language Resource and Evaluation Conference LREC 2010*, 17-23 May 2010, Valletta, Malta.
- van Mol, M., & Paulussen, H. (2001). AraLat: a relational database for the development of bi-lingual Arabic dictionaries. In S. Lee (Ed.), *Proceedings of Asialex 2001, Asian Bilingual-ism and the Dictionary* (pp. 206-211). Seoul, August 2001.
- van Mol, M. (2000). The development of a new learner's dictionary for Modern Standard Arabic: the linguistic corpus approach. In U. Heid, S. Evert, E. Lehmann & C. Rohrer (Eds.), *Proceedings of the ninth EURALEX International Congress* (pp. 831-836). Stuttgart, 8-12 August.

- Zaghouani, W. (2014) Critical Survey of the Freely Available Arabic Corpora, *Workshop on Free/Open-Source Arabic Corpora and Corpora Processing Tools, LREC'14 workshop*, Reykjavik, Iceland.
- Zarrouki, T., & Kebdani, M. (2009). مشروع أية -سبل القاموس العربي للتدقيق الإملائي مفتوح المصدر، تجربة وآفاق Aya-Spell Project, An Open-source Arabic Spell Checker Dictionary, experience and Future Work. *Proceedings of the workshop of morphological analyzer experts for Arabic language, organized by Arab League Educational, Cultural and Scientific Organization (ALECSO), King Abdul-Aziz City of Technology (KACT) and Arabic Language Academy.*, Damascus - Syria.
- Zarrouki, T., & Balla, A. (2009). Implementation of infixes and circumfixes in the spellcheckers. *2nd International Conference on Arabic Language Resources and Tools*, Cairo - Egypt.
- Zemanek, P. (2001). Clara (Corpus Linguae Arabicae): An Overview. In ELSNET (Ed.), *Proceedings of ACL/EACL workshop on Arabic language processing*. Toulouse, France.

Notes

- Note 1. شبكة مشكاة الإسلامية Meshkat Islamic Network <http://www.almeshkat.net>
- Note 2. A web interface for searching the traditional Arabic lexicons for a certain root http://www.comp.leeds.ac.uk/cgi-bin/scmss/arabic_roots.py
- Note 3. Frequency list of the Arabic Wikipedia corpus (wiki-ar) is found on <http://corpus.leeds.ac.uk/frqc/wiki-ar.num>
- Note 4. Leeds collection of Internet corpora: Arabic Internet Corpus <http://corpus.leeds.ac.uk/internet.html>
- Note 5. The text of the Qur'an used in this experiment was represented in MSA script.
- Note 6. Jordanian Arabic Language Academy: Word lists of technical terms <http://www.majma.org.jo/?cat=53>

Copyrights

Copyright for this article is retained by the author(s), with first publication rights granted to the journal.

This is an open-access article distributed under the terms and conditions of the Creative Commons Attribution license (<http://creativecommons.org/licenses/by/4.0/>).