

Construction and Annotation of the Jordan Comprehensive Contemporary Arabic Corpus (JCCA)

Majdi Sawalha,[†] Faisal Alshargi,^{*} Abdallah Alshdaifat,[†]
Sane Yagi,[‡] and Mohammad A. Qudah[†]

[†]University of Jordan, Jordan ^{*}Universität Leipzig, Germany

[‡]University of Sharjah, UAE

sawalha.majdi@ju.edu.jo, alshargi@informatik.uni-leipzig.de,

a.shdaifat@ju.edu.jo, saneyagi@yahoo.com, m.qudah@ju.edu.jo

Abstract

To compile a modern dictionary that catalogues the words in currency, and to study linguistic patterns in the contemporary language, it is necessary to have a corpus of authentic texts that reflect current usage of the language. Although there are numerous Arabic corpora, none claims to be representative of the language in terms of the combination of geographical region, genre, subject matter, mode, and medium. This paper describes a 100-million-word corpus that takes the British National Corpus (BNC) as a model. The aim of the corpus is to be balanced, annotated, comprehensive, and representative of contemporary Arabic as written and spoken in Arab countries today. It will be different from most others in not being heavily-dominated by the news or in mixing the classical with the modern. In this paper is an outline of the methodology adopted for the design, construction, and annotation of this corpus. DIWAN (Al-Shargi and Rambow, 2015) was used to annotate a one-million-word snapshot of the corpus. DIWAN is a dialectal word annotation tool, but we upgraded it by adding a new tag-set that is based on traditional Arabic grammar and by adding the roots and morphological patterns of nouns and verbs. Moreover, the corpus we constructed covers the major spoken varieties of Arabic.

1 Introduction

A collection of texts in machine-readable format is called a corpus. The creation of a corpus is often motivated by interest in linguistic phenomena. Therefore, the design and creation of a corpus is always linked to purpose of usage. Thousands of corpora have been created and many are freely available. These corpora vary in size, type, format, usage, and purpose of creation. They are usually annotated with morphological, syntactic, semantic, discursal, or prosodic information. Individual texts in a corpus often have meta-data in the

header that give information about such attributes as genre of the text, author, source, date and country of publication, etc. (Baker et al., 2006).

Building a balanced and representative corpus remains an ideal goal for corpus creators. A balanced corpus includes a wide range of texts from the different genres and domains that the corpus claims to depict. Sometimes, this type of corpus is referred to as a reference, general, or core corpus. Similarly, a corpus is claimed to be representative if it contains the major linguistic variation in the concerned language. Although it is not an easy task to achieve balanceness and representiveness in a corpus, it can be done with a level of approximation and scalability (McEnery and Hardie, 2012; Baker et al., 2006).

The web provides a massive collection of texts which is growing rapidly. Constructing corpora by harvesting web pages is usually referred to as web-crawling. The web is an excellent information source with large amounts of data which one can select, organize, and compile into corpora of all types (McEnery and Hardie, 2012). Since the late 1980s, Arabic corpora have been constructed. However, not many of them are freely available as open-source. Most are for written Modern Standard Arabic (MSA). Morphosyntactically annotated Arabic corpora are very rare and not freely available to researchers.

This paper reports on the construction and annotation of a comprehensive 100-million-word corpus of contemporary Arabic. The purpose is to provide an open-source corpus of contemporary Arabic which is balanced, representative of the language, and comparable to the internationally recognized British National Corpus. The text of the corpus was selected from a wide range of genres, domains, and types. It consists of 83% written language and 17% spoken language. The texts of the corpus were collected primarily from text materials available online but also from the

transcripts of purpose-made recordings (see Section 3). The corpus was automatically annotated both morphologically and syntactically. A sample of one million words was manually and semi-manually verified; it was additionally annotated for sentiment and glossed in English. To accomplish this annotation, we used DIWAN (Al-Shargi and Rambow, 2015) but had to specifically develop for it morphological and syntactic annotation schemes on the basis of the long-established Arabic linguistic tradition (see Section 5). We also added new features to the DIWAN annotation tool to facilitate our semi-manual annotation process (see Section 6).

2 Literature Review

Arabic corpora vary in size, type, purpose, design, text type, etc. (Al-Sulaiti and Atwell, 2006). Zaghouni, 2017 surveyed freely available Arabic corpora and classified 66 of them into six main categories, namely (i) raw text corpora, (ii) annotated corpora, (iii) lexicons, (iv) speech corpora, (v) handwriting recognition corpora and (vi) miscellaneous corpora.

The Corpus of Contemporary Arabic (Al-Sulaiti and Atwell, 2006) was the first freely available Arabic corpus. Around one million words were collected from newspapers and magazines. Since then, most monolingual Arabic corpora have been constructed by collecting texts from news sources (i.e. newspaper articles). Examples of such corpora are: the Open Source Arabic Corpora (OSAC) which contain around 18 million words of written MSA and Classical Arabic (CA) texts (Saad and Ashour, 2010); Akhbar Al Khaleej 2004 Corpus consists of 3 million words of newspaper texts (Abbas and Smaïli, 2005); Al-Watan 2004 Corpus contains 10 million words of newspaper texts as well (Abbas et al., 2011); KACST Arabic Corpus includes more than 700 million words collected from 10 text source types such as newspapers, magazines, books, old manuscripts, university theses, refereed periodicals, websites, curricula, news agencies, and official prints (Al-Thubaity, 2015). There is also the International Corpus of Arabic (ICA) which was constructed by Bibliotheca Alexandrina and it contains 100 million words that were collected from the press, net articles, books, and academic text sources (Alansary and Nagi, 2014). The ArabiCorpus at Brigham Young University is one of the most pop-

ular web-based corpora. It consists of around 174 million words, 77% of which is from newspapers. It does, however, include around 9 million words of premodern literature, 1 million words of modern literature, 28 million words of non-fiction, and a token of colloquial Egyptian (0.164 million words).

The King Saud University Corpus of Classical Arabic (KSUCCA) consists of around 50 million words (Alrabia et al., 2014). The corpus includes texts of six genres, namely religion, linguistics, literature, science, sociology, and biography. The arTenTen corpus used web crawlers to automatically harvest 5.8 billion words from Arabic websites (Belinkov et al., 2013). Its purpose was linguistic and lexicographic in nature. It was automatically annotated using MADAMIRA and it is available on Sketch Engine.

The Historical Arabic Corpus (HAC) has 45 million words that were organized into primary and secondary resources, seven genres, and 100-year eras in the Gregorian calendar. Its intended purpose is historical semantics and etymological lexicography (Ismail et al., 2014).

Two specialized Arabic corpora use the Quran as a source of their textual content; hence, each consists of the same number of words in the Quran, 77430 words. The Quranic Arabic Corpus is morphologically and syntactically annotated. Its annotation was done automatically and verified collaboratively by the wider community (Dukes et al., 2013). The second corpus is the Boundary Annotated Quran Corpus. It is annotated with prosodic information and phrase boundaries (Brierley et al., 2012; Sawalha et al., 2012). It took advantage of boundary markups that flag starts and stops in the Quran (Sawalha et al., 2014; Brierley et al., 2016). Interest in dialectal Arabic corpora has recently surged. An example of such corpora is the Curras Palestinian Arabic corpus, a corpus of more than 56K tokens, which are annotated with morphological and lexical features (Jar-rar et al., 2017). There are Arabic corpora that are only available for a fee, such as the Linguistic Data Consortium's¹ *The Penn Arabic Treebank*² and the European Language Resources Association's³ *An-Nahar Newspaper Text Corpus*⁴.

¹<https://www ldc.upenn.edu/>

²<https://catalog ldc.upenn.edu/LDC2016T02>

³<http://catalogue.elra.info/en-us/>

⁴<catalogue.elra.info/en-us/repository/browse/ELRA-W0027/>

This brief review, which is based on a more extensive survey of the literature, points to the absence of resources that make the claim that they represent **in a comprehensive manner** the Arabic of today **as written and spoken** by contemporary native speakers. There is a great need for a corpus of modern Arabic as used by present-day native speakers of the language. The corpus must be truly representative of the language that the current inhabitants of the Arab World use, regardless of whether it is of the high or low variety. It must also be balanced in its representation of the written and spoken language, and of the various discourse genres. It must truly depict the language of the curricula and academia.

3 Methodology

To ensure that this corpus of modern Arabic is representative, balanced, comprehensive, and for general purposes, we followed the model of the British National Corpus (BNC)⁵. That is why this corpus contains slightly more than 100-million words of the same text types, domains, and genres. The corpus contains 87% of texts from written sources and 13% of transcribed spoken language. The written part includes texts from Applied Sciences, Arts, Belief and Thought, Commerce and Finance, Imaginative works, Leisure, Natural and Pure Sciences, Social Sciences, and World Affairs. The spoken subcorpus includes transcripts of Spontaneous Conversations (4.2%) and Context-Governed Spoken Language (6.2%) from the categories of Educational/Informative, Business, Public/Institutional, and Leisure. Tables 1 and 2 show the text categories of the corpus of the written and spoken subcorpora respectively.

Twenty million words of the category of World Affairs were selected from newspapers published in 20 Arab countries where around one million words were collected for each country from one or two newspapers published in that country. The different genres of newspaper articles include Politics; Arts and Culture; Economics; Local News; Opinions; Regional and International News; Sports; and Others (e.g., Weather Forecasts, News about Technology, Health, Tourism, etc.). The subcategory of Social Sciences includes around 14 million words of texts from books and online sources. It contains texts of the genres: Languages and Linguistics; Modern Arabic Dic-

tionaries; Philosophy; Islamic Studies and Quran Interpretation; History; Geography; Anthropology and Sociology; Law; Education; Food and Nutrition; Travel; Lectures; Sports; etc. The subcategory of Belief and Thought consists of about three million words of texts of sacred books such as: the *Quran*; Quran Interpretation; the Hadith including *Hadith Qudsi*; the *Old Testament*; the *New Testament*; *Dictionary of the Bible*; and Interpretations of the Testaments, etc.

More than seven million words were collected from online sources to fill the subcategory of Commerce and Finance. These articles belong to a variety of topics within the commerce and finance genre. They include Accounting; Taxes; Investment; Finance; Financial Legal Issues; Inventory; Currency, etc. The subcategory of Imaginative Language consists of 16 million words. The texts were collected from written sources that include; stories; novels; poetry; plays; translations of international stories and novels. The subcategory of Leisure consists of 12 million words which include articles on topics such as Animals; Cars; Technology; Health; Women; Tourism; Cooking Recipes; How to; Arabian Cities; Jordanian Stories and Traditions; and Fitness. The subcategory of Arts was collected from web sources and comprises around seven million words. The texts of this category contain articles on Arts; Digital Photography; Film and Video Production; Printing; Area Planning and Landscaping; Sculpture; Ceramics and Metals; Computer Graphic Arts; Entertainment and Performance; Cinema and Theater; Photography; Music; Architecture; Fine Arts; Decorative Arts; International Arts; Arabic Calligraphy, etc. Around seven million words were collected from books and web resources for the category of Applied Sciences. The topics included in this category are Medicine; Engineering; Information Technology; Energy, etc. Finally, the Natural and Pure Sciences subcorpus consists of around four million words that come from Mathematics, Physics, Chemistry, Biology, etc.

The corpus is designed to have detailed metadata about each article. This is valuable knowledge that can be used to guide the search within the corpus. It can also be used in text classification and text data mining. Moreover, the corpus and its metadata constitute an excellent dataset for training machine learning algorithms on such tasks as genre identification. The metadata include infor-

⁵<http://www.natcorp.ox.ac.uk/>

‘and her lovers’ will have the lex حبيب *Hbyb* ‘lover’ 3) **BWhash**: In this field, the Buckwalter rendition of the lemma is split into prefix, stem, and suffix. The stem is marked by the symbol # on both sides, 4) **Gloss**: the English translation of the lemma appears in this field.

There are features in DIWAN that indicate the proclitics and enclitics of words. The clitics are assigned slots: prc3, prc2, prc1, and prc0 for proclitics, and enc0; enc1, and enc2 for enclitics. A lower index indicates closer proximity to the stem. Additionally, there are features that mark the part of speech (POS), functional number and gender of nouns, and aspect of verbs. Functional number and gender refer to the function of a word, rather than its form. For example قادة *qAdp* ‘leaders’ is functionally masculine and plural, even though it ends in ة, which is the marker of feminine singular nouns.

We added three new features to DIWAN, (i) **root** which is a base form, for example لمس *lms* to touch is the root of these two words سيلمسونها *sylmswnhA* they will touch it and يلمس *ylms* ‘he touches’, (ii) **sentiment** which shows the attitude towards a word as to whether it is negative, positive, or neutral; for example, the sentiment annotation of the word ‘sabba’ in سب العدو *sb AlEdw* ‘he cursed the enemy’ is negative while that of the word ‘ahabba’ in أحب المرأة *>Hb Almr>p* ‘he loved the woman’ is positive and that of عمان *EmAn* ‘Amman’ is neutral. And (iii) **pattern** the morphological mold that the root is formed by; e.g., the word كاسر *kAsir* breaker is derived by the mold فاعل *fAEil* doer and the root كَسَرَ *kasara* he broke. To show the details of the annotation, we present table 3.

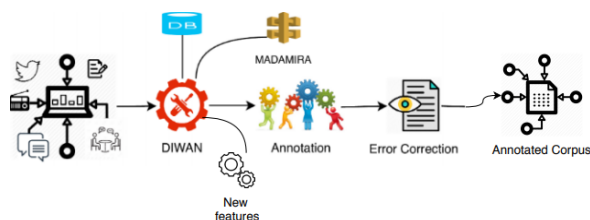


Figure 2: Steps to Creating a Comprehensive Corpus for Contemporary Arabic

6 Morphology

Morphological annotation of the whole corpus was automatically performed using MADAMIRA

(Pasha et al., 2014). We isolated a one-million word snapshot of the corpus for manual verification. Twenty-five B.A. students of Arabic at the University of Jordan carried out the manual verification and two professors of linguistics supervised their work and vetted their annotation. The annotators used DIWAN (Al-Shargi and Rambow, 2015) to review and verify MADAMIRA’s analysis. The morphological annotation required (1) Development of a new tag-set with detailed morphological description. Fourteen new noun-tags were added to Madamira. These new tags fall into three groups: i) derived nouns: Active participle, Passive participle, Exaggeration, Qualificative adjective, Noun of time/place, Noun of Instrument, and Elative noun; ii) underived nouns: Concrete noun and Abstract noun; and iii) gerunds: Original gerund, Gerund with initial miim, Gerund of instance, Gerund of state, and Gerund of profession. (2) Providing the roots of the nouns and verbs, since such a root conveys the core lexical meaning of a word. It normally consists of three consonants, and less frequently of two or four consonants. The majority of Arabic words (nouns and verbs) are derived from trilateral roots, uncommonly from biliteral or quadrilateral roots. For instance, the consonantal root د . ر . س *d.r.s* has the basic lexical meaning of studying, from which these words are derived: دَرَسَ *darosN* ‘lesson’, مُدَرِّس *mudar~is* ‘teacher’, دِرَاسَة *diraAsap* ‘studying’, مَدْرَسَة *madorasap* ‘school’, دَارِس *daAris* ‘student’. In all these derived words, the consonants d-r-s constitute their root (McCarthy, John, 1981; Prunet et al., 2000; Davis and Zawaydeh, 2001). (3) Providing the morphological pattern of each noun and verb. This pattern constitutes a canonical template that consists of a series of discontinuous consonants including those of the root, a series of discontinuous vowels, and a templatic pattern. It carries a schematic meaning and grammatical information together including the word’s part of speech. For instance, the morphological pattern C1VVC2VC3 together with the vowel melody - a - i - represents the active participle of Form I verbs (Bat-El, 1994, 2001; Ratcliffe, Robert , 1998; Ussishkin, Adam, 1999, 2005).

7 Spoken vs Written Language

Languages often have a low variety that is used in everyday communication and a high variety that is used in formal settings. The spoken language

Analyze	Sentence					
<i>Sentence</i> <i>BW</i> <i>gloss</i> <i>lex:</i> <i>pfx:</i> <i>stm:</i> <i>sfx:</i> <i>gen., num:</i> <i>root</i> <i>sntmnt</i> <i>ptrn</i>	الدولي Aldwly international dawoliy .1 AI/DET dwI/NOUN_RELATIVE (null)/CASE_DEF_GEN m,s دول neutral فُعَلِي	بالقانون bAlqAnwn law qAnuwn.1 b/PREP+AI/DET qAnwn/NOUN_ABSTRACT - m,s قَن positive فَأْتُول	الاستهتار AlAsthtAr negligence AisothAr.1 AI/DET AsthtAr/NOUN - m,s هَتَر negative اِسْتَهْتَال	في fy in - fy/PREP - none neutral none	الشركة Al\$rkp company AI/DET \$rk/NOUN p/NSUFF_FEM_SG f,s none neutral فُعَلَة	أمنت <mEnt insisted <imoEn.1 - >mEn/PV t/PVSUFF.SUBJ:2FS f,s معن neutral أَفْعَل
<i>Sentence</i> <i>BW</i> <i>gloss</i> <i>lex:</i> <i>pfx:</i> <i>stm:</i> <i>sfx:</i> <i>gen., num:</i> <i>root</i> <i>sntmnt</i> <i>ptrn</i>		فرجت frjt opened far~aj.1 - frj/PV t/PVSUFF.SUBJ:3FS f,s فرج positive فُعَل	حلقاتها HlqAthA rings Haloqap.1 Hlq/NOUN_ABSTRACT At/NSUFF_FEM_PL+ (null)/CASE_DEF_GEN+ hA/POSS_PRON_3FS f,p حلِق neutral فُعَلَة	استحكمت >sHkmt completed AstHkm.1 - AstHkm/PV t/PVSUFF.SUBJ:2FS f,s حَكَم negative اِسْتَحْكَمَل	فلما flmA when lam~A.1 f/SUB_CONJ lmA/ADV none , none none, none none neutral none	ضاق DAqt intensified dAq.1 - DAq/PV t/PVSUFF.SUBJ:3FS f,s ضيق negative فُعَل
<i>Sentence</i> <i>BW</i> <i>gloss</i> <i>lex:</i> <i>pfx:</i> <i>stm:</i> <i>sfx:</i> <i>gen., num:</i> <i>root</i> <i>sntmnt</i> <i>ptrn</i>		العلم AlElym all-knowing AI/DET Elym/ADJ_INTENS - m,s علم positive فُعِيل	السميع AlsmYE all-hearing AI/DET smYE/ADJ_INTENS 222 m,s سَمِع positive فُعِيل	وهو whw he w/CONJ hw/PRON_3MS - m,s none neutral none	الله AlIAh God All~h1 - Allh/NOUN_PROP - m,s أله positive غَال	فسيكفكيهم fsykfykhm will suffice <imoEAn.1 f/CON+s/FUT_PART+ y/IV3MS kfy/IV k/IVSUFF.DO:2MS+ hm/IVSUFF.DO:3MP m,s كني positive يَفْعِيل
<i>Sentence</i> <i>BW</i> <i>gloss</i> <i>lex:</i> <i>pfx:</i> <i>stm:</i> <i>sfx:</i> <i>gen., num:</i> <i>root</i> <i>sntmnt</i> <i>ptrn</i>	دقيقة dqyqp closely dirAsap.1 dqyq/ADJ_QUALIT f,s دق positive فُعِيلَة	دراسة drAsp studying dirAsap.1 - drAs/GERUND p/NSUFF_FEM_SG f,s درس positive فُعَالَة	المُرصودة AlmrSwdp observed maroSwd.1 AI/DET mrSwd/NOUN_ PASSEIVE_PART p/NSUFF_FEM_SG f,s رصد neutral مَفْعُولَة	الظاهرة AlmAhrp phenomenon ZAhir.1 AI/DET Zahr/NOUN_ ACTIVE_PART - f,s ظهر neutral فَاعِلَة	الباحثون AlbAHvwn researchers bAHir.1 AI/DET bAHv/NOUN_ ACTIVE_PART p/NSUFF_FEM_SG - m,p بَحْث positive فَاعِل	دُرِس drs studied darasa.1 - drs/PV - m,s درس positive فُعَل
<i>Sentence</i> <i>BW</i> <i>gloss</i> <i>lex:</i> <i>pfx:</i> <i>stm:</i> <i>sfx:</i> <i>gen., num:</i> <i>root</i> <i>sntmnt</i> <i>ptrn</i>	الشارع AlSArE public \$AriE.1 AI/DET \$ArE/NOUN_ CONCRETE - m,s شَرع neutral فَاعِل	نبض nbD pulse naboD.1 - nbD/GERUND - m,s نبض positive فُعَل	مع mE with maE.1 - mE/ADV - none, none none neutral none	يتماهى ytmAhY identify tamahY.1 y/IV3MS tmAhY/IV - m,s مهَي positive يَتَمَآغَل	الرسمي AlrsmY~ official rasomiy~.1 AI/DET rsmY/NOUN_ RELATIVE - m,s رِسم neutral فُعَلِي	الموقف Almwqf position mawoqif.1 AI/DET mwqf/GERUND_ MEEM - m,s وقف neutral مَفْعِيل
<i>Sentence</i> <i>BW</i> <i>gloss</i> <i>msa</i> <i>lex:</i> <i>pfx:</i> <i>stm:</i> <i>sfx:</i> <i>gen., num:</i> <i>root</i> <i>sntmnt</i> <i>ptrn</i>		المأخوذ AlmAxx* the thingy المأخوذ maAxuw*.1 AL/DET mAxx*/NOUN_ CONCRETE - m,s أَخَذ neutral مَفْعُول	يحب yjb to get يحب jaAb.1 y/IV3MS yjb/IV - m,s حَبِي neutral يَفْعِيل	عَآخاله ExwAlh to his uncles على آخواله xaAl.1 EIY/PREP AxwAl/NOUN_ CONCRETE h/POSS_PRON_3MS m,p خول positive فُعَال	القاروط AlqArwT the kid القاروط qaArwT.1 AL/DET qArwT/NOUN_ CONCRETE - m,s قرط negative فَأْتُول	وَدَيْت wd~yt i sent وَأَدَيْتُ wd~Y.1 - wd~Y/PV - t/PVSUFF.SUBJ:1S m,s أَدِي neutral فُعَل

Table 3: Annotated sentences of JCCA Corpus. In this table, the abbreviation *BW* represents Buckwalter transliteration, *gloss* the English meaning, *lex* the lexical entry, *pfx* the prefix, *stm* the stem, *sfx* the suffix, *gen* the gender, *root* the consonantal roots, *sntmnt* the sentiment designation, and *ptrn* the morphological pattern.

tends to be more liberal and more prone to change, the written variety more coded and more conservative. Arabic has three major varieties, two written

and one spoken: Classical Arabic, the language of scholarship until the end of the eighteenth century; Modern Standard Arabic, the language of ed-

Tag	Description	Arabic
DET	Definite Article	أداة تعريف
PREP	Prepositions	حرف جر
CONJ	Conjunction	حرف عطف
INTERROG	Interrogative particles	حرف استفهام
FUT.PART	Particles of futurity	حرف استقبال
PREFIX	Prefix	زيادة في أول الكلمة
CV.PREF	Imperative prefix	حرف أمر
IMPERF.PREF	Imperfect prefix	حرف مضارعة
INF.PART	Infinitive particle	حرف مصدرى
INF.SUBJUNC.PART	Infinitive/Subjunctive particle	حرف مصدرى ونصب
INF.ANNUL.PART	Infinitive/Annulling particle	حرف مصدرى ناسخ
NON.GOVERN	Non-Governing particle	حرف غير عامل
NEG.PART	Negative particle	حرف نفي
OTHER	Non-Governing particle	سابقة أخرى

Table 4: Prefix Tags (Bold is new)

ucation and formal written communication from the Arab renaissance in the nineteenth century onward; and the dialects, the colloquial regional varieties that are spoken in everyday communication. Since the corpus constructed here is comprehensive and since it claims to be representative of contemporary Arabic, it has to exclude Classical Arabic, but include Modern Standard Arabic, and the regional dialects. We define Contemporary Arabic as the language both written and spoken by living native speakers of Arabic; therefore, the dialects need to be represented. We are not alone in this view, check out *A Frequency Dictionary of Arabic* (Buckwalter and Parkinson, 2011) and the *Oxford Arabic Dictionary* (Arts et al., 2014).

The major spoken varieties are, therefore, represented in the corpus: North Africa is represented by the Moroccan dialect; the Nile region by Egyptian; the Arabian Peninsula by Taizi, Sanaani, and Najdi; Greater Syria by Shami, Jordanian, and Palestinian. The data in the form of contextualized sentences were collected from (1) personal communication in Facebook and Whatsapp family groups; (2) jokes, songs, videoclips, movie scripts, and TV interviews in the local dialects; and (3) personal interviews of old speakers, especially those with minimal education. The data were collected by students who came from these regions. Like any other language, Arabic has differences between the dialects and the standard variety, between the spoken and written varieties. There is variation in the pronunciation of some consonants and vowels (e.g., q, D, Z, v, *, A); suppression of word final inflections; fixed word-order (i.e., subject-verb-object (SVO)); contracted forms (e.g., *ma Zal~i\$* for *ma Zal~a \$ay'N* 'nothing remains'); use of high frequency lexical items (e.g., *قاعد qAEid* rather than

Tag	Description	Arabic
GERUND	Gerund	المصدر
GERUND.MEEM	Gerund with initial miim	المصدر الميمي
GERUND.INSTANT	Gerund of instance	مصدر المرة
GERUND.STATE	Gerund of state	مصدر الهيئة
GERUND.PROFESSION	Gerund profession	مصدر صناعي
NOUN.CONCRETE	Concrete noun	ام ذات
NOUN.ABSTRACT	Abstract noun	ام معنى
NOUN.ACTIVE.PART	Active participle	ام فاعل
NOUN.PASSIVE.PART	Passive participle	ام مفعول
ADJ.INTENS	Form of exaggeration	صيغة المبالغة
ADJ.QUALIT	Adjective	الصفة المشبهة
NOUN.TIME.PLACE	Noun of time/place	ام الزمان والمكان
NOUN.INSTRUMENT	Instrumental noun	ام الآلة
ADJ.COMP	Elative noun	ام التفضيل
NOUN.RELATIVE	Relative noun	ام منسوب
NOUN.PROP	Proper noun	ام علم
NOUN.PROP.FOREIGN	Foreign proper noun	ام علم أجنبي
ADV	Adverb	الظرف
PRON	Pronoun	الضمير المنفصل
DEM.PRON	Demonstrative pronoun	ام الإشارة
REL.PRON	Relative pronoun	ام موصول
INTERROG.PRON	Interrogative pronoun	ام استفهام
REL.ADV	Conditional noun	ام شرط
NOUN.VERB.LIKE	Verb-like noun	ام الفعل
NOUN.FIVE	Five nouns	الأسماء الخمسة
NOUN.DIMINUTIVE	Diminutive	ام تصغير
NOUN.BLEND	Blend noun	ام منحوث
NOUN.NUM	Numeral	ام عدد
EXCEPT.NOUN	Exceptive Noun	ام استثناء
COMP.NOUN	compound noun	ام مركب
FOREIGN	Foreign word	كلمة أجنبية
ABBREV	Abbreviation	اختصار
PV	Perfect verb	فعل ماض
PV.PASS	Passive Perfect v.	فعل ماض مجهول
IV	Imperfect v.	فعل مضارع
IV.PASS	Passive Imperfect v.	فعل مضارع مجهول
UNINFLECTED.VERB	Uninflected Verb	فعل جامد
CV	Imperative verb	فعل أمر
PREP	Preposition	حرف جر
NEG.PART	Preposition	حرف نفي
CONJ	Conjunction	حرف عطف
INTERROG.PART	Interrogative particle	حرف استفهام
SUBJUNC.PART	Subjunctive particle	حرف نصب
JUSSIVE.PART	Jussive particle	حرف جزم
ANNUL.PART	Annulling particle	حرف ناسخ
VOC.PART	Vocative particle	حرف نداء
EXCEPT.PART	Exceptive par.	حرف استثناء
FUTUR.PART	Par. of futurity	حرف استقبال
YES.NO.RESP.PART	Yes/No particle	حرف جواب
CONDITION.PART	conditional particle	حرف شرط
CERT.PART	Certain/Uncertain particle	حرف تحقيق
PART	other particles	حروف أخرى
PUNC	Punctuation mark	علامة ترقيم
NUMBER	Number	رقم
CURRENCY	Currency	عملة
DATE	Date	تاريخ
NON-ARABIC	Non-Arabic word	كلمة غير عربية
OTHER	OTHER	أخرى

Table 5: Stem Tags (Bold is new)

جالس *jAlis* 'sitting'); use of some lexical items that are archaic in MSA (e.g., *AifliH* 'Partake of food' in Jordanian Arabic in addition to the senses in Standard Arabic of Plough! and Succeed!); liberal incorporation of foreign words (e.g., *mas~aj* 'sent a message'); abandonment of the dual and the passive voice (e.g., *إنكسر* *inkasar* 'broke' rather than *كُسِرَ kusira* 'it got broken'); abandonment of the yes-no question

Tag	Type	Arabic	Tag	Type	Arabic
POSS.PRON	Proclitic	ضمير متصل بالاسم	SUBJ.PRON	Suffix	ضمير متصل بالفعل
OBJ.PRON	Proclitic	ضمير متصل بالفعل (مفعول به)	SUFF.FEM.TA	Proclitic	تاء التانيث
NSUFF.FEM.SG	Proclitic	تاء مربوطة	RELATIVE.YA	Proclitic	ياء النسبة
CASE.INDEF.ACC.GEN	Suffix	التنوين	SUFF	Suffix	زيادة في آخر الكلمة
NSUFF.FEM.PL	Proclitic	حروف جمع المؤنث	NSUFF.MASC.PL.NOM	Proclitic	حروف جمع مذکر مرفوع
NSUFF.MASC.PL.ACC	Proclitic	حروف جمع مذکر منصوب	NSUFF.MASC.PL.GEN	Proclitic	حروف جمع مذکر مجرور
NSUFF.MASC.DU.NOM	Proclitic	حروف المثنى مذکر مرفوع	NSUFF.MASC.DU.ACC	Proclitic	حروف المثنى مذکر منصوب
NSUFF.MASC.DU.GEN	Proclitic	حروف مثنى مذکر مجرور	NSUFF.FEM.DU.NOM	Proclitic	حروف المثنى مؤنث مرفوع
NSUFF.FEM.DU.ACC	Proclitic	حروف مثنى مؤنث منصوب	NSUFF.FEM.DU.GEN	Proclitic	حروف مثنى مؤنث مجرور
EMPHATIC.NUN	Suffix	نون التوكيد	PROTECT.NUN	Suffix	نون الوقاية
REL.PRON	Relative Pronoun	اسم موصول	ADV	Adverb	ظرف
SINGLAR	Number/Singular	مفرد	DUAL	Number/Dual	مثنى
PLURAL	Number/Plural	جمع سائر	BROKEN.PLR	Number/Broken plural	جمع تكسير
COLCV.NOUN	Number/Collective noun	اسم الجمع			

Table 6: Tags for suffixes (Bold is new)

particles *هل hal* and *أ >*; use of the suffix *شـ* \$ at the end of a verb (e.g., *ما قعدش mA qaEadi\$* rather than *ما قعدَ mA qaEadahe* did not sit); loss of gender distinction, especially in the language of females (e.g., *إجو البنات <ijw AlbanAt* rather than *جاءت البنات jA'at AlbanAt* 'the girls came'). Arabic has a free word order because of grammatical inflections. When all words' grammatical functions are marked with appropriate inflections, it is not necessary to restrict the arrangement of words in a sentence; hence, Classical Arabic exhibits a totally free word order. Modern Standard Arabic shows preference for verb-subject-object even though inflections are amongst its distinctive features. The spoken varieties continue a historical tradition that we suspect had started as early as Islamic times, where case inflection had lost grounds to fixed word order. Preference in Classical Arabic for the default word order (i.e., verb-subject-object) in an otherwise free word order system was a portent of developments to come. As Islamic conquest brought Arabs in contact with foreigners who soon adopted the language, and as the diglossic gap widened, grammatical inflection lost favor in the low variety while it retained its glamour in the high variety, under the influence of the Quran. The spoken, the low, variety started to favor the subject-verb-object word order as a result of the loss of case inflections and to set apart the agent from the patient of the predicate. The written variety manifested in MSA, on the other hand, used the verb-subject-object order as the unmarked default and retained other combinations for special purposes. All modern regional varieties are descendants of old spoken varieties of Arabic in much the same way as Modern Standard Arabic is a successor of Classical Arabic, the written variety. Regional varieties of Arabic share

great many syntactic features. For example, they have two negation patterns: single negation and discontinuous negation (Alqassas, 2015). The first uses the negative particle *ما mA* followed by the verb phrase, whilst the second adds the negative marking suffix *شـ* \$ to the verb in addition to the negative particle that precedes it. Thus, I didnt say may be expressed as *ما قلتش mA qult-i\$* or *ما قلت mA qult*. To negate the future, however, there are three options: (1) the negative particle followed by the imperfect verb as in *ما أسافر mA >asAfir* 'I will not travel'; (2) or followed by the imperfect inflected with the negative marking suffix as in *ما أسافرش mA >asAfr-i\$*; (3) or followed by the future particle *رح raH* and the imperfect verb as in *ما رح أسافر mA raH >asAfir*. JCCA consists in part of a spoken language component that is annotated morphologically and syntactically, glossed with MSA forms, and translated into English. This is especially useful with contractions, the hallmarks of spoken Arabic. The gloss is often the non-contracted equivalent in MSA as demonstrated in Table 7.

8 Conclusion and Future Work

This paper outlined the methodology for the design, construction, and annotation of the Jordan Comprehensive Contemporary Arabic Corpus (JCCA). The corpus is balanced, comprehensive, and representative of contemporary Arabic as written and spoken in Arab countries today. It consists of 100 million words that reflect current usage of the language. The corpus consists of 87% written and 13% spoken language. The text of the corpus was selected such that it would be representative of a wide range of geographical regions, genres, subject matters, modes, and media. DI-

Contracted	BW	Full Form	Gloss
شلونك	\$lwnk	أي شيء لونك	how are you?
اصطفل	ASTfl	اصطف الذي تريد	whatever you want
إيش	<y\$	أي شيء	pardon me?
لبش	ly\$	لأي شيء	why?
شو	\$w	أي شيء هو	what?
بيش	by\$	بأي شيء	for how much?
قدش	qdy\$	قدر أي شيء	how much?
معلش	mEly\$	ما عليك شيء	it's OK!
مظلش	mZl\$	ما ظل شيء	nothing left
إللي	<lly	الذي، التي	that/which

Table 7: Contracted words in colloquial Arabic, In this table, the abbreviation *Contracted* represents examples of spoken words (i.e. contractions), *BW* is Buckwalter transliteration, *Full Form* is the non-contracted equivalent in MSA, *gloss* the English meaning.

WAN was upgraded and used to annotate and manually verify the annotation of a one-million-word snapshot of the corpus, making it a gold standard of superior quality that can serve as a resource against which automatic annotation may be compared. JCCA construction made these additional contributions: (i) Development of a new and elaborate tag-set that is based on the morphology of traditional Arabic grammar; (ii) Addition of the roots and morphological patterns of nouns and verbs; (iii) Coverage of the major spoken varieties of Arabic: North Africa; the Nile; the Arabian Peninsula; and Levant. Future work is to make this corpus a monitor corpus where new texts are added proportionally every year. This will facilitate tracking language change and will render the corpus more amiable to lexicography.

9 Acknowledgment

The research reported here was supported by the Scientific Research Fund of the Ministry of Higher Education and Scientific Research, Jordan (Grant No. Soci/2/1/2016).

References

Mourad Abbas and Kamel Smaïli. 2005. [Comparison of Topic Identification methods for Arabic Language](#). In *International Conference on Recent Advances in Natural Language Processing - RANLP 2005*, 14-17, Borovets, Bulgaria.

Mourad Abbas, Kamel Smaili, and Berkani. 2011. Evaluation of topic identification methods on Arabic corpora. *Journal of Digital Information Management*, 9(5):185–192.

Faisal Al-Shargi and Owen Rambow. 2015. [DIWAN: A dialectal word annotation tool for Arabic](#). In *Proceedings of the Second Workshop on Arabic Natural Language Processing*, pages 49–58, Beijing, China. Association for Computational Linguistics.

Latifa Al-Sulaiti and Eric Steven Atwell. 2006. [The design of a corpus of Contemporary Arabic](#). *International Journal of Corpus Linguistics*, 11(2):135–171.

Abdulmohsen Al-Thubaity. 2015. [A 700m+ Arabic corpus: KACST Arabic corpus design and construction](#). *Language Resources and Evaluation*, 49(3):721–751.

Sameh Alansary and Magdy Nagi. 2014. The International Corpus of Arabic: Compilation, Analysis and Evaluation. In *Proceedings of the EMNLP 2014 Workshop on Arabic Natural Language Processing (ANLP)*, pages 8–17, Doha, Qatar.

Ahmad Alqassas. 2015. [Negation, tense and npis in jordanian arabic](#). *Lingua*, 156:101–128.

Maha Sulaiman Alrabia, AbdulMalik Al-Salman, Eric Atwell, and Nawal Alhelewh. 2014. [KSUCCA: A Key To Exploring Arabic Historical Linguistics](#). *International Journal of Computational Linguistics (IJCL)*, 5(2):27–36.

Faisal AlShargi, Aidan Kaplan, Ramy Eskander, Nizar Habash, and Owen Rambow. 2016. [Morphologically annotated corpora and morphological analyzers for moroccan and sanaani yemeni arabic](#). In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC 2016)*.

Tressy Arts, Yonatan Belinkov, Nizar Habash, Adam Kilgarriff, and Vit Suchomel. 2014. [artenten: Arabic corpus and word sketches](#). *Journal of King Saud University - Computer and Information Sciences*, 26(4):357 – 371. Special Issue on Arabic NLP.

Paul Baker, Andrew Hardie, and Tony McEnery. 2006. *A Glossary of Corpus Linguistics*. Edinburgh University Press.

Outi Bat-El. 1994. [Stem modification and cluster transfer in Mmodern Hebrew](#). *Natural Language Linguistic Theory*, 12(4), 571-596.

Outi Bat-El. 2001. [In search for the roots of the C-root: The essence of Semitic morphology](#). Workshop on Root and Template Morphology. Los Angeles: University of South California.

Yonatan Belinkov, Nizar Habash, Aadm Kilgarriff, Noam Ordan, Ryan Roth, and Vit Suchomel. 2013. [arTenTen12: A new, vast corpus for Arabic](#). In *Second Workshop on Arabic Corpus Linguistics, WACL'S*, Lancaster University, UK.

Claire Brierley, Majdi Sawalha, and Eric Atwell. 2012. [Open-Source Boundary-Annotated Corpus for Arabic Speech and Language Processing](#). In *LREC*, pages 1011–1016.

- Claire Brierley, Majdi Sawalha, Barry Heselwood, and Eric Atwell. 2016. A Verified Arabic-IPA Mapping for Arabic Transcription Technology, Informed by Quranic Recitation, Traditional Arabic Linguistics, and Modern Phonetics. *Journal of Semitic Studies*, 61(1):157–186.
- Tim Buckwalter and Dilworth Parkinson. 2011. *A frequency dictionary of Arabic: Core vocabulary for learners*. London: Routledge.
- Stuart Davis and Bushra Zawaydeh. 2001. *Arabic hypocoristics and the status of the consonantal root*. *Linguistic Inquiry*, 32(3): 512-520.
- Kais Dukes, Eric Atwell, and Nizar Habash. 2013. [Supervised collaboration for syntactic annotation of quranic arabic](#). *Language Resources and Evaluation*, 47(1):33–62.
- Thomas Eckart, Faisal Al-shargi, Uwe Quasthoff, and Dirk Goldhahn. 2014. Large arabic web corpora of high quality: The dimensions time and origin. In *Workshop on Free/Open-Source Arabic Corpora and Corpora Processing Tools, LREC, Reykjavk*.
- Omaima Ismail, Sane Yagi, and Basam Hammo. 2014. Corpus Linguistic Tools for Historical Semantics in Arabic. *International Journal of Arabic-English Studies (IJAES)*, 15:135–152.
- Mustafa Jarrar, Nizar Habash, Faeq Alrimawi, Diyam Akra, and Nasser Zalmout. 2017. [Curras: an annotated corpus for the palestinian arabic dialect](#). *Language Resources and Evaluation*, 51(3):745–775.
- McCarthy, John. 1981. *A prosodic theory of nonconcatenative morphology*. *Linguistic Inquiry*, 12, 373-418.
- Tony McEnery and Andrew Hardie. 2012. *Corpus Linguistics: Method, Theory and Practice*. Cambridge Textbooks in Linguistics. Cambridge University Press.
- Arfath Pasha, Mohamed Al-Badrashiny, Mona T Diab, Ahmed El Kholy, Ramy Eskander, Nizar Habash, Manoj Pooleery, Owen Rambow, and Ryan Roth. 2014. Madamira: A fast, comprehensive tool for morphological analysis and disambiguation of arabic. In *LREC*, volume 14, pages 1094–1101.
- Jean-Francois Prunet, Renee Bland, and Ali Idrissi. 2000. *The mental representation of Semitic words*. *Linguistic Inquiry*, 31(4). 609-648.
- Ratcliffe, Robert . 1998. *The broken plural problem in Arabic and comparative Semitic: allomorphy and analogy in non-concatenative morphology*. Amsterdam/Philadelphia: John Benjamins.
- Motaz Saad and Wesam Ashour. 2010. Osac: Open source arabic corpora. In *EEEECS10 the 6th International Symposium on Electrical and Electronics Engineering and Computer Science*, pages 118–123.
- Majdi Sawalha, Claire Brierley, and Eric Atwell. 2012. Prosody Prediction for Arabic via the Open-Source Boundary-Annotated Quran Corpus. *Journal of Speech Sciences*, 2(2):175–191.
- Majdi Sawalha, Claire Brierley, and Eric Atwell. 2014. Automatically generated, phonemic Arabic-IPA pronunciation tiers for the Boundary Annotated Qur’an Dataset for Machine Learning (version 2.0). In *proceedings of LRE-Rel 2: 2nd Workshop on Language Resource and Evaluation for Religious Texts, LREC 2014 post-conference workshop 31st May 2014, Reykjavik, Iceland*, page 42. LRA.
- Ussishkin, Adam. 1999. *The inadequacy of the consonantal root: Modern Hebrew denominal verbs and output-output correspondence*. *Phonology*, 16(3), 401-442.
- Ussishkin, Adam. 2005. *A Fixed prosodic theory of nonconcatenative templatic morphology*. *Natural Language Linguistic Theory*, 23(1), 169-218.
- Wajdi Zaghouani. 2017. Critical Survey of the Freely Available Arabic Corpora. *CoRR*, abs/1702.07835.