



# ArabLEX Comprehensive Arabic Full Form Lexicon

معجم اللغة العربية الكامل by Jack Halpern

**The CJK Dictionary Institute** (CJKI) is pleased to announce the release of the Arabic Full Form Lexicon, or *ArabLEX*. Covering approximately **530 million entries,** this is the most comprehensive Arabic computational lexicon ever created. *Full form* means that it includes all inflected forms. It covers not only general vocabulary, but also, for the first time, fully inflected proper nouns.

*ArabLEX* is, quite literally, the ultimate resource for Arabic NLP and AI, ideally suited for such applications as morphological analysis, machine translation, speech technology, deep learning, and cybersecurity. No other Arabic lexicon comes close to it the scope, coverage and comprehensiveness.

*ArabLEX* is rich in morphological, grammatical, phonological, and orthographical attributes (currently about 30). In addition, it maps all unvocalized forms to their vocalized counterparts and to the lemma, and provides precise phonemic and phonetic transcriptions. These features can significantly contribute to the training of language models for NLP and AI applications.

The **Burj Khalifa** (*Alburj*), soaring at 832 meters above the skies of Dubai, is the tallest structure ever created. Similarly, *ArabLEX* is the largest Arabic lexicon

1

ever created. We thus chose Alburj as an emblem to epitomizes the essence of *ArabLEX*.

This unparalleled computational lexicon, the fruit of nearly a decade of intense development and validation, is now available to the NLP, AI, and cybersecurity communities for research and product development

### 1. Full Form Lexicon

#### 1.1 What is a Full Form Lexicon?

A *full form lexicon* is a computational lexicon that contains all inflected, conjugated, declined, and cliticized forms that occur in a language (referred to as *wordforms*). Unlike ordinary dictionaries, which include only the canonical forms (base lexemes), a full form lexicon includes all wordforms. For example, the full set of wordforms for *eat* includes *eating*, *eaten* and *ate*, while for *boy* it includes *boys*, *boy's* and *boys'*. Arabic morphology is more complicated. Adding the *proclitics y* wa 'and' j li 'to' and the *enclitic* lixed *tihima* to the stem example, while for 'walikatibatihima 'and to the two female writers'.

In English, the number of such forms is quite limited, but an Arabic word can have thousands of inflected and cliticized forms. For example, the verb *kataba* has about 7250 forms (by comparison Japanese has about 2500), whereas the noun كَاتِب *kātib* has about 5270 forms. As a result, the number of entries in *ArabLEX* reaches about 530 million.

#### 1.2. Why a full form lexicon?

Traditionally, MT and other NLP applications have been (and some still are) based on rules (RBMT) or on statistical models (SMT). Recently, neural machine translation (NMT) is becoming the norm. Despite of the significant improvements that NMT has brought about, this new technology still has some

shortcomings, such as the handling of proper nouns and multiword expressions (MWE), as described in Halpern's papers on large-scale lexical resources [1] and MWEs [8]. An important issue in Arabic speech technology is the numerous complex morphological forms, like وَلِكَاتِبَاتِهِمَ *walikatibātíhima*, and the high level of orthographic ambiguity (due to the lack of vowels, as in رولكاتباتهما).

A full form lexicon can significantly contribute to the quality of Arabic MT and speech technology (both synthesis and recognition) by mapping unvocalized to vocalized forms, by providing detailed morphological information, and by providing phonemic/phonetic transcriptions for all wordforms.

## 2. Arabic Full Form Lexicon (ArabLEX)

*ArabLEX* is a full form Arabic lexicon that provides comprehensive coverage for inflected, conjugated and cliticized forms, and includes a rich set of attributes for natural language processing.

#### 2.1 Distinctive features

Various features of *ArabLEX* offer special benefits to developers of Arabic NLP and AI applications, especially speech technology and machine translation.

- Created by a team of specialists in Arabic morphology and computational lexicography
- The first release in early 2021 covers about **530 million** full form entries, including proper nouns.
- Includes all inflected, conjugated, declined, and cliticized wordforms, including plurals, dual, feminine, case endings, conjugated forms, as well as proclitics, enclitics, stems and roots.
- Unvocalized and precisely fully vocalized Arabic
- Accurate phonemic transcriptions and IPA for all entries

- Millions of orthographic variants for both vocalized and unvocalized Arabic
- A large variety of grammatical codes include part-of-speech, person, gender, and case codes -- currently about 30 attributes for each entry.
- All wordforms are cross-referenced to their lemma (canonical form)
- Constantly maintained and expanded

### 2.2. Modules and Subsets

The full set of *ArabLEX* consists of the following major four modules and several submodules

The full version of *ArabLEX* consists of four major modules:

DAG	Database of Arabic General Vocabuary	83 million entries
DAN	Database of Arabic Names	218 million entries
DAF	Database of Arabic Foreign Names	226 million entries
DAP	Database of Arabic Place Names	6 million entries

Each module can be provided in four different subsets for specific applications. DAX represents any of the four *ArabLEX* modules. Each of these modes can be fully customized to specific requirements.

DAX-XP	Excludes proclitics, which are ten times as numerous as enclitics and may be unnecessary
DAX-PH	Includes phonemic and phonetic attributes fine tuned for speech technology, such as IPA, SAMPA, vocalization, and CARS, a specially designed phonemic transcription.
DAX-XC	Excludes all clitics, inflections and declensions (canonical forms) but can includes plural and dual forms.

DAX-WL	A bare bone wordlist, with or without inflections, and a
	minimum set of attributes, such POS codes.

### 2.3 Compilation History

For about a decade, our team of computational lexicographers and language specialists have been engaged in the development of full form lexicons for Spanish, Arabic, and Japanese. To this end, we analyzed the grammar and morphology of these languages in great depth, to a degree well beyond comprehensive descriptive grammars for these languages. Initially we focused on Spanish (bilingual SFULEX) and Japanese (monolingual JFULEX) full form lexicons, which have significantly contributed to MT technology, such as support for a Spanish-English MT system that achieved a "human quality" translation [7].

In the last couple of years we have intensified our efforts to expand, proofread and validate *ArabLEX*, with the aim of covering nearly all wordforms in Modern Standard Arabic, which contains millions of full form proper nouns in bilingual format.

### 3. Enhancing Speech Technology

The quality of Arabic TTS lags considerably behind that of the major languages. One reason is that the Arabic script is highly ambiguous. For example, كاتب can represent as many as *seven* pronunciations. In addition, the complexity of such cliticized forms as ولكاتباتهما *walikatibātíhima*, and the absence of vowels, makes Arabic TTS especially challenging.

The extreme orthographic ambiguity of Arabic has led to unacceptably high error rates. In a survey we discovered that sometimes over 50%, and even

80%, of the words in a sentence are mispronounced. The time has come to dramatically improve the intolerably low quality of Arabic TTS.

*ArabLEX*, which maps all unvocalized wordforms, including all cliticized forms, to their vocalized counterparts and provides phonemic and phonetic transcriptions that include precise word stress and even vowel neutralization. These features can help developers significantly enhance the quality of Arabic TTS by training ASR systems to achieve higher recognition rates. A module specifically designed for ASR developers provides phonemic and phonetic variants necessary for recognition but not for synthesis (such as *katibūn* for Street in addition to the standard MSA *katibūna*).

To summarize, ArabLEX can bring the following benefits to speech technology:

- Hundreds of millions of full-form entries, including millions of proper nouns.
- Covers *all combinations* of proclitics and enclitics for inflected wordforms.
- Tens of millions of orthographic variants.
- Exhaustively lists alternative pronunciations for orthographical disambiguation
- Future versions will provide 'importance flags' to help indicate the most likely alternative.
- Highly accurate phonemic transcriptions for all wordforms, including precise stress and vowel neutralization.

### 3.1 Orthographical ambiguity

One reason that Arabic speech technology lags behind is that the Arabic script is highly ambiguous. Words are often written as a string of consonants with no indication of vowels. For example, کاتب can represent as many as *seven* 

pronunciations: *kāatib*, *kātibun*, *kātibin*, *kātaba*, *kātibi*, *kātiba* and *kātibu*. Many other characteristics of the Arabic script contribute to a high level of orthographic ambiguity, as described in Halpern's paper on Arabic named entities [5].

The morphological complexity of such cliticized forms as ولكاتباتهمـ*walikatibātíhima*, and the absence of vowel diacritics, makes Arabic TTS especially challenging. That is, determining the morphological composition of such forms, and the correct vowels for such consonants as ولكاتباتهما in ت often requires morphological, semantic and contextual analysis which tax the capabilities of state-of-the-art speech technology.

#### 3.2 Improving TTS accuracy

#The extreme orthographic ambiguity of Arabic has led to unacceptably high error rates, even by the TTS systems offered by major players such as Google, Apple and Microsoft. Our institute has conducted a survey, the TTS Survey Report [9], to determine the scope of this problem. Surprisingly, we discovered that it is not unusual for over 50%, and even 80%, of the words in a sentence to be mispronounced, and that there is a trend for cliticized words to be incorrectly pronounced. For example, the cliticized word يُوَلِلْكَاتِبِينَ , correctly pronounced *walilkatibîna*, is mispronounced as *walilkātibáyna*.

The tables in that report show that he error rate of Arabic TTS is unacceptably high. Such a high error rate would be unthinkable in the other major world languages. Another issue is **prosody** (stress and intonation) and **vowel neutralization** (e.g. *ing* is written as a long vowel in *i* but is shortened in actual pronunciation to *na*). This is a complex issue, described in detail in Halpern's paper on Arabic stress [4].

Speech synthesis, speech recognition and prosody in current Arabic speech technology are, on the whole, inaccurate, unnatural and often unpleasant to the ear. The time has come for developers to make serious efforts to make dramatic improvements.

#### 3.3 Improving ASR accuracy

*ArabLEX* includes features specifically designed to support automatic speech recognition (ASR). For speech synthesis (TTS), it is only necessary to *generate* one accurate pronunciation. For example, کاتبون 'writers' in standard Arabic is pronounced *katibūna*, but for ASR it is also necessary to *recognize* the less formal variant pronunciation *katibūn* Similarly, the standard pronunciation of 'I write' is '*áktubu*, but the final vowel is often omitted and it is pronounced '*áktub*.

The above alternatives are on a *phonemic* level. That is, the phoneme /na/ is being replaced by the phoneme /n/ as a result of vowel omission. There are also variations on the *phonetic* level; that is, certain phonemes have regional allophones. For example, ج in such words as جمل jamal is pronounced [gɛ̈mɛl] in Egypt, [dʒɛ̈mɛ̃l] in the Gulf region, and [ʒɛ̈mɛ̃l] in the Levant. It is important top note that this does not refer to the local dialects in those regions, but to a *regional varieties* of Modern Standard Arabic (MSA).

Thus *ArabLEX* not only represents  $\gtrsim$  in the standard IPA [dʒ] for TTS, it also lists the regional [ʒ] and [g] for ASR training. The goal is to enable the recognition of these allophones, but not to generate them.

### 3.4 Benefits to speech technology

One of the key components for training speech technology systems is the *pronunciation dictionary*. A major feature of *ArabLEX* is that it can serve as an extremely comprehensive pronunciation dictionary.

*ArabLEX* not only maps all unvocalized forms (including all cliticized forms) to their vocalized counterparts and to their lemmas, but also provides precise **phonemic transcriptions** (CARS system) [5] and **phonetic transcriptions** (IPA) that includes precise word stress and vowel neutralization for each entry. For example, in the IPA *wɛ̃likɛ̃:'tibikumɛ('*), the stressed syllable is indicated by (') (U+0C28), while (') (U+02D1) indicates that the final  $\varepsilon$  is neutralized vowel of optional half length. These features can help developers significantly enhance the quality of Arabic TTS, and can be used in training ASR systems to achieve higher recognition rates.

To summarize, *ArabLEX* can bring the following benefits to speech technology:

- Covers approximately 530 million entries, including millions of proper nouns.
- Covers *all combinations* of proclitics and enclitics for inflected wordforms (mostly verbs, nouns, adjectives and proper nouns).
- Tens of millions of orthographic variants for all wordforms.
- Provides an exhaustive list of alternative pronunciations of identical unvocalized strings to enable orthographical disambiguation (e.g. six alternatives for كاتباتك).
- Future versions will provide 'importance flags' to help determine the most likely alternative.
- Highly accurate phonemic transcriptions for all wordforms, including precise stress and vowel neutralization.

 Phonetic transcriptions (IPA) indicate the correct allophonic variants in context as well as regional variants for ASR.

### 4. Enhancing Machine Translation

Although neural machine translation (NMT) has achieved dramatic improvements in translation quality, it does have some shortcomings, as pointed out by Philipp Koehn [3] and in Halpern's paper [1].

Some issues in Arabic MT are (1) the high orthographic ambiguity, (2) the morphological complexity (forms like ولكاتباتهما are difficult to analyze), (3) the recognition of named entities (which are often cliticized), and (4) the large number of wordforms for Arabic nouns and verbs.

ArabLEX can significantly enhance the translation accuracy of Arabic MT. Not only can it be integrated into NMT systems to provide comprehensive coverage of cliticized forms, but it can also be be used as a special kind of corpus to train the language model and enable more accurate morphological, syntactic, and semantic analysis.

When NMT first appeared, it was believed that lexicons could not be integrated into NMT systems. Later it was shown that it is technically possible to do so by regarding a lexicon as a kind of sentence-aligned, parallel corpus and assigning a higher probability to lexicon lookup results so as to override the results of the normal NMT algorithms. By using such techniques, it should therefore be possible to integrate *ArabLEX* into Arabic NMT systems [1].

### 5. How ArabLEX Works

Let us demonstrate the broad scope of the information that *ArabLEX* provides on the morphology, phonology, grammar, and orthography of Arabic words and

their thousands of inflected and cliticized forms. The stem كَاتِب *kātib* 'writer', for example, combines with the proclitics وَ *wa* and *ل li* and the enclitic اتِهِمَا *tíhima* to yield وَلِكَاتِبَاتِهِمَا *walikaِtibātíhima*ِ.

Grammatical information		
Data field	Value	
Full form	ۅٙڸػؘٳؾؚڹؚػؙؗڡؘٳ	
Lemma	ػؘٵؾؚڹۨ	
Stem	كَاتِب	
Gender	С	
Case	GEN	
Number	D	
Person	2	
Definiteness	D	
Root	ك-ت-ب	

### Grammatical information

#### Phonological information

Data Field	Value	
Unvocalized	محمد	
Vocalized	مُحَمَّدُ	
Phonemic	muhammadun	
Phonetic	mu'ħɛ̈mmɛ̈dun	
Transliterated	muham~dN	

#### Morphological information

Data Field	Value	Transcription
Full form	ۅؘڸؚػؘٳؾؚؚؠػؙڡؘٳ	walikātibíkum <u>a</u>
Lemma	ػٙٳؾؚڹ۠	kātibun
Stem	كَاتِب	kātib
Proclitic	وَلِ	wali
Enclitic	كُمَاِ	(i)kúm <u>a</u>
Root	ك-ت-ب	k-t-b

Data Field	Value	Transcription
Variant 1	ألكسندرة	²aliksándara
Variant 2	الكسندرة	²aliksándara
Variant 3	ألكسندره	'aliksándara
Variant 4	الكسندره	'aliksándara

#### Orthographical information

#### 5.1 Grammatical information

This includes gender codes, number codes, case ending codes, persons codes, the stem, the state (definiteness), and the lemma. For وَلِكَاتِبَاتِهِمَا, *ArabLEX* provides as the following grammatical information.

Data field	Value	Description
Full form	ۅؘڸؚػؘٳؾؚؚؚڮؘؙؗؗؗڡؘٳ	walikātibíkum <u>a</u>
Lemma	ػؘٵؾؚڹۨ	kātibun
Stem	كَاتِب	kātib
Gender	С	common gender (masculine & feminine)
Case	GEN	genitive case
Number	D	dual
Person	2	second person
State	D	definite, indefinite or construct state
Root	ك-ت-ب	the triliteral root

Table 1: Grammatical information

Abundant grammatical information is useful for morphological analysis, orthographic disambiguation, semantic analysis, and pedagogical applications.

### 5.2 Phonological information

This includes full vocalization and phonemic transcription for all unvocalized Arabic, including the full form headword, the lemma, the stem, proclitics, and enclitics. The headword is also given in IPA, a precise phonetic transcription that includes word stress. For example, for each full form Arabic headword the following fields are given:

Data field	Value	Description	
Unvocalized	محمل	Unvocalized Arabic as it actually occurs	
Vocalized	مُحَمَّد	precise and full vocalization	
Phonemic	muhammadun	accurate phonemic transcription in CARS	
Phonetic	mu'ħëmmëdun	phonetic transcription in IPA (or optionally SAMPA ) with word stress	
Transliterated	muham~dN	orthographic transliteration in the Buckwalter system	

#### Table 2: Phonological information

The phonemic and phonetic transcriptions provide precise information that is useful for training speech technology systems, both TTS and ASR.

#### 5.3 Morphological information

This includes the lemma, stem, proclitics (prefixes) and enclitics (suffixes) for each full form headword, as shown below.

	_	
Data field	Value	Transcription
Full form	ۅؘڸؚػؘٳؾؚڹػؙڡؘٳ	walikātibíkuma

#### Table 3: Morphological information

Lemma	ػؘٳؾؚڹۨ	kātibun
Stem	كَاتِب	kātib
Proclitic	وَلِ	wali
Enclitic	كُمَاِ	(i)kúm <u>a</u>
Root	ك-ت-ب	k-t-b

The morphological information is useful for morphological analysis, semantic analysis, lemmatization, decliticization, verb conjugation, and dictionary lookup.

#### 5.4 Orthographical information

This includes the vocalized and unvocalized headwords as well as their orthographic variants, as shown below.

Data field	Value	Transcription
Vocalized headword	ٲؘٵڵػؘٳؾؚڹۨ	² <u>a</u> lkātibun
Unvocalized headword	آلكاتب	² <u>a</u> lkātibun
Unvocalized variant	أالكاتب	²alkātibun

#### Table 4: Orthographical information

The orthographical information is useful for word/entity recognition, word/entity extraction, normalization, and dictionary lookup.

#### 5.5 Orthographical disambiguation

A central issue in Arabic NLP applications, especially in speech technology, is identifying which of the possible wordforms an Arabic string like كاتباتـك represents. This can represent six wordforms shown in the table below, each

with a different meaning and different morphological or syntactic function. We will refer to this process as *orthographical disambiguation*.

POS	ARAB_V	ARAB_U	ARAB_T	GEN	NUM	CASE	PER	DEF
N	كَاتِبَاتُكَ	كاتباتك	katibātuka	F	Р	NOM	2SM	d
N	كَاتِبَاتُكِ	كاتباتك	katibatuki	F	Р	NOM	2SF	d
N	كَاتِبَاتِكَ	كاتباتك	kaౖtibâtika	F	Р	GEN	2SM	d
N	كَاتِبَاتِكِ	كاتباتك	katibatiki	F	Р	GEN	2SF	d
N	كَاتِبَاتِكَ	كاتباتك	kaౖtibātika	F	Р	ACU	2SM	d
N	كَاتِبَاتِكِ	كاتباتك	katibatiki	F	Р	ACU	2SF	d

Table 5: Orthographical disambiguation

The rich set of grammatical attributes shown above (and the morphological attributes not shown) can help train the language model to correctly disambiguate such ambiguous forms. That is, they provide the grammatical and morphological context in which کاتباتک can occur, helping determine the specific correct wordform for that context, and thus the correct pronunciation. For example, the attributes for کَاتِبَاتِک show that it refers to plural female writers in the definite state, who belong to second person singular feminine in the genitive case, which helps to determine the correct pronunciation of *katibatiki* in training speech technology models.

### 6. Conclusions

ArabLEX provides a rich set of morphological and phonological features. It brings the following benefits to speech technology and MT:

- Enhance the quality of MT, NLP and AI applications.
- Full support for accurate morphological analysis, including stemming, lemmatization, segmentation and tokenization.

- Supplements corpora for training speech technology models.
- Improved accuracy of word and entity recognition and extraction.
- Support for query processing in information retrieval applications.
- Support for automatic conjugators for pedagogical and NLP applications.
- Part-of-speech analysis and POS tagging.
- Accurate determination of the root for each wordform.

In summary, *ArabLEX* aims to serve as the ultimate resource for Arabic natural language processing.

### References

[1] Halpern, Jack. <u>Very Large-Scale Lexical Resources to Enhance Chinese and</u> <u>Japanese Machine Translation</u>

[2] Habash, Nizar Y. Introduction to Arabic Natural Language Processing

[3] Koehn, Philipp. The State of Neural Machine Translation (NMT)

[4] Halpern, Jack. Word Stress and Vowel Neutralization in Modern Standard Arabic

[5] Halpern, Jack. CJKI Arabic Romanization System

[6] Halpern, Jack. <u>Lexicon-Driven Approach to the Recognition of Arabic Named</u> <u>Entities</u>

[7] Carbonell, Jaime; Klein, Steve; Miller, David; Steinbaum, Michael. <u>Context-Based</u> <u>Machine Translation</u>

[8] Wajdan Algihab, Imam Muhammad bin Saud Islamic University. <u>Arabic Speech</u> <u>Recognition with Deep Learning: A Review</u>

[9] Halpern, Jack. <u>TTS Survey Report</u>

## **APPENDIX 1: DATA FIELDS**

### Table 6: Basic data fields

No	Field	Value	Field Description
1	IDENTIFIER	05373	unique ID for Arabic headword (full form) or variant thereof
2	SUBID	00	identifies variants of headword
3	POS	N	part of speech code
4	ARAB_V	ۅؘڸؚػؘٳؾؚؚؚػؙڡؘٳ	vocalized Arabic headword (full form)
5	ARAB_U	ولكاتبكما	unvocalized Arabic headword (full form)
6	ARAB_T	walikātibíkuma	Arabic headword in phonemic transcription (full form) in CARS system (including vowel neutralization)
7	LEMMA_V	ػؘٳؾؚڹۨ	lemma in vocalized Arabic
8	LEMMA_U	کاتب	lemma in unvocalized Arabic
9	LEMMA_T	kấtibun	lemma in phonemic transcription in CARS
10	GEN	М	gender code for stem:masculine

11	NUM	S	number code for stem: singular
12	CASE	GEN	case ending code for nouns and adjectives: genitive
13	PER	2DC	person code for full form (ARABIC_V): second person dual common gender
14	DEF	D	Definite, indefinite or construct state
15	TYPE	N/A	code for verb conjugation (for verbs only)
16	TENSE	N/A	code for verb tense (for verbs only)

### Table 7: Advanced data fields

No	Field	Value	Field Description
17	PROC_V	وَلِ	prefix or proclitic in vocalized Arabic
18	PROC_U	ول	prefix or proclitic in unvocalized Arabic
19	PROC_T	wáli	prefix or proclitic in phonemic transcription
20	STEM_V	كَاتِب	stem in vocalized Arabic
21	STEM_U	كاتب	stem in unvocalized Arabic
22	STEM_T	kấtib	stem in phonemic transcription in CARS

23	ENC_V	كُمَاِ	suffix or enclitic in vocalized Arabic
24	ENC_U	کما	suffix or enclitic in unvocalized Arabic
25	ENC_T	ikúm <u>a</u>	suffix or enclitic in phonemic transcription in CARS
26	ROOT	ك-ت-ب	the root of each headword
27	TRANSLIT	walikaAtibikumaA	graphemic transliteration in the Buckwalter system
28	IPA	wëlikë:'tibikumɛ(')	phonetic transcription of full-form headword in IPA, including word stress and vowel neutralization

## **APPENDIX 2: DATA SAMPLE**

Below is a small subset of *ArabLEX* for the noun كَاتِب*ُ kātibun,* which contains a total of about 5000 wordforms. The full sample can be found at *ArabLEX\_*sample.xls.

POS	ARAB_V	ARAB_U	ARAB_T	LEMMA_V	LEMMA_U	LEMMA_T
N	ڣؘڸػؘٳؾؚؚؚػؙڡؘٳ	فلكاتبكما	falikātibíkuma	ػٙٳؾؚڹۨ	کاتب	kấtibun
N	ڣؘڸػٙٳؾؚؠؚڡؚۣؗؗۿ	فلكاتبهم	falikātíbihim	ػٙٳؾؚؚؚٮ	کاتب	kấtibun
N	ڣؘڸػٙٳؾؚؚڣۣڹٞ	فلكاتبهن	falikātibihínna	ػٙٳؾؚؚؚٮ	كاتب	kấtibun
N	فَلِكَاتِبِهِمَا	فلكاتبهما	falikātibíhim <u>a</u>	ػؘٳؾؚؚڹۨ	کاتب	kấtibun
N	فَلِكَاتِبِهِمَا	فلكاتبهما	falikātibíhim <u>a</u>	ػٙٳؾؚڹۨ	کاتب	kấtibun
N	فَلِكَاتِبِي	فلكاتبي	falikātib <u>i</u>	ػٙٳؾؚؚڹۨ	کاتب	kấtibun
N	فَلِكَاتِبَيْنِ	فلكاتبين	falikātibáyni	ػٙٳؾؚؚؚؚ	کاتب	kấtibun
N	ڣؘڸؚػؘٳؾؚڹؾٞ	فلكاتبي	falikātibáyya	ػؘٳؾؚؚؚ	کاتب	kấtibun
N	فَلِكَاتِبَيْكَ	فلكاتبيك	falikātibáyka	ػؘٳؾؚڹۨ	کاتب	kấtibun
N	فَلِكَاتِبَيْكِ	فلكاتبيك	falikātibáyki	ػؘٳؾؚڹۨ	کاتب	kấtibun

Table 8: San	nple for nouns	(fields 3 to 9)
--------------	----------------	-----------------

ARAB_V	GEN	NUM	CASE	PER	DEF	TYPE	TENSE
فَلِكَاتِبِكُمَا	М	S	GEN	2DC	d	-	-
ڣؘڸػؘٳؾؚؚڣۣؗؗؗؗؗ	М	S	GEN	3PM	d	-	-
ڣؘڸػؘٳؾؚؚڣۣڹٞ	М	S	GEN	3PF	d	-	-
فَلِكَاتِبِهِمَا	М	S	GEN	3DM	d	-	-
فَلِكَاتِبِهِمَا	М	S	GEN	3DF	d	-	-
فَلِكَاتِبِي	М	S	ACU	1SC	d	-	_
ڣؘڸؚػٙٳؾؚؠٙؽ۠ڹۣ	М	D	GEN	0	D	-	-
ڣؘڸػٳؾؚڹؾۜ	М	D	GEN	1SC	d	-	_
فَلِكَاتِبَيْكَ	М	D	GEN	2SM	d	-	-
ڣؘڸػؘٳؾؚؠؘؽڮ	М	D	GEN	2SF	d	-	-

Table 9: Sample for nouns (fields 10-16)

Table 10: Sample for nouns (fields 17-22)

ARAB_V	PROC_V	PROC_U	PROC_T	STEM_V	STEM_U	STEM_T
فَلِكَاتِبِكُمَا	فَلِ	فل	fáli	کَاتِب	کاتب	kấtib
فَلِكَاتِبِهِمْ	فَلِ	فل	fáli	کَاتِب	کاتب	kấtib
ڣؘڸؚػؘٳؾؚؚڽؚڡۣڹۘٞ	فَلِ	فل	fáli	کَاتِب	کاتب	kấtib

فَلِكَاتِبِهِمَا	فَلِ	فل	fáli	کَاتِب	کاتب	kấtib
فَلِكَاتِبِهِمَا	فَلِ	فل	fáli	کَاتِب	کاتب	kấtib
فَلِكَاتِبِي	فَلِ	فل	fáli	کَاتِب	کاتب	kấtib
فَلِكَاتِبَيْنِ	فَلِ	فل	fáli	کَاتِب	کاتب	kấtib
ڣؘڸػؘٳؾؚڹؘۑٞۜ	فَلِ	فل	fáli	كَاتِب	کاتب	kấtib
فَلِكَاتِبَيْكَ	فَلِ	فل	fáli	كَاتِب	کاتب	kātib
فَلِكَاتِبَيْكِ	فَلِ	فل	fáli	كَاتِب	کاتب	kātib

Table 11: Sample for nouns (fields 23-28)

ARAB_V	ENC_V	ENC_U	ENC_T	ROOT	BW	IPA
فَلِكَاتِبِكُمَا	کُمَاِ	کما	ikúm <u>a</u>	ك-ت-ب	falikaAtibikumaA	
فَلِكَاتِبِهِمْ	هِمْ	هم	ihim	ك-ت-ب	falikaAtibihimo	
ڣؘڸؚػؘٳؾؚؚڽؚڡۣڹۜٞ	ۿؚڹۘٞ	هن	íhínna	ك-ت-ب	falikaAtibihin~a	
فَلِكَاتِبِهِمَا	هِمَا	هما	íhím <u>a</u>	ك-ت-ب	falikaAtibihimaA	
فَلِكَاتِبِهِمَا	هِمَا	هما	íhím <u>a</u>	ك-ت-ب	falikaAtibihimaA	
فَلِكَاتِبِي	ي	ي	<u>i</u>	ك-ت-ب	falikaAtibiy	
فَلِكَاتِبَيْنِ	ؽڹؚٙ	ين	áyni	ك-ت-ب	falikaAtibayoni	

ڣؘڸػٵؾؚڹؘۑٞٙ	ۑؖ	ي	áyya	ك-ت-ب	falikaAtibay~a	
فَلِكَاتِبَيْكَ	يْكَ	يك	áyka	ك-ت-ب	falikaAtibayoka	
فَلِكَاتِبَيْكِ	ؽڬؚ	يك	áyki	ك-ت-ب	falikaAtibayoki	