

A Comprehensive Full-Form Lexicon for Arabic NLP and Speech Technology

Jack Halpern

The CJK Dictionary Institute
34-14, 2-chome, Tohoku, Niiza-shi
Saitama 352-0001
JAPAN
jack@cjki.org

Yannis Haralambous

IMT Atlantique & UMR CNRS 6285 LabSTICC
Technopole Brest-Iroise, CS 83818
29238 Brest Cedex 3
FRANCE
yannis.haralambous@imt-atlantique.fr

Abstract

Natural Language Processing (NLP) applications require morphological data with precise grammatical attributes, while speech technology requires abundant phonemic and phonetic transcriptions. This presents a challenge for Arabic due to its abundant morphological, orthographic, and phonemic variation. Existing systems encounter challenges in processing incomplete and unstructured data from web sources, leading to suboptimal performance in morphological analysis and speech technology. ArabLEX, a comprehensive full-form lexicon for MSA, addresses these issues by providing a foundation for enhancing NLP precision. It comprises over 530 million entries with fully inflected, conjugated, declined, and cliticized forms accompanied by detailed morphological attributes as well as precise phonological transcriptions and orthographic variants. This combines exhaustive listing of forms with detailed descriptions that can significantly mitigate the inherent ambiguity of Arabic. It can serve as a foundation for developing accurate NLP and speech technology applications by providing accurate orthographic variants in both the Arabic script and in phonemic transcriptions.

1 Introduction

1.1 What is a Full-Form Lexicon

According to Crystal's *Dictionary of Linguistics* (Crystal, 2008), a word is a "unit of expression which has universal intuitive recognition by native speakers." Although this does not provide an objective criterion for "wordhood," words are an important notion in NLP. Ever since the first dictionary in history, the *Sumerian Lexicon* (2300 BCE), lexicographers have worked on the collection of

"words." Traditionally, the headwords of dictionaries have been canonical forms (lemmata). A rare dictionary format is the *full-form lexicon*. It explicitly includes all word forms of a language, i.e., fully inflected, conjugated, declined, or cliticized ("inflected" for short) members of a lexeme class, rather than just the lemmata. For example, the English lexemes *eat* and *boy* have the members *eat*, *eats*, *eating*, *eaten*, *ate* and *boy*, *boys*, *boy's*, *boys'* respectively. For highly inflected languages like Arabic, the abundance of combinatorics (stem, affixes, clitics) can result in full-form lexicons with hundreds of millions of entries.

Historically, Machine Translation (MT) and other NLP applications relied on rules or statistical models. In recent years, utilizing neural machine translation (NMT) and large language models (LLM) has become the norm. These tools rely on efficient disambiguation. Despite their remarkable achievements, challenges remain in Arabic NMT, such as the handling of proper nouns and multi-word expressions (MWE) (Halpern, 2019) and overcoming the lack of bilingual training corpora.

1.2 The Case of Arabic

Arabic is special in its morphology. On the one hand, Modern Standard Arabic (MSA), used in the media, government, and education, is the official language of 380 million people, but (practically) no one's mother tongue (Haugen, 1972; Mejdell, 2014). A kind of hierarchical polyglossia is the norm. On the other hand, the morphology of Arabic is based on roots and patterns (templatic morphology) (Ryding, 2005), so we are not just dealing with stems and affixes as in Roman languages but with tri- or quadriliteral consonantal roots with infixes, prefixes, suffixes and circumfixes. This morphological generative principle is omnipresent and

75 even applies to loanwords (Gadelli, 2015). It can
76 thus be considered to be an innate property of Ara-
77 bic. Therefore, a full-form lexicon should cover all
78 MSA root + pattern combinations (so that all gram-
79 matical word forms are available to the user),
80 which is necessary for both speech recognition and
81 written text.

82 1.3 Previous Work

83 Several Arabic modeling tools have been devel-
84 oped for morphological analysis, tokenization,
85 generation of inflected and conjugated forms, POS
86 tagging, and disambiguation. We refer to such tasks
87 as analysis and generation, and to such tools as
88 morphological engines. Popular tools include
89 AlKhalil (Boudchiche *et al.*, 2017), MADA (Ha-
90 bash, Rambow, and Roth, 2009), BAMA (Buck-
91 walter, 2002), PATB (Penn Arabic Treebank)
92 (Maamouri *et al.*, 2004), FARASA (Abdelali *et al.*,
93 2016), MADAMIRA (Pasha *et al.*, 2014), and
94 Elixir_FM (Smrž, 2007). A more recent, highly
95 ambitious tool is CALIMA Star (Taji *et al.*, 2018).

96 Despite the high performance of these tools (Taji
97 *et al.*, 2018), they have shortcomings, such as in-
98 consistency, ignoring lexical rationality, and lack-
99 ing phonological attributes. Naturally, the pro-
100 cessing performed by morphological engines is
101 supported by lexical databases, such as tables for
102 stems, clitics, and affixes (Halpern, 2018). Still, the
103 goal of these tools is to perform computational
104 tasks such as tokenization and disambiguation ra-
105 ther than serving as comprehensive lexicons for
106 enumerating all possible ambiguous sequences. A
107 notable outlier worth mentioning is the Arabic full-
108 form lexicon and Finite State Transducer (FST)
109 project by Souidi and Eisele (2004).

110 1.4 Introducing ArabLEX

111 Unlike morphological engines, ArabLEX is a
112 stand-alone lexical database that can be integrated
113 with such engines. It does not perform computa-
114 tional tasks itself. Its goal is to act as a comprehen-
115 sive database to support morphological engines
116 and NLP tools. In theory, an engine can query the
117 lexicon as an external module via function call or
118 API, but ideally, it should be integrated directly. If
119 a morphological engine is likened to the engine of
120 a car, then a full-form lexicon like ArabLEX is the
121 fuel – it drives the engine, not the car itself (e.g., a

122 text-to-speech (TTS) application that relies on a
123 pronunciation dictionary).

124 ArabLEX is a *full-form lexicon* aiming to be as
125 comprehensive as possible, though some word
126 classes, such as periphrastic elatives, have not yet
127 been included. In the first phase (May 2024) Arab-
128 LEX contains about 530 million entries for content
129 words (nouns, adjectives, and verbs) in the do-
130 mains of general vocabulary and (for the first time)
131 fully inflected and cliticized proper nouns for both
132 Arab and non-Arab personal names and place
133 names. It provides exhaustive coverage of all in-
134 flected, declined, conjugated and cliticized forms
135 and includes a rich set of grammatical, morpholog-
136 ical, phonological, and orthographic attributes, as
137 shown in detail by The CJK Dictionary Institute
138 (2020). This makes it suitable for NLP applications
139 such as machine translation, named entity recogni-
140 tion, and morphological analysis and generation.
141 For example, the verb *katabta* is one of 7,251 possi-
142 ble forms of *kataba*. It has tags such as 2SM,
143 meaning second person masculine singular. Special
144 emphasis is placed on speech technology by
145 providing such attributes as accurate phonemic
146 transcriptions as well as full diacritization.

147 Note that the phonemic¹ transcriptions in this
148 paper are italicized and given in the CARS system
149 (Halpern, 2009), designed by our institute for ped-
150 agogical and speech applications. Transliterations
151 are given in the Buckwalter transliteration system
152 (Buckwalter, 2002) and enclosed in forward-
153 slashes. Note also that ArabLEX is undergoing
154 maintenance and expansion, and it is expected to
155 exceed one billion entries, making it—to our
156 knowledge—the most comprehensive Arabic com-
157 putational lexicon ever created.

158 2 Levels of Ambiguity in Arabic

159 2.1 Morphological / Lexical Ambiguity

160 In templatic morphology, inflection is performed
161 by changing the vowel + consonant patterns by af-
162 fixation and cliticization. Not only can words be in-
163 flected, declined, and conjugated (“inflected” for
164 short), but they can also take many clitics. For ex-
165 ample, adding the proclitics *wa* ‘and’, *li* ‘to’, and
166 the enclitic *ātīhimā* to the stem *kātib* ‘writer’ yields
167 the complex form *walikātībātīhimā* (وَلِكَاتِبَاتِهِمَا)
168 ‘and to their (dual) female writers’. This results in

¹ Technically, CARS is a morpho-phonemic transcription system, as it encodes information such as vowel neutralization.

169 a very large number of word forms. For example, 170 the full paradigms for كَاتِبٌ *kātibun* ‘writer’ and 171 كَتَبَ *kataba* ‘write’ reach about 5,660 and 6,900 172 forms, respectively.

173 The difference between morphological and lexi- 174 cal ambiguity is analogous to the difference be- 175 tween inflection and derivation in Western lan- 176 guages: when a word is inflected, the forms we ob- 177 tain represent the same lexeme; when it is derived, 178 we move to a different lexeme. This happens also 179 in Arabic, e.g., the graphemic sequence كَتَبَ may 180 denote كَتَبْتُ ‘I wrote,’ or كُتِبَ ‘books.’ The lexeme 181 of the former is the verb ‘to write,’ and the lexeme 182 of the latter is the noun ‘book.’

183 Distinguishing between morphological and lex- 184 ical ambiguity is computationally relevant because 185 the latter implies multiple POS tags and, therefore, 186 also potentially multiple syntax trees.

187 2.2 Recognition of Arabic Patterns

188 Conventional wisdom has it that Arabic is ambigu- 189 ous “due to the non-representation of short vowels.” 190 In fact, a whole gamut of factors contributes to am- 191 biguity (Halpern, 2002), including (1) the absence 192 of short vowels (e.g., كَاتِبَ represents the seven 193 word forms *kātib*, *kātibun*, *kātibin*, *kātaba*, *kātibi*, 194 *kātiba*, *kātibu*), (2) representation of long *ā* by *ا* as 195 in سوريَا or by *أ* as in آسِيَا, but some bare alifs rep- 196 resenting *tanwiin* rather than long *ā*, as in شُكْرَا 197 *shukran*, (3) *ʿalif alfaa* *Sila* (otiose *alif*) (Ryding, 198 2005), orthographic conventions not being pro- 199 nounced (e.g., كَتَبُوا being realized as *katabu*²), (4) 200 the omission of *shadda* indicating consonant 201 gemination, e.g., مُحَمَّدَ (diacriticized مُحَمَّد), 202 which provides no clues that the /m/ is geminated, 203 and (5) vowel neutralization sometimes being lex- 204 ically determined and thus unpredictable from the 205 orthography, e.g., فِي الْقَاهِرَةِ ‘in Cairo’, the prepo- 206 sition /fyi/ is pronounced *fī*, not *fii*.

207 Examples (1)–(4) given above are cases of gra- 208 phemically under-represented patterns. Indeed, 209 patterns may contain short vowels or conso- 210 nants/long vowels that are written but must be rec- 211 ognized as being part of a pattern.

212 2.3 Orthographic Disambiguation

213 A central issue in Arabic NLP, especially speech 214 technology, is identifying which word form an am- 215 biguous string like كَاتِبَاتِك represents. This string 216 can represent any of six-word forms, each with a 217 different meaning, a different pronunciation and 218 potentially a different POS tag and/or syntactic 219 function.

220 The process of identifying the correct form is re- 221 ferred to as orthographic disambiguation (Halpern, 222 2008). The rich set of grammatical and morpholog- 223 ical attributes in ArabLEX can help language mod- 224 els to correctly disambiguate such forms.

225 2.4 Word Stress and Vowel Neutralization

226 Prosody (word stress) and vowel neutralization 227 play a critical role in ensuring that synthesized 228 speech sounds natural. نَا *naa*, for example, is writ- 229 ten as a long vowel in أَنَا but is shortened to *na*. This 230 complex issue is described in detail in Halpern's 231 (2009) paper on Arabic stress.

232 The morpho-phonemic and phonetic transcrip- 233 tions in ArabLEX explicitly indicate precise word 234 stress and vowel neutralization for each entry. For 235 example, in the IPA [wēlikēːˈtibikumə(ː)], the 236 stressed syllable is indicated by (ˈ) (U+0C28), 237 while (ː) (U+02D1) indicates that the final *ε* is a 238 neutralized vowel of optional half-length.

239 3 Enhancing Speech Technology

240 3.1 Arabic Speech Technology

241 Though advances in neural networks have dramati- 242 cally improved the quality of speech technology, 243 in a 2020 survey, we compared the TTS systems 244 provided by leading IT companies, showing that 245 Arabic significantly lags behind other major lan- 246 guages (Halpern, 2020). ArabLEX addresses these 247 shortcomings by serving as a comprehensive pro- 248 nunciation dictionary that enhances the quality of 249 both TTS and automatic speech recognition (ASR). 250 It includes an NLP-oriented morpho-phonemic 251 transcription called CARS (Halpern, 2009) and two 252 phonetic transcriptions: SAMPA (Wells, 1997) and 253 IPA (International Phonetic Association, 1999), to 254 support the training of TTS and ASR models. For 255 example, the entry وَلِكَاتِبَاتِهِمَا is transcribed as *wal-* 256 *ikātibātihima*, an accurate phonemic representation

² Pronouncing *ا* as *wa* is a grave mistake committed by at least one of the major engines.

257 to which we have added morphological information. It consists of the stem *kātib* ‘writer’ and the
 258 proclitics *wa* ‘and’ and *li* ‘to’ and the enclitic
 259 *ātihimā* ‘their’. The *ā* indicates two occurrences of
 260 *ā* that have been neutralized to short *a*, indicated by
 261 the underline.

263 3.2 TTS Accuracy

264 Due to the extreme orthographic ambiguity of Arabic, even major IT players struggle to synthesize
 265 speech accurately. The CJKI survey (Halpern,
 266 2020) revealed that it is not unusual for over 50%,
 267 and even 80%, of the words in a sentence, especially cliticized words, to be mispronounced. Errors
 268 are evaluated within the context of a sentence. Errors
 269 are considered erroneous if it includes mistakes such as incorrect case endings (e.g.,
 270 pronouncing الكاتب as *lkātibi* when it should be
 271 *lkātibu*), omitted shaddas (such as pronouncing عدد
 272 as *éádada* when it should be *éáddada* ‘to enumerate’), or other pronunciation errors that can be un-
 273 ambiguously identified. In Table 1, pronunciation
 274 errors are marked by an asterisk.

Unvo- cal- ized	Vo- cal- ized	Google (13%)	iOS (31%)	Bing (25%)	CJKI
عدد	عَدَدٌ	* <i>éádada</i>	* <i>éádad</i> <i>a</i>	* <i>éádad</i> <i>a</i>	<i>éáddad</i> <i>a</i>
الكاتب	الْكَاتِبُ	* <i>lkātibi</i>	<i>lkātibu</i>	<i>lkātibu</i>	<i>lkātibu</i>
ما	مَا	<i>mā</i>	<i>mā</i>	<i>mā</i>	<i>mā</i>
الحكام	الْحُكَمَاءُ	* <i>lhukkā</i> <i>mi</i>	* <i>lhukk</i> <i>āmi</i>	* <i>lhukk</i> <i>āmi</i>	<i>lhukkā</i> <i>ma</i>

279 Now let us look at the errors in context for a
 280 composed text. The original sentence

281 عدد الكاتب ما قال إن هؤلاء الحكام يفعلونه في
 282 الخارج مثل الهجمات الإلكترونية ومطاردة
 283 المعارضين اللاجئين في العواصم الغربية.

284 was mispronounced by Google TTS as

285 **éádada* [*éáddada*] *lkātibu mā qāla* ‘inna ha’ulā’i **lhukkāmi*
 286 [*lhukkāma*] *yaféalunahu fī lkhārīji* **mīthli* [*mīthla*] *lhajamāti*
 287 *l’ilikurūniyyati wamuṣārādati lmuḡriḏīna llaji’īna fī*
 288 *leawāšimi lgharbiyyati*.

289 Asterisks mark incorrectly pronounced words.
 290 The correct pronunciations are given in brackets.

291 A more recent test (December 2021) on *تَصَحَّبُوا*
 292 with Google produced *tašhábūwa*, instead of
 293 *tāshābu*. Not only is word stress incorrect, but the
 294 final *وا* (*wa* + otiose *alif*) (Ryding, 2005), which

296 must be silent, is pronounced (a major error). Table
 297 1 reveals that the error rate for the major TTS en-
 298 gines is notable, highlighting a substantial need for
 299 improvement. The word error rates (WER) in com-
 300 posed texts ranged from 13% to 25%, whereas in
 301 web-extracted texts, they ranged from 70% to 90%.

302 Since the formal tests were conducted, we con-
 303 tinued to conduct informal tests of TTS accuracy,
 304 such as with Google Translate, and did not observe
 305 any significant improvement in accuracy.

306 3.3 Enhancing TTS accuracy

307 The CJKI PATTS (Palestinian Arabic Text to
 308 Speech) white paper (The CJK Dictionary Institute,
 309 2023) presents samples of an early-development
 310 proprietary TTS solution for Palestinian Arabic. It
 311 utilizes phonetic data (specifically IPA) to ensure
 312 that the underlying TTS system generates accurate
 313 realizations. By presenting accurate phonetic data
 314 to a TTS system that supports supplementing such
 315 data, PATTS is able to generate accurate phonetic
 316 realizations without having to fine-tune the under-
 317 lying TTS system to the Palestinian Dialect. One
 318 such system is Amazon AWS Polly, which supports
 319 the “phoneme” SSML tag for its Arabic TTS
 320 voices, which allows the user to specify how a
 321 word should be pronounced using a subset of X-
 322 SAMPA or IPA. (Amazon Web Services, 2023)

323 3.4 ASR Accuracy

324 For TTS, it is necessary to generate one accurate
 325 pronunciation, but ASR systems must recognize al-
 326 ternative pronunciations, including informal ones.

327 For example, the standard pronunciations of كاتبون
 328 ‘writers’ and أكتب ‘I write’ are *katibūna* and
 329 *áktubu*, but the less formal variants *katibūn* and
 330 *áktub* are widespread and possibly even more com-
 331 mon.

332 Such alternatives include pausal forms and final
 333 vowel elision. The former refers to sentence-final
 334 forms causing final vowels to be elided in Classical
 335 Arabic, while the latter is the elision of certain final
 336 vowels in both medial and final forms, common in
 337 spoken MSA. For example, رَجَعْتُ إِلَى الْبَيْتِ ‘I re-
 338 turned home’, pronounced *rajáetu ‘ila_lbayti*, in
 339 pausal form becomes *rajáetu ‘ila_lbayt* and in
 340 spoken MSA becomes *rajáet ‘ila_lbayt*. Note how
 341 the final *ti* and *tu* are truncated to *t*.

342 The above alternatives are for standard MSA.
 343 There are also regional allophones. For example, /j/
 344 in words such as *jamal* ‘camel’ is pronounced [g]

345 in Egypt, [dʒ] in the Gulf region, and [ʒ] in the Le-
346 vant. These are regional variants of MSA. Arab-
347 LEX not only includes the IPA for the standard
348 MSA, namely [dʒ] for /j/, but also the regional al-
349 lophones [ʒ] and [g]. It aims to include transcrip-
350 tions of common non-standard regional allo-
351 phones.

352 In this context, applying a phonetic alphabet like
353 IPA or X-SAMPA becomes particularly relevant, as
354 utilizing phonetics in ASR systems has proven ben-
355 efiticial (Feng et al. 2023).

356 4 Machine Translation

357 Although NMT has dramatically improved transla-
358 tion quality, it has some shortcomings, as Philipp
359 Koehn (2020) and Halpern (2018) pointed out.
360 Some issues in Arabic are (1) the high orthographic
361 ambiguity, (2) the morphological complexity
362 (forms like *ولكاتباتها* are difficult to analyze), (3)
363 the recognition of named entities (often cliticized),
364 and (4) a large number of word forms for nouns and
365 verbs.

366 ArabLEX offers comprehensive coverage of in-
367 flected and cliticized forms and can be used to sup-
368 plement existing corpora or as a pseudo-corpus for
369 language model training, enhancing the accuracy
370 of morphological, syntactic, and semantic analysis.
371 Additionally, the proper noun modules of Arab-
372 LEX – DAN (Database of Arabic Names), DAF
373 (Database of Arabic Foreign Names), and DAP
374 (Database of Arabic Places) – are bilingual and ro-
375 manized, serving as a bilingual dictionary.

376 5 ArabLEX in Action

377 5.1 Scope and Coverage

378 The first release of ArabLEX in 2021 covered
379 about 530 million entries for general vocabulary
380 and proper nouns. ArabLEX consists of the follow-
381 ing four main modules: DAG (Arabic General Vo-
382 cabulary, 83M entries), DAN (Arabic Names,
383 218M entries), DAF (Arabic Foreign Names,
384 226M entries) and DAP (Arabic Place Names, 6M
385 entries). ArabLEX has 30 data fields with detailed
386 grammatical, phonological, morphological, and or-
387 thographic attributes (Halpern, 2020).

388 It can be argued that generating entries by rules
389 and templates can result in a large number of non-
390 existing or erroneous forms. We have taken ex-
391 treme care to ensure that only grammatically and,

392 as far as possible, semantically valid forms are in-
393 cluded. Though currently some forms may not
394 have been observed to exist, they are indeed valid
395 and could occur in the future. For most applica-
396 tions, the negative effects are negligible compared
397 to a lack of data (Koperniak, 2017).

398 5.2 ArabLEX Compared to Other Re- 399 sources

400 The comprehensiveness of a dataset holds signifi-
401 cant importance for NLP applications. This is par-
402 ticularly pronounced in morphological engines,
403 (Attia et al., 2011).

404 Previous efforts to compile extensive Arabic
405 lexicographical or lexical databases have yielded
406 datasets containing around 200,000 unique lemma
407 entries. These datasets tend to lack a diverse set of
408 attributes (Attia et al., 2011; Alshargi et al., 2019).
409 In contrast, the CALIMA dataset for Egyptian Ar-
410 abic comprises 36,000 distinct lemmata and con-
411 tains 40 fields with attributes such as morphology,
412 gender, and root (AlShuhayeb, 2023; Habash et al.,
413 2012). Detailed datasets like this typically contain
414 entries in the range of 30,000 headwords (Alshargi
415 et al., 2019).

416 ArabLEX, on the other hand, covers a combined
417 375.335 unique lemmata, including a large number
418 of named entities, while exceeding the level of de-
419 tail and versatility of its counterparts. Especially by
420 offering phonetic (IPA, XSAMPA) and morpho-
421 phonemic (CARS) transcriptions and fully diacrit-
422 ized Arabic, ArabLEX fills a gap in current lexi-
423 cal resources.

424 Another key difference is the total number of en-
425 tries accessible for explicit analysis; that is, entries
426 that are pre-generated as opposed to on-the-fly. The
427 Calima dataset contains approximately 48 million
428 entries that can be examined when all its lemmata
429 and affixes are exhaustively generated
430 (AlShuhayeb, 2023). By contrast, ArabLEX con-
431 sists of 530 million entries pre-compiled in TSV
432 format immediately accessible for use and analysis.

433 5.3 Comparison with CALIMA Star

434 ArabLEX’s model of Arabic morphology is
435 more refined than those of other systems. This re-
436 sults in high recall by covering almost all word
437 forms. Multiple layers of sanity-checking ensure
438 high precision and grammatical validity of each
439 form. To illustrate this, we compared some features
440 of ArabLEX and CALIMA Star (“Calima” below),
441 the most advanced morphological engine, using the

442 affirmative of the verb كَتَبَ ‘to write’. The results
443 are based on the Calima generator web interface.

444 (1) The coverage of inflected and cliticized
445 forms differs dramatically. Many conjugated forms
446 are missing in Calima, which also generates some
447 invalid forms. The table below shows the number
448 of forms for كَتَبَ.

Item	CALIMA Star	ArabLEX
Total forms	2,448	5,886
Uncliticized	104	124
Cliticized	2,344	5,762

Table 2: Coverage CALIMA Star vs. ArabLEX.

449

450 (2) The cliticized forms كَتَبْنَا, كَتَبْتَنِي, and كَتَبْتَنِكَ
451 are not given by Calima, whereas some forms it
452 provides, like لَا يَكْتُبُ, are grammatically invalid.
453 The number of cliticized forms provided by Ara-
454 bLEX (both proclitics and enclitics) for كَتَبَ ex-
455 ceeds that of Calima by 146%.

456 (3) The results of a preliminary investigation of
457 proclitic coverage by Calima (expanded on below)
458 shows that Calima does not support the proclitic
459 />a/ (أَ), even if selected from the menu. ArabLEX
460 provides more clitic combinations: 39 proclitic
461 combinations and over 2000 (to our knowledge
462 double that of Calima) proclitic-enclitic combina-
463 tions, which were carefully vetted to ensure their
464 validity. For example, the singleton proclitic se-
465 quence />awabi{lo/ is a valid combination for
466 nouns, but />awaka{lo/ is not, while any proclitic
467 in />a, wa, fa, >awa, >afa/ can combine with any
468 enclitic in /N, FA, FY/ for singular nouns.

469 (4) Whereas ArabLEX takes great care to in-
470 clude only grammatically valid forms, Calima
471 seems to generate agrammatical forms such as
472 لَا أَكْتُبُ and سَأَكْتُبُ, or invalid forms such as
473 لَا أَكْتُبُ (omitting the space after لَا).

474 (5) The verb conjugation paradigm is missing
475 important forms. For example, Calima does not re-
476 turn the active participle كَاتِبٌ, nor the passive par-
477 ticiples مَكْتُوبٌ for the verb lemma كَتَبَ.

478 (6) The imperative forms أَكْتُبْ, أَكْتُبِي, etc. are
479 not generated even when explicitly requested via
480 the user interface, which is a serious shortcoming.

481 In conclusion, the ArabLEX morphological
482 model of the verb كَتَبَ is significantly finer than

483 that of Calima. Whereas the former always pro-
484 vides grammatically valid forms, the latter some-
485 times generates agrammatical ones.

486 5.4 Real-World Applications

487 Amazon has acknowledged the significant contri-
488 bution of ArabLEX to its advanced Arabic speech
489 technology for Alexa. It’s comprehensiveness and
490 the in-depth morphological and phonological data
491 has helped Amazon reduce the error rates for both
492 recognition and generation; that is, to recognize Ar-
493 abic queries, including place and personal names,
494 as well as generate answers with greater precision
495 (The CJK Dictionary Institute, 2022).

496 5.5 Grammatical Attributes

497 The grammatical attributes of ArabLEX are useful
498 for morphological analysis, orthographic disam-
499 biguation, POS tagging, semantic analysis, and
500 more. These include codes for gender, number,
501 case endings and person, as well as the stem, defi-
502 niteness, lexical rationality, and the lemma.

Data field	Value
Full-form	وَلِكَاتِبِكُمْ
Lemma	كَاتِبٌ
Stem	كَاتِب
Gender	C (common)
Case	GEN (genitive)
Number	D (dual)
Person	2 (second)
Definiteness	D (definite)
Root	ك-ت-ب

Table 3: Grammatical attributes.

503

504 5.6 Phonological Attributes

505 The phonemic and phonetic transcriptions are use-
506 ful for improving speech technology, both TTS and
507 ASR (Tahon *et al.*, 2016; Feng *et al.*, 2023). These
508 include precise, fully diacriticized Arabic with ac-
509 curate phonemic and phonetic transcriptions as
510 well as word stress and vowel neutralization. The
511 main phonological attributes are shown in Table 4.

Data field	Value
Vocalized	مُحَمَّدٌ
Phonemic	muhammadun
Phonetic	[muˈhəmmədun]
X-SAMPA	muˈXE "mmE "dun
Transliterated	muham~adN

Table 4: Phonological attributes for محمد.

512

5.7 Morphological/Orthographic Attributes

The morphological attributes include all inflected, conjugated, declined, and cliticized word forms, such as plurals, duals, feminine, case endings, conjugated forms, as well as proclitics, enclitics, stems, and roots. They are useful for morphological analysis, semantic analysis, lemmatization, decliticization, deaffixation, verb conjugation, and dictionary lookup. Operations such as decliticization, deaffixation and tokenization (Carbonell et al., 2006) are easy to perform since clitics are given explicitly in their own fields (Enclitic, Proclitic, and Stem below). The main morphological attributes are shown in Table 5.

Data Field	Value	Transcription
Full-form	ولكاتبكما	walikātibikumā
Lemma	كاتب	kātibun
Stem	كاتب	kātib
Proclitic	ول	wali
Enclitic	كما	(i)kūmā
Root	ك-ت-ب	k-t-b

Table 5: Morphological attributes.

Orthographic attributes are useful for orthographic disambiguation, which is necessary for word and entity recognition, TTS, morphological analysis, word/entity extraction, normalization, and dictionary lookup. These include orthographic variants such as pausal and elided forms and even common typographical oddities. Here is an example of typical orthographic variants for the name Alexandra: الكسندرا، ألكسندرا، ألكسندره، ألكسندرة، ألكسندرة. As shown above, ة and ة are sometimes interchangeable in names. Orthographic variants also include allographs, for example the use of ي (alif maqsuura) as an alternative for ي (yaa) in Egypt, and the use of پ instead of ب for [p] in some regions.

5.8 Named Entity Recognition

The DAN module of ArabLEX covers about 100,000 vocalized personal names and their 6.5 million romanized variants. DAN is widely deployed in both security and NLP processing tools for NER and MT. Similarly, the DAF and DAP modules consist of about 240,000 names for places and non-Arab personal names. These modules account for about 450 million fully inflected and cliticized entries in ArabLEX (Halpern, 2009).

5.9 Accessing ArabLEX

ArabLEX is specifically targeted at researchers and software developers needing rich morphological

and phonological resources. It is available through The CJK Dictionary Institute and the European Language Resource Association, a non-profit repository of language resources (European Language Resource Association, 2022).

6 Compilation Methods

6.1 Quality Control

The ArabLEX team, comprising professional editors, translators, computational linguists, and university instructors, has conducted extensive research to ensure maximum accuracy and comprehensive coverage of all word forms and their variants. Many programs were developed for data validation and proofreading to ensure accuracy and consistency, such as programs for automatic error detection and correction and data validation. The following outlines one of the data validation processes our team employs to refine our vocalization validation module (VBW_INTEG) to ensure the accuracy of fully vocalized Arabic and phonemic transcriptions, critical for speech technology:

(1) A program validates the fact that inflections are correctly vocalized based on strictly defined rules such as *hamza* rules, presence of short vowels and many more. (2) The program then attempts to rectify the errors it encounters autonomously. (3) Errors that the program cannot rectify are presented to our proofreaders, who manually classify, analyze, and rectify them. (4) Based on the feedback of our proofreaders, the validation rules are then either adjusted or our database of exceptions is expanded. (5) The process is then repeated.

This iterative process has been applied over the course of many years, resulting in a system with a comprehensive set of rules and exceptions.

To illustrate, when validating one of our Arabic dictionaries using the same program, we identified entries such as /diywaAnN/ (ديوان) and /\$>owN/ (شأؤ). Subsequently, the spellings were automatically corrected to /diywaAnN/ (ديوان) and /\$a>owN/ (شأؤ). Our proofreaders then reviewed these modifications and confirmed them as correct.

Note that this process has refined our rule base to be highly sophisticated. True exceptions are rather rare, normally one-off isolated instances. If a trend or pattern is found in the exceptions, they are analyzed and codified as rules in the error detection program so they will no longer be considered errors. The accuracy of phonemic transcriptions is likewise ensured, as it uses the fully vocalized Arabic generated by this process and undergoes a similar validation process.

6.2 Inflection, Conjugation, Cliticization

Generating inflected forms involves many complex steps, including sanity checking and human proofreading. Nouns and adjectives are declined/inflected for feminine, dual, and plural forms. For example, for /bayotN/ ‘house,’ we derive /bayotaAni/, /buyuwtN/, and /buyuwtaAtN/. As for conjugation, the verb paradigms from the CJKI Arabic Verb Conjugator (CAVE) (The CJK Dictionary Institute, 2011) are used to acquire the verb conjugations for each subject pronoun for each tense. CAVE has 180 categories and fully explicit conjugated paradigms (generated by hand-vetted precise rules and exceptions) for each category. For example, for /kataba/ ‘he wrote’ we get /yakotubu/ (third person masculine singular imperfect), /Aukotubo/ (second person masculine singular imperative), etc. To encliticize, the correct enclitic template is selected based on the ending of the inflected form. For example, the noun /xirapu/ ‘the hereafter’ ends in /pu/, so the template in Table 6 is selected. Enclitics are then added to correspond to each case and subject pronoun. For /xirapu/, we generate such forms as /xiratiy/, /xiratuka/ and /xiratuki/. To procliticize, the appropriate proclitics are elected from the template. For example, for /bayotN/ ‘house’, the enclitic is /-N/ (tanwiin), so we refer to the appropriate row (row 2) in Table 7 and generate />abayotN/, /wabayotN/, etc.

Note that the clitics are not merely blindly concatenated to the base form – there are over 2,000 valid orthographic, grammatical, and semantic combinations of clitics that are defined by our human-vetted constraint-defining tables, as shown in Table 7, and several thousand that are invalid.

Per	Case	Enclitic	Rule
000	NOM	u	
1SC	NOM	iy	-p → -t
2SM	NOM	uka	-p → -t
2SF	NOM	uki	-p → -t

Table 6: Template for nouns that end in /p/ (ð).

Proclitic	Enclitic	Gen	Num
0,>a,wa,fa,>aw a,>afa,Aalo,...	a,u	M	S
0,>a,wa,fa,>aw a,>afa	N,FA,FY	M	S
0,>a,wa,fa,>a wa,>afa	uhaA,uhu,uhumaA, uhumo, uhun~a,uka, uki,ukumaA,...	M	S

Table 7: Possible combinations of clitics.

7 Future Work

We will continue expanding ArabLEX by adding new entries and data fields, including technical terms, and named entities, as well as more phonological attributes, orthographic variants, alternative pronunciations, and additional word classes (POS). Especially noteworthy are new headwords that consist of multiword expressions (Halpern, 2019) (inflections or conjugations consisting of space-delimited components), such as periphrastic elatives (أَقْلُّ أَكْثَرُ إِيلَام) ‘more painful’), negative elatives (with أَقْلُّ or أَخْفُّ), inflected numerical expressions, phrasal verbs, compound tenses, verb negation, and more. The addition of proclitics, enclitics and inflections lead to ArabLEX exceeding 500 million records (15 billion data points). It is expected to reach about one billion records in the near future.

In parallel to ArabLEX, we have been developing a series of full-form lexicons for the major Arabic dialects, called DiaLEX, based on the same methodology used for ArabLEX. Since there is no official orthography for the dialects, we conducted thorough research on various dialects by analyzing corpora and dictionaries and by collaborating with native-speaking experts. As a result, we have identified the most common conventions for each dialect, which made the creation of DiaLEX possible. DiaLEX currently (May 2024) covers the major Arabic dialects Egyptian, Emirati and Hijazi. The initial release of the first three has been completed, covering about 150 million entries, and the development of a Palestinian full-form lexicon (PA_LEX) is now in progress (May 2024).

ArabLEX, in combination with DiaLEX, can serve as a holistic resource for the development of NLP applications for MSA and its dialects.

8 Bibliographical References

- Abdelali, A., Darwish, K., Durrani, N., and Mubarak, H. (2016). *Farasa: A fast and furious segmenter for Arabic*. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Demonstrations*, pages 11–16, San Diego, California, June. Association for Computational Linguistics. <https://doi.org/10.18653/v1/N16-3003>
- Algihab, W., Alawwad, N., Aldawish, A., and Al-Humoud, S. (2019). *Arabic Speech Recognition with Deep Learning: A Review*. In G. Meiselwitz (Ed.) *Social Computing and Social Media. Design, Human Behavior and Analytics*, Springer International Publishing (pages 15–31). https://doi.org/10.1007/978-3-030-21902-4_2
- Alshargi, F., Dibas, S., Alkhereyf, S., Faraj, R., Abdulkareem, B., Yagi, S., Kacha, O., Habash, N., & Rambow, O. (2019). *Morphologically Annotated Corpora for Seven Arabic Dialects: Taizi, Sanaani, Najdi, Jordanian, Syrian, Iraqi and Moroccan*. In *Proceedings of the Fourth Arabic Natural Language Processing Workshop* (pp. 137–147). <https://doi.org/10.18653/v1/W19-4615>.
- AlShuhayeb, H., Minaei-Bidgoli, B., Shenassa, M. E., Hossayni, S (2023). *Noor-Ghateh: A Benchmark Dataset for Evaluating Arabic Word Segmenters in Hadith Domain*. arXiv preprint arXiv:2307.09630. <https://doi.org/10.48550/arXiv.2307.09630>
- Amazon Web Services, Inc. (2023). *Amazon Polly Developer Guide*. <https://docs.aws.amazon.com/polly/>
- Attia, M., Pecina, P., Toral, A., Tounsi, L., & Genabith, J. (2011). *An Open-Source Finite State Morphological Transducer for Modern Standard Arabic*. In *Proceedings of the Fourth International Conference on Arabic Language Resources and Tools (MEDAR)*, pp. 125–133. <https://aclanthology.org/W11-4417>
- Boudchiche, M., Mazroui, A., Ould Abdallahi Ould Bebah, M., Lakhouaja, A., and Boudlal, A. (2017). “AlKhalil Morpho Sys 2: A robust Arabic morpho-syntactic analyzer,” *Journal of King Saud University – Computer and Information Sciences* 29 (2), 141–146. <https://doi.org/10.1016/j.jksuci.2016.05.002>
- Buckwalter, T. (2002). *Buckwalter Arabic Morphological Analyzer Version 1.0*. Linguistic Data Consortium, University of Pennsylvania, 2002. LDC Catalog No.: LDC2002L49. <https://doi.org/10.35111/7vzm-mb15>
- Carbonell, J., Klein, S., Miller, D., Steinbaum, M., Grassiany, T., and Frei, J. (2006). *Context-Based Machine Translation*. In *Proceedings of the 7th Conference of the Association for Machine Translation in the Americas: Technical Papers*, pages 19–28, Cambridge, Massachusetts, USA, August. Association for Machine Translation in the Americas. <https://aclanthology.org/2006.amta-papers.3/>
- CJK Dictionary Institute, The (2011). The CJKI Arabic Verb Conjugator. Downloaded from <http://cjk.org/arabic/cave/cavehelp.htm> on 13 January 2022.
- CJK Dictionary Institute, The (2020). ArabLEX Technical Specifications. Downloaded from https://www.cjk.org/wp-content/uploads/Arablex_specs.pdf on 13 January 2022.
- CJK Dictionary Institute, The (2020). Enhancing Arabic Speech Technology with comprehensive Arabic training lexicon. Downloaded from https://www.cjk.org/wp-content/uploads/TTS_Report.pdf on 13 January 2022.
- CJK Dictionary Institute, The (2022). ArabLEX Arabic Full-Form Lexicon. Retrieved from <https://www.cjk.org/data/arabic/nlp/arablex-arabic-full-form-lexicon/>
- CJK Dictionary Institute, The (2023). *Palestinian Arabic Text-to-Speech system (PATTS)*. Retrieved from <https://www.cjk.org/wp-content/uploads/patts.html>
- Crystal, D. 2008. *A Dictionary of Linguistics and Phonetics*. New York: Wiley-Blackwell. ISBN 9781405152969. <https://doi.org/10.1002/9781444302776>.
- Diab, N. (2021). *Out of the BLEU: An Error Analysis of Statistical and Neural Machine Translation of WikiHow Articles from English into Arabic*. *CDELTA Occasional Papers in the Development of English Education*. 75. 181–211. <https://doi.org/10.21608/opde.2021.208437>.
- European Language Resources Association (ELRA). (2022). *The Multilingual European Language Corpus* [Datasets: ELRA-L0131, ELRA-M0105, ELRA-M0106 and ELRA-M0107]. ELRA. <https://catalog.elra.info>
- Feng, S., Tu, M., Xia, R., Huang, C., & Wang, Y. (2023). *Language-universal phonetic encoder for low-resource speech recognition*. arXiv. <https://arxiv.org/abs/2305.11576>
- Gadelli, N. 2015. *The morphological integration of loanwords into Modern Standard Arabic: Towards a morphological categorization of loanwords*. Bachelor thesis in General Linguistics, Lund University, Sweden. <https://lup.lub.lu.se/luur/download?func=downloadFile&recordId=5053170&fileId=5053172>
- Habash, N., Rambow, O., and Roth, R. (2009). *MADA+TOKAN: A toolkit for Arabic tokenization, diacritization, morphological disambiguation, POS*

- 784 150 tagging, stemming and lemmatization. In *Proceedings of the Second International Conference on*
785 *Arabic Language Resources and Tools*, Cairo, Egypt, April. The MEDAR Consortium. [https://api.semanticscholar.org/Cor-](https://api.semanticscholar.org/CorpusID:10124099)
786 [pusID:10124099](https://api.semanticscholar.org/CorpusID:10124099)
- 790 Habash, N., Eskander, R., Hawwari, A. (2012). **A Morphological Analyzer for Egyptian Arabic**. In *Proceedings of the Twelfth Meeting of the Special Interest Group on Computational Morphology and Phonology*, pages 1–9, Montréal, Canada. Association for Computational Linguistics. <https://aclanthology.org/W12-2301>
- 797 Halpern, J. (2002). The Challenges and Pitfalls of Arabic Romanization and Arabization. Downloaded from <https://www.cjki.org/arabic/arannana.pdf> on 13 January 2022.
- 801 Halpern, J. (2008). **Exploiting Lexical Resources for Disambiguating CJK and Arabic Orthographic Variants**. In *Proceedings of the Sixth International Conference on Language Resources and Evaluation (LREC'08)*, Marrakech, Morocco, May. European Language Resources Association (ELRA). [http://www.lrec-conf.org/proceed-](http://www.lrec-conf.org/proceedings/lrec2008/pdf/109_paper.pdf)
802 [ings/lrec2008/pdf/109_paper.pdf](http://www.lrec-conf.org/proceedings/lrec2008/pdf/109_paper.pdf)
- 809 Halpern, J. (2009). **Word stress and vowel neutralization in modern standard Arabic**. In *Proceedings of the Second International Conference on Arabic Language Resources and Tools*, Cairo, Egypt, April. The MEDAR Consortium. <http://www.elda.org/medar-conference/pdf/16.pdf>
- 815 Halpern, J. (2009). **CJKI Arabic Romanization System (CARS)**. Downloaded from https://www.cjki.org/cjk/arabic/cars/cars_paper.pdf on 13 January 2022.
- 819 Halpern, J. (2009). **Lexicon-Driven Approach to the Recognition of Arabic Named Entities**. In *Proceedings of the Second International Conference on Arabic Language Resources and Tools*, Cairo, Egypt, April. The MEDAR Consortium. <https://api.semanticscholar.org/CorpusID:43949907>
- 825 Halpern, J. (2016). **Compilation Techniques for Pedagogically Effective Bilingual Learners' Dictionaries**. *International Journal of Lexicography*, Volume 29(3): 323–338. <https://doi.org/10.1093/ijl/ecw023>
- 829 Halpern, J. (2018). **Very large-scale lexical resources to enhance Chinese and Japanese machine translation**. In *Proceedings of the eleventh international conference on language resources and evaluation (LREC 2018)*, Miyazaki, Japan, May. European Language Resources Association (ELRA). <https://aclanthology.org/L18-1137>
- 836 Halpern, J. (2019). **Lexicographic Criteria for Selecting Multiword Units for MT Lexicons**. Downloaded from https://www.cjki.org/cjk/reference/Lexicographic_criteria_Halpern.pdf on 13 January 2022.
- 840 Halpern, J. (2020). **Enhancing Arabic Speech Technology with Comprehensive Arabic Training Lexicon**. Downloaded from https://www.cjk.org/wp-content/uploads/TTS_Report.pdf on 07 August 2023.
- 844 Haugen, E. (1972). **Schizoglossia and the Linguistic Norm**, in *Studies by Einar Haugen*, Berlin, Boston: De Gruyter Mouton, pages 441–445. doi: <https://doi.org/10.1515/9783110879124.441>
- 848 International Phonetic Association (1999). **Handbook of the International Phonetic Association: A Guide to the Use of the International Phonetic Alphabet**. Cambridge University Press. doi: <https://doi.org/10.1017/9780511807954>
- 853 Koehn, P. (2020). **Neural Machine Translation**. Cambridge University Press, Cambridge, UK. doi: <https://doi.org/10.1017/S1351324920000650>
- 856 Koperniak, S. (2017). **Artificial data give the same results as real data — without compromising privacy**. Institute for Data, Systems, and Society. Retrieved from <https://news.mit.edu/2017/artificial-data-give-same-results-as-real-data-0303>
- 861 Maamouri, M., Bies, A., Buckwalter, T., and Mekki, W. (2004). **The Penn Arabic Treebank: Building a Large-Scale Annotated Arabic Corpus**. In *NEMLAR Conference on Arabic Language Resources and Tools*, pages 102–109, Cairo, Egypt. https://www.researchgate.net/publication/228693973_The_penn_arabic_treebank_Building_a_large-scale_annotated_arabic_corpus
- 870 Mejdell, G. (2014). **Lugát al-’umm and al-luga al-’umm - the ‘mother tongue’ in the Arabic context**, in *Arabic and Semitic Linguistics Contextualized*, Harrassovitz, pages 214–226. <https://doi.org/10.2307/j.ctvc2rmgq.16>
- 875 Pasha, A., Al-Badrashiny, M., Diab, M., El Kholy, A., Eskander, R., Habash, N., Pooleery, M., Rambow, O., and Roth, R.M. (2014). **MADAMIRA: A Fast, Comprehensive Tool for Morphological Analysis and Disambiguation of Arabic**. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, Reykjavik, Iceland, May. European Language Resources Association (ELRA). http://www.lrec-conf.org/proceedings/lrec2014/pdf/593_Paper.pdf
- 885 Ryding, K.C. (2005). **A Reference Grammar of Modern Standard Arabic**, Cambridge University Press, Cambridge, UK. doi: <https://doi.org/10.1017/CBO9780511486975>
- 889 Smrž, O. (2007). **ElixirFM — Implementation of Functional Arabic Morphology**. In *Proceedings of the 2007 Workshop on Computational Approaches to*

- 892 *Semitic Languages: Common Issues and Resources*,
893 pages 1–8, Prague, Czech Republic, June. Association
894 for Computational Linguistics. <https://aclanthology.org/W07-0801>
895
- 896 Soudi, A., Eisele, A. (2004). *Generating an Arabic*
897 *Full-form Lexicon for Bidirectional Morphology*
898 *Lookup*. In *Proceedings of the Fourth International*
899 *Conference on Language Resources and Evaluation*
900 *(LREC'04)*, Lisbon, Portugal, May. European Lan-
901 *guage Resources Association (ELRA)*.
902 [http://www.lrec-conf.org/proceed-](http://www.lrec-conf.org/proceedings/lrec2004/pdf/567.pdf)
903 [ings/lrec2004/pdf/567.pdf](http://www.lrec-conf.org/proceedings/lrec2004/pdf/567.pdf)
- 904 Taji, D., Khalifa, S., Obeid, O., Eryani, F., and Habash,
905 N. (2018). *An Arabic Morphological Analyzer and*
906 *Generator with Copious Features*. In *Proceedings of*
907 *the Fifteenth Workshop on Computational Research*
908 *in Phonetics, Phonology, and Morphology*, pages
909 140-150, Brussels, Belgium, October. Association
910 for Computational Linguistics. [https://aclanthol-](https://aclanthology.org/W18-5816/)
911 [ogy.org/W18-5816/](https://aclanthology.org/W18-5816/)
- 912 Wells, J.C, (1997). SAMPA computer readable pho-
913 netic alphabet. In Gibbon, D., Moore, R. and Win-
914 ski, R. (eds.), 1997. *Handbook of Standards and Re-*
915 *sources for Spoken Language Systems*. Berlin and
916 New York: Mouton de Gruyter. Part IV, section B.
917 <https://www.phon.ucl.ac.uk/home/sampa/>
- 918 Tahon, M., Qader, R., Lecorvé, G., Lolive, D. 2016.
919 *Improving TTS with corpus-specific pronunciation*
920 *adaptation*. *Interspeech*, San Francisco, United
921 States [https://www.isca-archive.org/inters-](https://www.isca-archive.org/interspeech_2016/tahon16_interspeech.pdf)
922 [speech_2016/tahon16_interspeech.pdf](https://www.isca-archive.org/interspeech_2016/tahon16_interspeech.pdf)

923