

ArabLEX

برج المعاجم
Tower of Lexicons



ArabLEX

Arabic Full Form Lexicon

معجم اللغة العربية الكامل

1. Introduction

1.1 Synopsis

The CJK Dictionary Institute (CJKI) is pleased to announce the release of the **Comprehensive Arabic Full Form Lexicon** (*ArabLEX*), covering approximately **538 million entries** inflected, conjugated and cliticized wordforms. This includes both general vocabulary as well as fully inflected Arab and non-Arab proper nouns (personal and place names).

ArabLEX is not only comprehensive in coverage, but is also rich in morphological, grammatical, phonological, and orthographical attributes. In addition, it maps all unvocalized forms to their vocalized counterparts and to the lemma, and provides precise phonemic (and eventually phonetic) transcriptions. These features can significantly contribute to the training of language models for NLP, deep learning and AI for such applications as speech technology (TTS and ASR) and machine translation. It aims to serve as the ultimate resource for Arabic natural language processing.

This document includes specification and data samples for all four *ArabLEX*:

DAG	Database of Arabic General Vocabulary	88 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

1.2 Approximate Quantities

Currently the DAG module (general vocabulary) of *ArabLEX* is estimated at 88 million entries. If we include proper nouns (DAF, DAP and DAN), the total is approx 538 million. In future expansion, we plan to include negated forms, periphrastic elatives, compound verb tenses and other categories, and thus expect the total entries to exceed one billion.

1.3 Database Structure and Schema

Each of the four modules (DAG, DAF, DAP and DAN) are stored in four data tables:

(1) **Canonical**, (2) **Morphological**, (3) **Orthographical**, and (4) **Phonological**. The database is stored in 15 data tables interlinked by IDs and SubIDs within each module. Each module consists of the following functionally distinct data tables:

Data Tables for Each Module

		DAG	DAP	DAF	DAN
1	Canonical		DAP-C	DAF-C	DAN-C
2	Morphological	DAG-M	DAP-M	DAF-M	DAN-M
3	Orthographical	DAG-O	DAP-O	DAF-O	DAN-O
4	Phonological	DAG-P	DAP-P	DAF-P	DAN-P

(1) **Canonical**: This is the primary table that provides basic data such as the lemmata in several formats, romanization, ID links to other modules, and other fields as necessary such as gender codes and frequency.

(2) **Morphological**: Comprehensive morphological attributes (max over 20) and wordforms (often thousands) for each lemma, including inflections (possessives, plurals), declensions (case endings), conjugations (tense and person), and cliticized forms. It is the core, often the largest, component of *ArabLEX*.

(3) **Orthographical**: Various attributes such as orthographical and spelling variants, including vocalized and unvocalized variants for each lemma and their wordforms.

(4) **Phonological**: Various phonological attributes, such as phonemic and phonetic transcriptions and syllabification for all lemmata and their wordforms (inflections, conjugations, etc.) ideal for speech technology.

NOTE: *ArabLEX* is constantly evolving and expanding. Some modules don't necessarily have all the fields shown in the field descriptions.

1.4 More Information

Detailed whitepaper

<https://www.cjk.org/wp-content/uploads/ArabLEX.pdf>

White paper summary

https://www.cjk.org/wp-content/uploads/ArabLEX_summary_.pdf

Academic paper

<https://www.cjk.org/wp-content/uploads/Halpern-LREC2022Paper.pdf>

Detailed sample

<https://www.cjk.org/wp-content/uploads/2020/11/ArabLEX.xls>

The cover: *The Burj Khalifa*, the tallest building in the world, soars 832 meters above the skies of Dubai. This is symbolic of *ArabLEX*, which is the largest Arabic lexicon ever created.

2. DATABASE OF ARABIC GENERAL VOCABULARY (DAG)

2.1 Synopsis

DAG is a **comprehensive full form lexicon of Arabic general vocabulary** that contains a rich set of morphological, grammatical, and phonological attributes. Proper nouns are in principle excluded as they are given in the modules DAN, DAF and DAP (described below).

2.2 More Information

Website: <http://www.cjk.org/data/arabic/nlp/arabic-full-form-lexicon/>

White paper: ArabLEX.pdf

Data sample: ArabLEX.xls

2.3 Approximate Quantities

Type	Lemmas	Inflected forms
Nouns	14,293	30,679,202
Adjectives	6,228	23,874,505
Verbs	9,509	33,377,031
TOTAL	30,030	87,930,738

2.4 The Morphological Component (DAG-M)

2.4.1 Field Description

No.	Field	Description
1	ID	uniquely identifies the lemma
2	SUBID	uniquely identifies inflected forms
3	POS	part of speech code [V] verb [N] noun [A] adjective
4	SP	POS subclassification verbs [T] transitive [I] intransitive [D] ditransitive [B] transitive & intransitive adjectives [P] positive [E] elative

		[e] elative with non-elative meaning
5	ARAB_V	vocalized cliticized Arabic entry
6	ARAB_BW	vocalized cliticized Arabic entry in Buckwalter transliteration
7	ARAB_U	unvocalized cliticized Arabic entry
8	LEMMA_V	lemma in vocalized Arabic
9	LEMMA_U	lemma in unvocalized Arabic
10	GEN	gender code for stem [M] masculine [F] feminine
11	NUM	number code for stem [S] singular [D] dual [P] plural [B] broken plural
12	CASE	case ending code for nouns and adjectives [ACU] accusative [GEN] genitive [NOM] nominative
13	PER	(person code) person, number, gender [000] zero person [1SC] first person singular [2SM] second person singular masculine [2SF] second person singular feminine [3SM] third person singular masculine [3SF] third person singular feminine [1PC] first person plural [2PM] second person plural masculine [2PF] second person plural feminine [2DC] second person dual [3PM] third person plural masculine [3PF] third person plural feminine [3DM] third person dual masculine [3DF] third person dual feminine [XSM] singular masculine [XDM] dual masculine [XPM] plural masculine [XSF] singular feminine [XDF] dual feminine [XPF] plural feminine
14	DEF	definiteness [D] explicitly definite by cliticizing of article Al [d] definite by adding pronomial enclitics [I] indefinite [G] first term of genitive construct (Idhafa, e.g. baytu in baytu-l-kaatibi)
15	RAT <small>(under construction)</small>	rationality (animate or not) [R] rational [I] irrational

16	FORM	verb form (see cavehelp.htm §12) [01] Form I [02] Form II [03] Form III [04] Form IV [05] Form V [06] Form VI [07] Form VII [08] Form VIII [09] Form IX [10] Form X [Q1] Form QI [Q2] Form QII [Q3] Form QIII [Q4] Form QIV
17	TYPE	code for verb conjugation type (see cavehelp.htm §13) [S] sound [Z] hamzated [G] geminate [A] assimilated [H] hollow [D] defective [W] doubly week
18	TENSE	code for verb tense [01] active perfect indicative [02] active imperfect indicative [03] active imperfect subjunctive [04] active imperfect jussive [05] active imperfect imperative [06] active imperfect energetic I OMITTED because obsolete [07] active imperfect energetic II OMITTED because obsolete [08] passive perfect indicative [09] passive imperfect indicative [10] passive imperfect subjunctive [11] passive imperfect jussive [12] active participle [13] passive participle [14] verbal noun [52] future indicative [53] future passive
19	PROC_V	prefix or proclitic in vocalized Arabic
20	PROC_U	prefix or proclitic in unvocalized Arabic
21	STEM_V	stem in vocalized Arabic
22	STEM_U	stem in unvocalized Arabic
23	ENC_V	suffix or enclitic in vocalized Arabic
24	ENC_U	suffix or enclitic in unvocalized Arabic

2.4.2 Data Sample

The fields in the sample below are a subset of the fields in the morphological component. There are many other fields useful for natural language processing. IMPORTANT: Only a small subset of available fields is shown here. Be sure to check out the full sample at ArabLEX.XLS for all fields.

Morphological Component (nouns)

SUBID	POS	ARAB_V	ARAB_BW	LEMMA_V	GEN	NUM	CASE	PER	DEF	ROOT
0001	N	كَاتِبٌ	kaAtibN	كَاتِبٌ	M	S	NOM	000	I	كتب
0002	N	كَاتِبٌ	kaAtibu	كَاتِبٌ	M	S	NOM	000	G	كتب
0003	N	كَاتِبِي	kaAtibiy	كَاتِبٌ	M	S	NOM	1SC	d	كتب
0004	N	كَاتِبُكَ	kaAtibuka	كَاتِبٌ	M	S	NOM	2SM	d	كتب
0005	N	كَاتِبُكِ	kaAtibuki	كَاتِبٌ	M	S	NOM	2SF	d	كتب
0006	N	كَاتِبُهُ	kaAtibuhu	كَاتِبٌ	M	S	NOM	3SM	d	كتب
0007	N	كَاتِبُهَا	kaAtibuhaA	كَاتِبٌ	M	S	NOM	3SF	d	كتب
0008	N	كَاتِبَنَا	kaAtibunaA	كَاتِبٌ	M	S	NOM	1PC	d	كتب
0009	N	كَاتِبُكُمْ	kaAtibukumo	كَاتِبٌ	M	S	NOM	2PM	d	كتب
0010	N	كَاتِبُكُنَّ	kaAtibukun~a	كَاتِبٌ	M	S	NOM	2PF	d	كتب
0011	N	كَاتِبُكُمَا	kaAtibukumaA	كَاتِبٌ	M	S	NOM	2DC	d	كتب
0012	N	كَاتِبُهُمْ	kaAtibuhumo	كَاتِبٌ	M	S	NOM	3PM	d	كتب
0013	N	كَاتِبُهُنَّ	kaAtibuhun~a	كَاتِبٌ	M	S	NOM	3PF	d	كتب
0014	N	كَاتِبُهُمَا	kaAtibuhumaA	كَاتِبٌ	M	S	NOM	3DM	d	كتب
0015	N	كَاتِبُهُمَا	kaAtibuhumaA	كَاتِبٌ	M	S	NOM	3DF	d	كتب
0016	N	كَاتِبٍ	kaAtibK	كَاتِبٌ	M	S	GEN	000	I	كتب
0017	N	كَاتِبٍ	kaAtibi	كَاتِبٌ	M	S	GEN	000	G	كتب
0018	N	كَاتِبِي	kaAtibiy	كَاتِبٌ	M	S	GEN	1SC	d	كتب

0019	N	كَاتِبٌ	kaAtibika	كَاتِبٌ	M	S	GEN	2SM	d	كتب
0020	N	كَاتِبٌ	kaAtibiki	كَاتِبٌ	M	S	GEN	2SF	d	كتب
0021	N	كَاتِبٍهُ	kaAtibihi	كَاتِبٌ	M	S	GEN	3SM	d	كتب
0022	N	كَاتِبِهَا	kaAtibihaA	كَاتِبٌ	M	S	GEN	3SF	d	كتب
0023	N	كَاتِبِنَا	kaAtibinaA	كَاتِبٌ	M	S	GEN	1PC	d	كتب
0024	N	كَاتِبِكُمْ	kaAtibikumo	كَاتِبٌ	M	S	GEN	2PM	d	كتب
0025	N	كَاتِبِكُنَّ	kaAtibikun~a	كَاتِبٌ	M	S	GEN	2PF	d	كتب
0026	N	كَاتِبِكُمَا	kaAtibikumaA	كَاتِبٌ	M	S	GEN	2DC	d	كتب
0027	N	كَاتِبِهِمْ	kaAtibihimo	كَاتِبٌ	M	S	GEN	3PM	d	كتب
0028	N	كَاتِبِهِنَّ	kaAtibihin~a	كَاتِبٌ	M	S	GEN	3PF	d	كتب
0029	N	كَاتِبِهِمَا	kaAtibihimaA	كَاتِبٌ	M	S	GEN	3DM	d	كتب
0030	N	كَاتِبِهِمَّا	kaAtibihimaA	كَاتِبٌ	M	S	GEN	3DF	d	كتب

Morphological Component (verbs)

SUBID	POS	SP	ARAB_V	ARAB_BW	LEMMA_V	PER	FORM	TYPE	TENSE	ROOT
0229	V	T	كَتَبَ	kataba	كَتَبَ	3SM	01	S1	01	كتب
0230	V	T	كَتَبْنِي	katabaniy	كَتَبَ	3SM	01	S1	01	كتب
0231	V	T	كَتَبَكَ	katabaka	كَتَبَ	3SM	01	S1	01	كتب
0232	V	T	كَتَبَكِ	katabaki	كَتَبَ	3SM	01	S1	01	كتب
0233	V	T	كَتَبَهُ	katabahu	كَتَبَ	3SM	01	S1	01	كتب
0234	V	T	كَتَبَهَا	katabahaA	كَتَبَ	3SM	01	S1	01	كتب
0235	V	T	كَتَبَنَا	katabanaA	كَتَبَ	3SM	01	S1	01	كتب
0236	V	T	كَتَبْكُمْ	katabakumo	كَتَبَ	3SM	01	S1	01	كتب
0237	V	T	كَتَبْكُنَّ	katabakun~a	كَتَبَ	3SM	01	S1	01	كتب
0238	V	T	كَتَبْكُمَا	katabakumaA	كَتَبَ	3SM	01	S1	01	كتب
0239	V	T	كَتَبَهُمْ	katabahumo	كَتَبَ	3SM	01	S1	01	كتب
0240	V	T	كَتَبَهُنَّ	katabahun~a	كَتَبَ	3SM	01	S1	01	كتب

0241	V	T	كَتَبُهُمَا	katabahumaA	كَتَبْ	3SM	01	S1	01	كتب
0255	V	T	وَكَتَبَ	wakataba	كَتَبْ	3SM	01	S1	01	كتب
0256	V	T	وَكَتَبَنِي	wakatabaniy	كَتَبْ	3SM	01	S1	01	كتب
0257	V	T	وَكَتَبَكَ	wakatabaka	كَتَبْ	3SM	01	S1	01	كتب
0258	V	T	وَكَتَبَكِ	wakatabaki	كَتَبْ	3SM	01	S1	01	كتب
0259	V	T	وَكَتَبَهُ	wakatabahu	كَتَبْ	3SM	01	S1	01	كتب
0260	V	T	وَكَتَبَهَا	wakatabahaA	كَتَبْ	3SM	01	S1	01	كتب
0261	V	T	وَكَتَبَنَا	wakatabanaA	كَتَبْ	3SM	01	S1	01	كتب
0262	V	T	وَكَتَبَكُمْ	wakatabakumo	كَتَبْ	3SM	01	S1	01	كتب
0263	V	T	وَكَتَبَكُنْ	wakatabakun~a	كَتَبْ	3SM	01	S1	01	كتب
0264	V	T	وَكَتَبَكُمَا	wakatabakumaA	كَتَبْ	3SM	01	S1	01	كتب
0265	V	T	وَكَتَبَهُمْ	wakatabahumo	كَتَبْ	3SM	01	S1	01	كتب
0266	V	T	وَكَتَبَهُنْ	wakatabahun~a	كَتَبْ	3SM	01	S1	01	كتب
0267	V	T	وَكَتَبَهُمَا	wakatabahumaA	كَتَبْ	3SM	01	S1	01	كتب
0268	V	T	فَكَتَبَ	fakataba	كَتَبْ	3SM	01	S1	01	كتب
0269	V	T	فَكَتَبَنِي	fakatabaniy	كَتَبْ	3SM	01	S1	01	كتب
0270	V	T	فَكَتَبَكَ	fakatabaka	كَتَبْ	3SM	01	S1	01	كتب
0271	V	T	فَكَتَبَكِ	fakatabaki	كَتَبْ	3SM	01	S1	01	كتب

2.5 The Orthographical Component (DAG-O)

2.5.1 Field Description

No.	Field	Description
1	ID	uniquely identifies the lemma
2	SUBID	uniquely identifies inflected forms
3	VARID	identifies variants of headword [00] no variants [01] main form [02] variant form
4	ARAB_V	vocalized cliticized Arabic entry

5	VAR_V	variant in vocalized Arabic
6	VAR_U	variant in unvocalized Arabic

2.5.2 Data Sample

Under construction.

2.6 The Phonological Component (DAG-P)

2.6.1 Field Description

No.	Field	Description
1	ID	uniquely identifies the lemma
2	SUBID	uniquely identifies inflected forms
3	VARID	identifies variants of headword [00] no variants [01] main form [02] variant form
4	ARAB_V	vocalized cliticized Arabic entry
5	CARS	cliticized Arabic entry in phonemic transcription (CARS)
6	IPA	phonetic transcription in the International Phonetic Alphabet (IPA)
7	XSAMPA	phonetic transcription in X-SAMPA

2.6.2 Data Sample

Phonological Component (nouns)

SUBID	ARAB_V	CARS	IPA	XSAMPA
0001	كَاتِبٌ	kátibun	'ka:.ti.bun	"ka:.ti.bun
0002	كَاتِبٌ	kátibu	'ka:.ti.bu	"ka:.ti.bu
0003	كَاتِبِي	kátibí	'ka:.ti.bi	"ka:.ti.bi
0004	كَاتِبَكَ	kátíbuka	ka:. 'ti.bu.ka	ka:. "ti.bu.ka
0005	كَاتِبَكَ	kátíbuki	ka:. 'ti.bu.ki	ka:. "ti.bu.ki
0006	كَاتِبَهُ	kátíbuhu	ka:. 'ti.bu.hu	ka:. "ti.bu.hu
0007	كَاتِبَهَا	kátíbuha	ka:. 'ti.bu.ha	ka:. "ti.bu.ha
0008	كَاتِبَنَا	kátíbuna	ka:. 'ti.bu.na	ka:. "ti.bu.na

0009	كَاتِبُكُمْ	kātibukum	ka:. 'ti.bu.kum	ka:. "ti.bu.kum
0010	كَاتِبُكُنَّ	kātibukúnna	ka:.ti.bu.'ku.n:a	ka:.ti.bu."ku.n:a
0011	كَاتِبُكُمَا	kātibúkumá	ka:.ti. 'bu.ku.ma	ka:.ti."bu.ku.ma
0012	كَاتِبُهُمْ	kātíbuhum	ka:. 'ti.bu.hum	ka:. "ti.bu.hum
0013	كَاتِبُهُنَّ	kātibuhúnna	ka:.ti.bu.'hu.n:a	ka:.ti.bu."hu.n:a
0014	كَاتِبُهُمَا	kātibúhumá	ka:.ti. 'bu.hu.ma	ka:.ti."bu.hu.ma
0015	كَاتِبُهُمَا	kātibúhumá	ka:.ti. 'bu.hu.ma	ka:.ti."bu.hu.ma
0016	كَاتِبٍ	kátibin	'ka:.ti.bin	"ka:.ti.bin
0017	كَاتِبٍ	kátibi	'ka:.ti.bi	"ka:.ti.bi
0018	كَاتِبِي	kátibj	'ka:.ti.bi	"ka:.ti.bi
0019	كَاتِبَكَ	kátibika	ka:. 'ti.bi.ka	ka:. "ti.bi.ka
0020	كَاتِبَكَ	kátibiki	ka:. 'ti.bi.ki	ka:. "ti.bi.ki
0021	كَاتِبِهِ	kátibih	ka:. 'ti.bi.hi	ka:. "ti.bi.hi
0022	كَاتِبِهَا	kátibihá	ka:. 'ti.bi.ha	ka:. "ti.bi.ha
0023	كَاتِبِنَا	kátibiná	ka:. 'ti.bi.na	ka:. "ti.bi.na
0024	كَاتِبُكُمْ	kātibikum	ka:. 'ti.bi.kum	ka:. "ti.bi.kum
0025	كَاتِبُكُنَّ	kātibikúnna	ka:.ti.bi.'ku.n:a	ka:.ti.bi."ku.n:a
0026	كَاتِبُكُمَا	kātibíkumá	ka:.ti. 'bi.ku.ma	ka:.ti."bi.ku.ma
0027	كَاتِبُهُمْ	kátibihim	ka:. 'ti.bi.him	ka:. "ti.bi.him
0028	كَاتِبُهُنَّ	kátibihínna	ka:.ti.bi.'hi.n:a	ka:.ti.bi."hi.n:a
0029	كَاتِبُهُمَا	kátibíhimá	ka:.ti. 'bi.hi.ma	ka:.ti."bi.hi.ma
0030	كَاتِبُهُمَا	kátibíhimá	ka:.ti. 'bi.hi.ma	ka:.ti."bi.hi.ma

Phonological Component (verbs)

SUBID	ARAB_V	CARS	IPA	XSAMPA
0229	كَتَبَ	kátaba	'ka.ta.ba	"ka.ta.ba
0230	كَتَبَنِي	katábaní	ka.'ta.ba.ni	ka."ta.ba.ni

0231	كَتَبَكَ	katábaka	ka.'ta.ba ka	ka."ta.ba ka
0232	كَتَبَكِي	katábaki	ka.'ta.ba ki	ka."ta.ba ki
0233	كَتَبَهُ	katábahu	ka.'ta.ba.hu	ka."ta.ba.hu
0234	كَتَبَهَا	katábahā	ka.'ta.ba.ha	ka."ta.ba.ha
0235	كَتَبَنَا	katábana	ka.'ta.ba.na	ka."ta.ba.na
0236	كَتَبْكُمْ	katábakum	ka.'ta.ba.kum	ka."ta.ba.kum
0237	كَتَبْكُنَّ	katabakúnna	ka.ta.ba.'ku.n:a	ka.ta.ba."ku.n:a
0238	كَتَبْكُمَا	katabákumā	ka.ta.'ba.ku.ma	ka.ta."ba.ku.ma
0239	كَتَبْهُمْ	katábahum	ka.'ta.ba.hum	ka."ta.ba.hum
0240	كَتَبْهُنَّ	katabahúnna	ka.ta.ba.'hu.n:a	ka.ta.ba."hu.n:a
0241	كَتَبْهُمَا	katabáhumā	ka.ta.'ba.hu.ma	ka.ta."ba.hu.ma
0255	وَكَتَبَ	wakátaba	wa.'ka.ta.ba	wa."ka.ta.ba
0256	وَكَتَبْنِي	wakatábaní	wa.ka.'ta.ba.ni	wa.ka."ta.ba.ni
0257	وَكَتَبَكَ	wakatábaka	wa.ka.'ta.ba ka	wa.ka."ta.ba ka
0258	وَكَتَبَكِ	wakatábaki	wa.ka.'ta.ba ki	wa.ka."ta.ba ki
0259	وَكَتَبَهُ	wakatábahu	wa.ka.'ta.ba.hu	wa.ka."ta.ba.hu
0260	وَكَتَبَهَا	wakatábahā	wa.ka.'ta.ba.ha	wa.ka."ta.ba.ha
0261	وَكَتَبَنَا	wakatábana	wa.ka.'ta.ba.na	wa.ka."ta.ba.na
0262	وَكَتَبْكُمْ	wakatábakum	wa.ka.'ta.ba.kum	wa.ka."ta.ba.kum
0263	وَكَتَبْكُنَّ	wakatabakúnna	wa.ka.ta.ba.'ku.n:a	wa.ka.ta.ba."ku.n:a
0264	وَكَتَبْكُمَا	wakatabákumā	wa.ka.ta.'ba.ku.ma	wa.ka.ta."ba.ku.ma
0265	وَكَتَبْهُمْ	wakatábahum	wa.ka.'ta.ba.hum	wa.ka."ta.ba.hum
0266	وَكَتَبْهُنَّ	wakatabahúnna	wa.ka.ta.ba.'hu.n:a	wa.ka.ta.ba."hu.n:a
0267	وَكَتَبْهُمَا	wakatabáhumā	wa.ka.ta.'ba.hu.ma	wa.ka.ta."ba.hu.ma
0268	فَكَتَبَ	fakátaba	fa.'ka.ta.ba	fa."ka.ta.ba
0269	فَكَتَبْنِي	fakatábaní	fa.ka.'ta.ba.ni	fa.ka."ta.ba.ni

0270	فَكَاتَبَكَ	fakatábaka	fa.ka.'ta.ba.ka	fa.ka."ta.ba.ka
0271	فَكَاتَبَكِ	fakatábaki	fa.ka.'ta.ba.ki	fa.ka."ta.ba.ki

3. DATABASE OF ARABIC PLACE NAMES Part I (DAP-C)

3.1 Synopsis

This **Arabic-English place name database** provides worldwide coverage of common place names, given in standard MSA orthography.

3.2 More Information

Website: <http://www.cjk.org/data/arabic/proper/dapna/>

3.3 Approximate Quantities

Type	Quantities
Total entries	21,000
Bilingual entries	7,000
Monolingual entries	14,000

3.4 The Canonical Component (DAP-C)

3.4.1 Field Description

No.	Field	Description
1	ID	unique ID
2	LEMMA_V	vocalized Arabic place name
3	LEMMA_BW	vocalized Arabic place name in Buckwalter transliteration
4	LEMMA_U	unvocalized Arabic place name
5	ENGLISH	English equivalent
6	REGION	indicates if place is in the Middle East [M] country or region in Arab League [m] country or region is one of the following countries - Afghanistan - Iran - Israel - Pakistan - Turkey

3.4.2 Data Sample

ID	LEMMA_V	LEMMA_BW	LEMMA_U	ENGLISH	REGION
P008289	عَسْقَلَانٌ	EasoqalaAnu	عَسْقَلَانٌ	Ashkelon	m
P016507	أَبْهَا	>abohaA	أَبْهَا	Abha	M
P017340	الْقَاهِرَةُ	AaloqaAhirapu	الْقَاهِرَةُ	Cairo	M
P018064	قِطَاعُ غَزَّة	qiTaAEu gaz~ap	قِطَاعُ غَزَّة	Gaza Strip	m
P018462	الْعَرَاقُ	AaloEiraAqu	الْعَرَاقُ	Iraq	M
P018961	لَندَنُ	lanodanu	لَندَنُ	London	
P019537	نيویورک	niyuwyuwroku	نيویورک	New York	
P019764	اوْسَاكَا	>uwsaAkaA	اوْسَاكَا	Osaka	
P030090	أَبُوظَبِي	>abuwZaboyi	أَبُوظَبِي	Abu Dhabi	M

4. DATABASE OF ARABIC PLACE NAMES Part II (DAP-M, DAP-O, DAP-P)

4.1 Synopsis

Part II of DAP includes the **cliticized forms for the place names** in DAP Part I (DAP-C), both enclitics and proclitics (e.g. وَكَعَمَانٍ).

4.2 More Information

Website: <http://www.cjk.org/data/arabic/nlp/arabic-full-form-lexicon/>

4.3 Approximate Quantities

Type	Quantities
Seed names	21,000
Cliticized forms	6,400,000

4.4 The Morphological Component (DAP-M)

4.4.1 Field Description

If the same schema (field structure) as DAG-M is required, all the fields shown in section 2.4 will be given. Many fields are irrelevant and will get a value of "-". Below are only the fields relevant to DAP.

No.	Field	Description
1	ID	unique foreign key to DAP-C
2	SUBID	unique identifier of cliticized form
3	ARAB_V	vocalized cliticized Arabic place name
4	ARAB_BW	vocalized cliticized Arabic place name in Buckwalter transliteration
5	ARAB_U	unvocalized cliticized Arabic place name
6	LEMMA_V	vocalized Arabic place name linked to DAP-C
7	GEN	gender code [M] masculine [F] feminine
8	NUM	number code [S] singular [D] dual [P] plural [B] broken plural
9	CASE	case ending code [ACU] accusative [GEN] genitive [NOM] nominative
10	PER2	person, number, gender for enclitics [000] zero person [1SC] first person singular [2SM] second person singular masculine [2SF] second person singular feminine [3SM] third person singular masculine [3SF] third person singular feminine [1PC] first person plural [2PM] second person plural masculine [2PF] second person plural feminine [2DC] second person dual [3PM] third person plural masculine [3PF] third person plural feminine [3DM] third person dual masculine [3DF] third person dual feminine
11	DEF	definiteness code [D] explicitly definite by cliticizing of article Al [d] definite by adding pronomial enclitics

		[I] indefinite [G] first term of genitive construct (Idhafa, e.g. <i>baytu in baytu-l-kaatibi</i>)
--	--	--

4.4.2 Data Sample

SUBID	ARAB_V	ARAB_BW	ARAB_U	LEMMA_V	GEN	NUM	CASE	PER	DEF
0085	وَمِصْرُ	wamiSoru	ومصر	مِصْرُ	F	S	NOM	000	G
0086	وَمِصْرِي	wamiSoriy	ومصري	مِصْرُ	F	S	NOM	1SC	d
0087	وَمِصْرَكَ	wamiSoruka	ومصرك	مِصْرُ	F	S	NOM	2SM	d
0088	وَمِصْرَكِ	wamiSoruki	ومصرك	مِصْرُ	F	S	NOM	2SF	d
0089	وَمِصْرُهُ	wamiSoruhu	ومصره	مِصْرُ	F	S	NOM	3SM	d
0090	وَمِصْرُهَا	wamiSoruhA	ومصرها	مِصْرُ	F	S	NOM	3SF	d
0091	وَمِصْرُنَا	wamiSorunaA	ومصرنا	مِصْرُ	F	S	NOM	1PC	d
0092	وَمِصْرُكُمْ	wamiSorukumo	ومصركم	مِصْرُ	F	S	NOM	2PM	d
0093	وَمِصْرُكُنَّ	wamiSorukun~a	ومصركن	مِصْرُ	F	S	NOM	2PF	d
0094	وَمِصْرُكُمَا	wamiSorukumaA	ومصركمـا	مِصْرُ	F	S	NOM	2DC	d
0095	وَمِصْرُهُمْ	wamiSoruhumo	ومصرهم	مِصْرُ	F	S	NOM	3PM	d
0096	وَمِصْرُهُنَّ	wamiSoruhun~a	ومصرهنـا	مِصْرُ	F	S	NOM	3PF	d
0097	وَمِصْرُهُمَا	wamiSoruhumaA	ومصرهما	مِصْرُ	F	S	NOM	3DM	d
0098	وَمِصْرُهُمَا	wamiSoruhumaA	ومصرهما	مِصْرُ	F	S	NOM	3DF	d
0099	وَمِصْرَ	wamiSora	ومصر	مِصْرُ	F	S	GEN	000	G
0100	وَمِصْرِي	wamiSoriy	ومصري	مِصْرُ	F	S	GEN	1SC	d
0101	وَمِصْرَكَ	wamiSoraka	ومصركـا	مِصْرُ	F	S	GEN	2SM	d
0102	وَمِصْرَكِ	wamiSoraki	ومصركــا	مِصْرُ	F	S	GEN	2SF	d
0103	وَمِصْرَهُ	wamiSorahu	ومصرهـا	مِصْرُ	F	S	GEN	3SM	d
0104	وَمِصْرَهَا	wamiSorahaA	ومصرهاـا	مِصْرُ	F	S	GEN	3SF	d
0105	وَمِصْرَنَا	wamiSoranaA	ومصرناـا	مِصْرُ	F	S	GEN	1PC	d

0106	وَمِصْرَ كُمْ	wamiSorakumo	ومصركم	مِصْرُ	F	S	GEN	2PM	d
0107	وَمِصْرَ كُنَّ	wamiSorakun~a	ومصركن	مِصْرُ	F	S	GEN	2PF	d
0108	وَمِصْرَ كُمَا	wamiSorakumaA	ومصركما	مِصْرُ	F	S	GEN	2DC	d
0109	وَمِصْرَ هُمْ	wamiSorahumo	ومصرهم	مِصْرُ	F	S	GEN	3PM	d
0110	وَمِصْرَ هُنَّ	wamiSorahun~a	ومصرهن	مِصْرُ	F	S	GEN	3PF	d
0111	وَمِصْرَ هُمَا	wamiSorahumaA	ومصرهما	مِصْرُ	F	S	GEN	3DM	d
0112	وَمِصْرَ هُمَا	wamiSorahumaA	ومصرهما	مِصْرُ	F	S	GEN	3DF	d
0113	وَمِصْرَ	wamiSora	ومصر	مِصْرُ	F	S	ACU	000	G
0114	وَمِصْرِي	wamiSoriy	ومصري	مِصْرُ	F	S	ACU	1SC	d
0115	وَمِصْرَك	wamiSoraka	ومصرك	مِصْرُ	F	S	ACU	2SM	d
0116	وَمِصْرَكِ	wamiSoraki	ومصرك	مِصْرُ	F	S	ACU	2SF	d
0117	وَمِصْرَهُ	wamiSorahu	ومصره	مِصْرُ	F	S	ACU	3SM	d
0118	وَمِصْرَهَا	wamiSorahaA	ومصرها	مِصْرُ	F	S	ACU	3SF	d
0119	وَمِصْرَنَا	wamiSoranaA	ومصرنا	مِصْرُ	F	S	ACU	1PC	d
0120	وَمِصْرَ كُمْ	wamiSorakumo	ومصركم	مِصْرُ	F	S	ACU	2PM	d
0121	وَمِصْرَ كُنَّ	wamiSorakun~a	ومصركن	مِصْرُ	F	S	ACU	2PF	d
0122	وَمِصْرَ كُمَا	wamiSorakumaA	ومصركما	مِصْرُ	F	S	ACU	2DC	d
0123	وَمِصْرَ هُمْ	wamiSorahumo	ومصرهم	مِصْرُ	F	S	ACU	3PM	d
0124	وَمِصْرَ هُنَّ	wamiSorahun~a	ومصرهن	مِصْرُ	F	S	ACU	3PF	d
0125	وَمِصْرَ هُمَا	wamiSorahumaA	ومصرهما	مِصْرُ	F	S	ACU	3DM	d
0126	وَمِصْرَ هُمَا	wamiSorahumaA	ومصرهما	مِصْرُ	F	S	ACU	3DF	d

4.5 The Orthographical Component (DAP-O)

4.5.1 Field Description

No.	Field	Description
1	ID	uniquely identifies the lemma
2	SUBID	uniquely identifies inflected forms

3	VARID	identifies variants of headword [00] no variants [01] main form [02] variant form
4	ARAB_V	vocalized cliticized Arabic entry
5	VAR_V	variant in vocalized Arabic
6	VAR_U	variant in unvocalized Arabic

4.5.2 Data Sample

Under construction.

4.6 The Phonological Component (DAP-P)

4.6.1 Field Description

No.	Field	Description
1	ID	uniquely identifies the lemma
2	SUBID	uniquely identifies inflected forms
3	VARID	identifies variants of headword [00] no variants [01] main form [02] variant form
4	ARAB_V	vocalized cliticized Arabic entry
5	CARS	cliticized Arabic entry in phonemic transcription (CARS)
6	IPA	phonetic transcription in the International Phonetic Alphabet (IPA)
7	XSAMPA	phonetic transcription in X-SAMPA

4.6.2 Data Sample

SUBID	ARAB_V	CARS	IPA	XSAMPA
0085	وَمَصْرُ	wamíṣru	wa.'miṣ.ru	wa."mis_-.ru
0086	وَمَصْرِي	wamíṣri	wa.'miṣ.ri	wa."mis_-.ri
0087	وَمَصْرُكَ	wamíṣruka	wa.'miṣ.ru.ka	wa."mis_-.ru.ka
0088	وَمَصْرُكِ	wamíṣruki	wa.'miṣ.ru.ki	wa."mis_-.ru.ki
0089	وَمَصْرُهُ	wamíṣruhu	wa.'miṣ.ru.hu	wa."mis_-.ru.hu

0090	وَمِصْرُهَا	wamışruha	wa.'miş.ru.ha	wa."mis_-.ru.ha
0091	وَمِصْرُنَا	wamışruna	wa.'miş.ru.na	wa."mis_-.ru.na
0092	وَمِصْرُكُمْ	wamışrukum	wa.'miş.ru.kum	wa."mis_-.ru.kum
0093	وَمِصْرُكُنَّ	wamışrukúnna	wa.mış.ru.'ku.n:a	wa.mis_-.ru."ku.n:a
0094	وَمِصْرُكُمَا	wamışrukumə	wa.mış.'ru.ku.ma	wa.mis_-.ru.ku.ma
0095	وَمِصْرُهُمْ	wamışruhum	wa.'miş.ru.hum	wa."mis_-.ru.hum
0096	وَمِصْرُهُنَّ	wamışruhúnna	wa.mış.ru.'hu.n:a	wa.mis_-.ru."hu.n:a
0097	وَمِصْرُهُمَا	wamışruhumə	wa.mış.'ru.hu.ma	wa.mis_-.ru.hu.ma
0098	وَمِصْرُهُمَا	wamışruhumə	wa.mış.'ru.hu.ma	wa.mis_-.ru.hu.ma
0099	وَمِصْرَ	wamışra	wa.'miş.ra	wa."mis_-.ra
0100	وَمِصْرِي	wamışri	wa.'miş.ri	wa."mis_-.ri
0101	وَمِصْرَكَ	wamışraka	wa.'miş.ra.ka	wa."mis_-.ra.ka
0102	وَمِصْرَكِ	wamışraki	wa.'miş.ra.ki	wa."mis_-.ra.ki
0103	وَمِصْرَهُ	wamışrahu	wa.'miş.ra.hu	wa."mis_-.ra.hu
0104	وَمِصْرَهَا	wamışrahə	wa.'miş.ra.ha	wa."mis_-.ra.ha
0105	وَمِصْرَنَا	wamışrana	wa.'miş.ra.na	wa."mis_-.ra.na
0106	وَمِصْرَكُمْ	wamışrakum	wa.'miş.ra.kum	wa."mis_-.ra.kum
0107	وَمِصْرُكُنَّ	wamışrakúnna	wa.mış.ra.'ku.n:a	wa.mis_-.ra."ku.n:a
0108	وَمِصْرُكُمَا	wamışrákumə	wa.mış.'ra.ku.ma	wa.mis_-.ra.ku.ma
0109	وَمِصْرُهُمْ	wamışrahum	wa.'miş.ra.hum	wa."mis_-.ra.hum
0110	وَمِصْرُهُنَّ	wamışrahúnna	wa.mış.ra.'hu.n:a	wa.mis_-.ra."hu.n:a
0111	وَمِصْرُهُمَا	wamışráhumə	wa.mış.'ra.hu.ma	wa.mis_-.ra.hu.ma
0112	وَمِصْرُهُمَا	wamışráhumə	wa.mış.'ra.hu.ma	wa.mis_-.ra.hu.ma
0113	وَمِصْرَ	wamışra	wa.'miş.ra	wa."mis_-.ra
0114	وَمِصْرِي	wamışri	wa.'miş.ri	wa."mis_-.ri
0115	وَمِصْرَكَ	wamışraka	wa.'miş.ra.ka	wa."mis_-.ra.ka

0116	وَمِصْرَكٌ	wamísraki	wa.'mis.ra.ki	wa."mis_-.ra.ki
0117	وَمِصْرَهُ	wamísrahu	wa.'mis.ra.hu	wa."mis_-.ra.hu
0118	وَمِصْرَهَا	wamísrahā	wa.'mis.ra.ha	wa."mis_-.ra.ha
0119	وَمِصْرَنَا	wamísranā	wa.'mis.ra.na	wa."mis_-.ra.na
0120	وَمِصْرَكُمْ	wamísrakum	wa.'mis.ra.kum	wa."mis_-.ra.kum
0121	وَمِصْرَكُنَّ	wamísrakúnna	wa.mis.ra.'ku.n:a	wa.mis_-.ra."ku.n:a
0122	وَمِصْرَكُمَا	wamísrákumā	wa.mis.'ra.ku.ma	wa.mis_-.ra."ra.ku.ma
0123	وَمِصْرَهُمْ	wamísrahum	wa.'mis.ra.hum	wa."mis_-.ra.hum
0124	وَمِصْرَهُنَّ	wamísrahúnna	wa.mis.ra.'hu.n:a	wa.mis_-.ra."hu.n:a
0125	وَمِصْرَهُمَا	wamísráhumā	wa.mis.'ra.hu.ma	wa.mis_-.ra.hu.ma
0126	وَمِصْرَهُمَا	wamísráhuma	wa.mis.'ra.hu.ma	wa.mis_-.ra.hu.ma

5. DATABASE OF ARABIC FOREIGN NAMES Part I (DAF-C)

5.1 Synopsis

This database covers non-Arab personal names, their Arabic equivalents, some Arabic script variants and English equivalents. All DAF data can be provided in vocalized or unvocalized formats.

5.2 More Information

Website: <http://www.cjk.org/data/arabic/proper/dafna/>

5.3 Approximate Quantities

The approximate number of entries is as follows:

Type	Quantities
Total entries	400,000
Seed names	223,367
Arabic variants	176,633

5.4 The Canonical Component (DAF-C)

5.4.1 Field Description

No.	Field	Description
1	ID	unique ID identifying entry block (unique on English)
2	LEMMA_V	vocalized non-Arab personal name
3	LEMMA_BW	vocalized non-Arab personal name in Buckwalter transliteration
4	LEMMA_U	unvocalized non-Arab personal name
5	ENGLISH	English original
6	TYPE	type of name [G] given name [S] surname [GS] given name and surname
7	GEN	gender of name [M] male [F] female [MF] male and female [-] not applicable
8	RS_FREQ	frequency of Latin surname based on US Census
9	RG_FREQ	frequency of Latin given name based on Social Security statistics

5.4.2 Data Sample

ID	LEMMA_V	LEMMA_BW	LEMMA_U	ENGLISH	TYPE	GEN	RS_FREQ	RG_FREQ
F086620	هالبرن	haAlobiron	هالبرن	Halpern	S	-	0004121	-
F098222	إزاربيلا	<izaAbiylaA	إزاربيلا	Izabella	G	F	-	0025717
F098611	جاك	jaAk	جاك	Jack	GS	MF	0015256	0696625
F100877	جانيت	jaAniyt	جانيت	Janet	GS	MF	0000437	0557605
F107324	جوليت	juwliyt	جوليت	Juliet	G	F	-	0030202
F170757	بيترسون	biytirosuwn	بيترسون	Peterson	GS	M	0278297	0000756
F193601	شمیت	\$omiyt	شمیت	Schmidt	S	-	0147034	-
F204232	سمیث	somiyy	سمیث	Smith	GS	MF	2442977	0004733
F234748	ویلیام	wiyliyaAm	ویلیام	William	GS	MF	0013373	4133327

6. DATABASE OF ARABIC FOREIGN NAMES Part II (DAF-M, DAF-O, DAF-P)

6.1 Synopsis

Part II of DAF includes the **cliticized forms for the non-Arab personal names**. Some cliticized forms may be relatively rare but they are valid.

6.2 More Information

Website: <http://www.cjk.org/data/arabic/nlp/arabic-full-form-lexicon/>

6.3 Approximate Quantities

Type	Quantities
Seed names	223,367
Cliticized forms	226,784,907

6.4 The Morphological Component (DAF-M)

6.4.1 Field Description

If the same schema (field structure) as DAG-M is required, all the fields shown in section 2.4 will be given. Some fields are irrelevant and will get a value of "-". Below are only the fields relevant to DAF.

No.	Field	Description
1	ID	unique foreign key to DAF-C
2	SUBID	unique identifier of cliticized form
3	ARAB_V	vocalized cliticized non-Arab personal name
4	ARAB_BW	vocalized cliticized non-Arab personal name in Buckwalter transliteration
5	ARAB_U	unvocalized cliticized non-Arab personal name
6	LEMMA_V	vocalized non-Arab personal name linked to DAF-C
7	GEN	gender code [M] masculine [F] feminine [C] common gender
8	NUM	number code [S] singular [D] dual [P] plural [B] broken plural
9	CASE	case ending code [ACU] accusative [GEN] genitive

		[NOM] nominative
10	PER2	<p>person, number, gender for enclitics</p> <p>[000] zero person [1SC] first person singular [2SM] second person singular masculine [2SF] second person singular feminine [3SM] third person singular masculine [3SF] third person singular feminine [1PC] first person plural [2PM] second person plural masculine [2PF] second person plural feminine [2DC] second person dual [3PM] third person plural masculine [3PF] third person plural feminine [3DM] third person dual masculine [3DF] third person dual feminine</p>
11	DEF	<p>definiteness code</p> <p>[D] explicitly definite by cliticizing of article Al [d] definite by adding pronomial enclitics [I] indefinite [G] first term of genitive construct (Idhafa, e.g. <i>baytu in baytu-l-kaatibi</i>)</p>

6.4.2 Data Sample

SUBID	ARAB_V	ARAB_BW	ARAB_U	LEMMA_V	GEN	NUM	CASE	PER2	DEF
0083	وَجَأْكُ	wajaAku	وَجَأْكُ	وَجَأْكُ	C	S	NOM	000	G
0084	وَجَأْكِي	wajaAkiy	وَجَأْكِي	وَجَأْكِي	C	S	NOM	1SC	d
0085	وَجَأْكُكُ	wajaAkuka	وَجَأْكُكُ	وَجَأْكُكُ	C	S	NOM	2SM	d
0086	وَجَأْكُكِي	wajaAkuki	وَجَأْكُكِي	وَجَأْكُكِي	C	S	NOM	2SF	d
0087	وَجَأْكُكُهُ	wajaAkuhu	وَجَأْكُكُهُ	وَجَأْكُكُهُ	C	S	NOM	3SM	d
0088	وَجَأْكُكُهَا	wajaAkuhaA	وَجَأْكُكُهَا	وَجَأْكُكُهَا	C	S	NOM	3SF	d
0089	وَجَأْكُكُنَا	wajaAkunaA	وَجَأْكُكُنَا	وَجَأْكُكُنَا	C	S	NOM	1PC	d
0090	وَجَأْكُكُمْ	wajaAkukumo	وَجَأْكُكُمْ	وَجَأْكُكُمْ	C	S	NOM	2PM	d
0091	وَجَأْكُكُنْ	wajaAkukun~a	وَجَأْكُكُنْ	وَجَأْكُكُنْ	C	S	NOM	2PF	d
0092	وَجَأْكُكُمَا	wajaAkukumaA	وَجَأْكُكُمَا	وَجَأْكُكُمَا	C	S	NOM	2DC	d
0093	وَجَأْكُكُهُمْ	wajaAkuhumo	وَجَأْكُكُهُمْ	وَجَأْكُكُهُمْ	C	S	NOM	3PM	d

0094	وَجَأْكُهُنَّ	wajaAkuhun~a	وجاکهن	جاک	C	S	NOM	3PF	d
0095	وَجَأْكُهُمَا	wajaAkuhumaA	وجاکهما	جاک	C	S	NOM	3DM	d
0096	وَجَأْكُهُمَا	wajaAkuhumaA	وجاکهما	جاک	C	S	NOM	3DF	d
0097	وَجَأِكِي	wajaAkiy	وجاکي	جاک	C	S	GEN	1SC	d
0098	وَجَأِكِاڭ	wajaAkika	وجاکاڭ	جاک	C	S	GEN	2SM	d
0099	وَجَأِكِاڭ	wajaAkiki	وجاکاڭ	جاک	C	S	GEN	2SF	d
0100	وَجَأِكِە	wajaAkihi	وجاکە	جاک	C	S	GEN	3SM	d
0101	وَجَأِكِەها	wajaAkihaA	وجاکەها	جاک	C	S	GEN	3SF	d
0102	وَجَأِكِىنا	wajaAkinaA	وجاکىنا	جاک	C	S	GEN	1PC	d
0103	وَجَأِكِىمْ	wajaAkikumo	وجاکىمْ	جاک	C	S	GEN	2PM	d
0104	وَجَأِكِىنْ	wajaAkikun~a	وجاکىنْ	جاک	C	S	GEN	2PF	d
0105	وَجَأِكِىمَا	wajaAkikumaA	وجاکىمَا	جاک	C	S	GEN	2DC	d
0106	وَجَأِكِىھِمْ	wajaAkihimo	وجاکىھِمْ	جاک	C	S	GEN	3PM	d
0107	وَجَأِكِىھُنَّ	wajaAkihin~a	وجاکىھُنَّ	جاک	C	S	GEN	3PF	d
0108	وَجَأِكِىھُمَا	wajaAkihimaA	وجاکىھُمَا	جاک	C	S	GEN	3DM	d
0109	وَجَأِكِىھُمَا	wajaAkihimaA	وجاکىھُمَا	جاک	C	S	GEN	3DF	d
0110	وَجَأِكِى	wajaAka	وجاک	جاک	C	S	ACU	000	G
0111	وَجَأِكِي	wajaAkiy	وجاکي	جاک	C	S	ACU	1SC	d
0112	وَجَأِكِىڭ	wajaAkaka	وجاکاڭ	جاک	C	S	ACU	2SM	d

6.5 The Orthographical Component (DAF-O)

6.5.1 Field Description

No.	Field	Description
1	ID	uniquely identifies the lemma
2	SUBID	uniquely identifies inflected forms
3	VARID	identifies variants of headword [00] no variants [01] main form [02] variant form

4	ARAB_V	vocalized cliticized Arabic entry
5	VAR_V	variant in vocalized Arabic
6	VAR_U	variant in unvocalized Arabic

6.5.2 Data Sample

Under construction.

6.6 The Phonological Component (DAF-P)

6.6.1 Field Description

No.	Field	Description
1	ID	uniquely identifies the lemma
2	SUBID	uniquely identifies inflected forms
3	VARID	identifies variants of headword [00] no variants [01] main form [02] variant form
4	ARAB_V	vocalized cliticized Arabic entry
5	CARS	cliticized Arabic entry in phonemic transcription (CARS)
6	IPA	phonetic transcription in the International Phonetic Alphabet (IPA)
7	XSAMPA	phonetic transcription in X-SAMPA

6.6.2 Data Sample

SUBID	ARAB_V	CARS	IPA	XSAMPA
0083	وَجَأْكُ	wajāku	wa.'ʒa..ku	wa."Za:.ku
0084	وَجَأْكِي	wajāki	wa.'ʒa..ki	wa."Za:.ki
0085	وَجَأْكُكَ	wajākuka	wa.'ʒa..ku.ka	wa."Za:.ku.ka
0086	وَجَأْكِكَ	wajākuki	wa.'ʒa..ku.ki	wa."Za:.ku.ki
0087	وَجَأْكُهُ	wajākuhu	wa.'ʒa..ku.hu	wa."Za:.ku.hu
0088	وَجَأْكُهَا	wajākuha	wa.'ʒa..ku.ha	wa."Za:.ku.ha
0089	وَجَأْكُنَا	wajākuna	wa.'ʒa..ku.na	wa."Za:.ku.na
0090	وَجَأْكُمْ	wajākum	wa.'ʒa..ku.kum	wa."Za:.ku.kum

0091	وَجَأْكُنْ	wajākukúnna	wa.ʒa:.ku.'ku.n:a	wa.Za:.ku."ku.n:a
0092	وَجَأْكُمَا	wajākúkumā	wa.ʒa:.ku.ku.ma	wa.Za:.ku.ku.ma
0093	وَجَأْكُهُمْ	wajākuhum	wa.'ʒa:.ku.hum	wa."Za:.ku.hum
0094	وَجَأْكُهُنْ	wajākuhúnna	wa.ʒa:.ku.'hu.n:a	wa.Za:.ku."hu.n:a
0095	وَجَأْكُهُمَا	wajākúhumā	wa.ʒa:.ku.hu.ma	wa.Za:.ku.hu.ma
0096	وَجَأْكُهُمَا	wajākúhuma	wa.ʒa:.ku.hu.ma	wa.Za:.ku.hu.ma
0097	وَجَأِكِي	wajāki	wa.'ʒa:.ki	wa."Za:.ki
0098	وَجَأِكِكِ	wajākika	wa.'ʒa:.ki ka	wa."Za:.ki ka
0099	وَجَأِكِكِ	wajākiki	wa.'ʒa:.ki.ki	wa."Za:.ki.ki
0100	وَجَأِكِهِ	wajākihi	wa.'ʒa:.ki.hi	wa."Za:.ki.hi
0101	وَجَأِكِهَا	wajākiha	wa.'ʒa:.ki.ha	wa."Za:.ki.ha
0102	وَجَأِكِنَا	wajākinā	wa.'ʒa:.ki.na	wa."Za:.ki.na
0103	وَجَأِكِكُمْ	wajākikum	wa.'ʒa:.ki.kum	wa."Za:.ki.kum
0104	وَجَأِكِكُنْ	wajākikúnna	wa.ʒa:.ki.'ku.n:a	wa.Za:.ki."ku.n:a
0105	وَجَأِكِكُمَا	wajākíkumā	wa.ʒa:.ki.ku.ma	wa.Za:.ki.ku.ma
0106	وَجَأِكِهِمْ	wajākihim	wa.'ʒa:.ki.him	wa."Za:.ki.him
0107	وَجَأِكِهِنْ	wajākihínna	wa.ʒa:.ki.'hi.n:a	wa.Za:.ki."hi.n:a
0108	وَجَأِكِهِمَا	wajākíhimā	wa.ʒa:.ki.hi.ma	wa.Za:.ki.hi.ma
0109	وَجَأِكِهِمَا	wajākíhima	wa.ʒa:.ki.hi.ma	wa.Za:.ki.hi.ma
0110	وَجَأِكِ	wajāka	wa.'ʒa:.ka	wa."Za:.ka
0111	وَجَأِكِي	wajāki	wa.'ʒa:.ki	wa."Za:.ki
0112	وَجَأِكِكِ	wajākaka	wa.'ʒa:.ka ka	wa."Za:.ka ka

7. DATABASE OF ARABIC NAMES Part I (DAN-C)

7.1 Synopsis

Very comprehensive database of **Arab personal names** and romanized name variants

with a variety of supplementary information.

7.2 More Information

- Website:** <http://www.cjk.org/data/arabic/proper/dan/>
White paper: <http://www.cjk.org/cjk/arabic/DAN3.pdf>
Research paper: <http://www.cjk.org/cjk/reference/danpaper.pdf>
Online demo: <http://cjk1.dyndns.org/ante/ante.php> (login "dandana")

7.3 Approximate Quantities

Type	Quantities
Seed names	unvocalized: 76,000 vocalized: approx 100,000
Romanized variants	6,595,008

7.4 The Canonical Component (DAN-C)

7.4.1 Field Description

No.	Field	Description
1	ID	unique identifier for vocalized Arab personal name
2	U_ID	unique identifier for unvocalized Arab personal name
3	SUBID	identifies members of a variants group
4	R_NAME	romanized variant of ARABIC
5	LEMMA_V	vocalized Arab personal name
6	LEMMA_BW	vocalized Arab personal name in Buckwalter transliteration
7	LEMMA_U	unvocalized Arab personal name
8	TYPE	type of name [G] given name [S] surname [GS] given name and surname [U] unknown
9	GEN	gender of name [M] male [F] female [MF] male and female [U] unknown [-] not applicable
10	R_TYPE	type of romanization [I] IC (Intelligence Community) Standard [V] variant of undetermined Romanization system
11	R_FREQ	web frequency as a string, not necessarily representative of web frequency as a name

7.4.2 Data Sample

SUBID	R_NAME	LEMMA_V	LEMMA_BW	TYPE	GEN	R_TYPE	R_FREQ
001683	Abdurrahman	عبد الرَّحْمَن	Eabod Aalr~aHom`n	G	M	V	0002680000
001421	Abderrahmane	عبد الرَّحْمَن	Eabod Aalr~aHom`n	G	M	V	0000603000
001647	Abdulrahman	عبد الرَّحْمَن	Eabod Aalr~aHom`n	G	M	V	0000589000
001617	Abdul Rahman	عبد الرَّحْمَن	Eabod Aalr~aHom`n	G	M	V	0000478000
001387	Abdelrahman	عبد الرَّحْمَن	Eabod Aalr~aHom`n	G	M	V	0000348000
001420	Abderrahman	عبد الرَّحْمَن	Eabod Aalr~aHom`n	G	M	V	0000206000
001366	Abdel Rahman	عبد الرَّحْمَن	Eabod Aalr~aHom`n	G	M	V	0000136000
000946	Abd Al Rahman	عبد الرَّحْمَن	Eabod Aalr~aHom`n	G	M	V	0000120000
001671	Abdurahman	عبد الرَّحْمَن	Eabod Aalr~aHom`n	G	M	V	0000074000
001407	Abderahmane	عبد الرَّحْمَن	Eabod Aalr~aHom`n	G	M	V	0000061900
001660	Abdur Rahman	عبد الرَّحْمَن	Eabod Aalr~aHom`n	G	M	V	0000056400
001472	Abdirahman	عبد الرَّحْمَن	Eabod Aalr~aHom`n	G	M	V	0000043700
001625	Abdul Rehman	عبد الرَّحْمَن	Eabod Aalr~aHom`n	G	M	V	0000040700
001013	Abd Elrahman	عبد الرَّحْمَن	Eabod Aalr~aHom`n	G	M	V	0000035800
001665	Abdur Rehman	عبد الرَّحْمَن	Eabod Aalr~aHom`n	G	M	V	0000023600
001388	Abdelrahmane	عبد الرَّحْمَن	Eabod Aalr~aHom`n	G	M	V	0000017300
000996	Abd El Rahman	عبد الرَّحْمَن	Eabod Aalr~aHom`n	G	M	V	0000013970
000968	Abd Ar Rahman	عبد الرَّحْمَن	Eabod Aalr~aHom`n	G	M	V	0000012100
001406	Abderahman	عبد الرَّحْمَن	Eabod Aalr~aHom`n	G	M	V	0000010000
001654	Abdulrehman	عبد الرَّحْمَن	Eabod Aalr~aHom`n	G	M	V	0000009200
001574	Abdourahmane	عبد الرَّحْمَن	Eabod Aalr~aHom`n	G	M	V	0000008740
001681	Abdurrachman	عبد الرَّحْمَن	Eabod Aalr~aHom`n	G	M	V	0000008570
001674	Abdurakhman	عبد الرَّحْمَن	Eabod Aalr~aHom`n	G	M	V	0000008000

001668	Abdurachman	عبد الرَّحْمَن	Eabod Aalr~aHom`n	G	M	V	0000007330
001514	Abdolrahman	عبد الرَّحْمَن	Eabod Aalr~aHom`n	G	M	V	0000007220
001481	Abdirahman	عبد الرَّحْمَن	Eabod Aalr~aHom`n	G	M	V	0000007070
001573	Abdourahman	عبد الرَّحْمَن	Eabod Aalr~aHom`n	G	M	V	0000006600
001613	Abdul Rachman	عبد الرَّحْمَن	Eabod Aalr~aHom`n	G	M	V	0000005240
001542	Abdorrahman	عبد الرَّحْمَن	Eabod Aalr~aHom`n	G	M	V	0000004510
001645	Abdulrachman	عبد الرَّحْمَن	Eabod Aalr~aHom`n	G	M	V	0000004350

8. DATABASE OF ARABIC NAMES Part II (DAN-M, DAN-O, DAN-P)

8.1 Synopsis

Part II of DAN that includes the **cliticized forms for Arab personal names**. Some encliticized forms may be relatively rare but they are valid.

8.2 More Information

Website: <http://www.cjk.org/data/arabic/nlp/arabic-full-form-lexicon/>
White paper: ArabLEX.pdf

8.3 Approximate Quantities

Type	Quantities
Seed names	unvocalized: 76,000 vocalized: approx 100,000
Criticized forms	218,215,875

8.4 The Morphological Component (DAN-M)

8.4.1 Field Description

If the same schema (field structure) as DAG-M is required, all the fields shown in section 2.4 will be given. Many fields are irrelevant and will get a value of "-". Below are only the fields relevant to DAN.

No.	Field	Description
1	ID	unique foreign key to DAN-C

2	SUBID	unique identifier for cliticized forms
3	ARAB_V	vocalized cliticized Arab personal name
4	ARAB_BW	vocalized cliticized Arab personal name in Buckwalter transliteration
5	ARAB_U	unvocalized cliticized Arab personal name
6	LEMMA_V	vocalized Arab personal name linked to DAN-C
7	GEN	gender code [M] masculine [F] feminine [C] common gender
8	NUM	number code [S] singular [D] dual [P] plural [B] broken plural
9	CASE	case ending code [ACU] accusative [GEN] genitive [NOM] nominative
10	PER2	person, number, gender for enclitics [000] zero person [1SC] first person singular [2SM] second person singular masculine [2SF] second person singular feminine [3SM] third person singular masculine [3SF] third person singular feminine [1PC] first person plural [2PM] second person plural masculine [2PF] second person plural feminine [2DC] second person dual [3PM] third person plural masculine [3PF] third person plural feminine [3DM] third person dual masculine [3DF] third person dual feminine
11	DEF	definiteness code [D] explicitly definite by cliticizing of article Al [d] definite by adding pronomial enclitics [I] indefinite [G] first term of genitive construct (Idhafa, e.g. baytu in baytu-l-kaatibi)

8.4.2 Data Sample

SUBID	ARAB_V	ARAB_BW	ARAB_U	LEMMA_V	GEN	NUM	CASE	PER2	DEF
0089	وَمُحَمَّدٌ	wamuHam~adN	ومحمد	محمد	C	S	NOM	000	I

0090	وَمُحَمَّدُ	wamuHam~adu	ومحمد	مُحَمَّدٌ	C	S	NOM	000	G
0091	وَمُحَمَّدِي	wamuHam~adiy	ومحمدي	مُحَمَّدٌ	C	S	NOM	1SC	d
0092	وَمُحَمَّدُكَ	wamuHam~aduka	ومحمدك	مُحَمَّدٌ	C	S	NOM	2SM	d
0093	وَمُحَمَّدُكِ	wamuHam~aduki	ومحمدكِ	مُحَمَّدٌ	C	S	NOM	2SF	d
0094	وَمُحَمَّدُهُ	wamuHam~aduhu	ومحمده	مُحَمَّدٌ	C	S	NOM	3SM	d
0095	وَمُحَمَّدُهَا	wamuHam~aduhaA	ومحمدها	مُحَمَّدٌ	C	S	NOM	3SF	d
0096	وَمُحَمَّدُنَا	wamuHam~adunaA	ومحمدنا	مُحَمَّدٌ	C	S	NOM	1PC	d
0097	وَمُحَمَّدُكُمْ	wamuHam~adukumo	ومحمدكم	مُحَمَّدٌ	C	S	NOM	2PM	d
0098	وَمُحَمَّدُكُنْ	wamuHam~adukun~a	ومحمدكن	مُحَمَّدٌ	C	S	NOM	2PF	d
0099	وَمُحَمَّدُكُمَا	wamuHam~adukumaA	ومحمدكمَا	مُحَمَّدٌ	C	S	NOM	2DC	d
0100	وَمُحَمَّدُهُمْ	wamuHam~aduhumo	ومحمدهم	مُحَمَّدٌ	C	S	NOM	3PM	d
0101	وَمُحَمَّدُهُنْ	wamuHam~aduhun~a	ومحمدهن	مُحَمَّدٌ	C	S	NOM	3PF	d
0102	وَمُحَمَّدُهُمَا	wamuHam~aduhumaA	ومحمدهمَا	مُحَمَّدٌ	C	S	NOM	3DM	d
0103	وَمُحَمَّدُهُمَا	wamuHam~aduhumaA	ومحمدهمَا	مُحَمَّدٌ	C	S	NOM	3DF	d
0104	وَمُحَمَّدٍ	wamuHam~adK	ومحمد	مُحَمَّدٌ	C	S	GEN	000	I
0105	وَمُحَمَّدٍ	wamuHam~adi	ومحمد	مُحَمَّدٌ	C	S	GEN	000	G
0106	وَمُحَمَّدِي	wamuHam~adiy	ومحمدي	مُحَمَّدٌ	C	S	GEN	1SC	d
0107	وَمُحَمَّدِكَ	wamuHam~adika	ومحمدكَ	مُحَمَّدٌ	C	S	GEN	2SM	d
0108	وَمُحَمَّدِكِ	wamuHam~adiki	ومحمدكِ	مُحَمَّدٌ	C	S	GEN	2SF	d
0109	وَمُحَمَّدِهِ	wamuHam~adihi	ومحمدهِ	مُحَمَّدٌ	C	S	GEN	3SM	d
0110	وَمُحَمَّدِهَا	wamuHam~adihaA	ومحمدها	مُحَمَّدٌ	C	S	GEN	3SF	d
0111	وَمُحَمَّدِنَا	wamuHam~adinaA	ومحمدنا	مُحَمَّدٌ	C	S	GEN	1PC	d
0112	وَمُحَمَّدِكُمْ	wamuHam~adikumo	ومحمدكم	مُحَمَّدٌ	C	S	GEN	2PM	d
0113	وَمُحَمَّدِكُنْ	wamuHam~adikun~a	ومحمدكن	مُحَمَّدٌ	C	S	GEN	2PF	d
0114	وَمُحَمَّدِكُمَا	wamuHam~adikumaA	ومحمدكمَا	مُحَمَّدٌ	C	S	GEN	2DC	d
0115	وَمُحَمَّدِهِمْ	wamuHam~adihimo	ومحمدهمِ	مُحَمَّدٌ	C	S	GEN	3PM	d

0116	وَمُحَمَّدٌ هِنْ	wamuHam~adihin~a	ومحمدـهـن	مُحَمَّد	C	S	GEN	3PF	d
0117	وَمُحَمَّدٌ هِمَا	wamuHam~adihimaA	ومحمدـهـمـا	مُحَمَّد	C	S	GEN	3DM	d
0118	وَمُحَمَّدٌ هِمَا	wamuHam~adihimaA	ومحمدـهـمـا	مُحَمَّد	C	S	GEN	3DF	d
0119	وَمُحَمَّدًا	wamuHam~adFA	ومحمدـا	مُحَمَّد	C	S	ACU	000	I
0120	وَمُحَمَّدَ	wamuHam~ada	ومحمدـ	مُحَمَّد	C	S	ACU	000	G

8.5 The Orthographical Component (DAN-O)

8.5.1 Field Description

No.	Field	Description
1	ID	uniquely identifies the lemma
2	SUBID	uniquely identifies inflected forms
3	VARID	identifies variants of headword [00] no variants [01] main form [02] variant form
4	ARAB_V	vocalized cliticized Arabic entry
5	VAR_V	variant in vocalized Arabic
6	VAR_U	variant in unvocalized Arabic

8.5.2 Data Sample

Under construction.

8.6 The Phonological Component (DAN-P)

8.6.1 Field Description

No.	Field	Description
1	ID	uniquely identifies the lemma
2	SUBID	uniquely identifies inflected forms
3	VARID	identifies variants of headword [00] no variants [01] main form [02] variant form
4	ARAB_V	vocalized cliticized Arabic entry
5	CARS	cliticized Arabic entry in phonemic transcription (CARS)

6	IPA	phonetic transcription in the International Phonetic Alphabet (IPA)
7	XSAMPA	phonetic transcription in X-SAMPA

8.6.2 Data Sample

SUBID	ARAB_V	CARS	IPA	XSAMPA
0089	وَمُحَمَّدٌ	wamuḥámmadun	wa.mu.'ḥa.m:a.dun	wa.mu."Xla.m:a.dun
0090	وَمُحَمَّدٌ	wamuḥámmadu	wa.mu.'ḥa.m:a.du	wa.mu."Xla.m:a.du
0091	وَمُحَمَّدِي	wamuḥámmadi	wa.mu.'ḥa.m:a.di	wa.mu."Xla.m:a.di
0092	وَمُحَمَّدِكَ	wamuḥammáduka	wa.mu.ḥa.'m:a.du.ka	wa.mu.Xla."m:a.du.ka
0093	وَمُحَمَّدِكَ	wamuḥammáduki	wa.mu.ḥa.'m:a.du.ki	wa.mu.Xla."m:a.du.ki
0094	وَمُحَمَّدُهُ	wamuḥammáduhu	wa.mu.ḥa.'m:a.du.hu	wa.mu.Xla."m:a.du.hu
0095	وَمُحَمَّدُهَا	wamuḥammáduhā	wa.mu.ḥa.'m:a.du.ha	wa.mu.Xla."m:a.du.ha
0096	وَمُحَمَّدُنَا	wamuḥammáduna	wa.mu.ḥa.'m:a.du.na	wa.mu.Xla."m:a.du.na
0097	وَمُحَمَّدُكُمْ	wamuḥammádukum	wa.mu.ḥa.'m:a.du.kum	wa.mu.Xla."m:a.du.kum
0098	وَمُحَمَّدُكُنْ	wamuḥammadukúnna	wa.mu.ḥa.m:a.du.'ku.n:a	wa.mu.Xla.m:a.du."ku.n:a
0099	وَمُحَمَّدُكُمَا	wamuḥammadúkuma	wa.mu.ḥa.m:a.'du.ku.ma	wa.mu.Xla.m:a."du.ku.ma
0100	وَمُحَمَّدُهُمْ	wamuḥammáduhum	wa.mu.ḥa.'m:a.du.hum	wa.mu.Xla."m:a.du.hum
0101	وَمُحَمَّدُهُنْ	wamuḥammaduhúnna	wa.mu.ḥa.m:a.du.'hu.n:a	wa.mu.Xla.m:a.du."hu.n:a
0102	وَمُحَمَّدُهُمَا	wamuḥammadúhumā	wa.mu.ḥa.m:a.'du.hu.ma	wa.mu.Xla.m:a."du.hu.ma
0103	وَمُحَمَّدُهُمَا	wamuḥammadúhumā	wa.mu.ḥa.m:a.'du.hu.ma	wa.mu.Xla.m:a."du.hu.ma
0104	وَمُحَمَّدٍ	wamuḥámmadin	wa.mu.'ḥa.m:a.din	wa.mu."Xla.m:a.din
0105	وَمُحَمَّدٍ	wamuḥámmadi	wa.mu.'ḥa.m:a.di	wa.mu."Xla.m:a.di
0106	وَمُحَمَّدِي	wamuḥámmadi	wa.mu.'ḥa.m:a.di	wa.mu."Xla.m:a.di
0107	وَمُحَمَّدِكَ	wamuḥammádika	wa.mu.ḥa.'m:a.di.ka	wa.mu.Xla."m:a.di.ka
0108	وَمُحَمَّدِكَ	wamuḥammádiki	wa.mu.ḥa.'m:a.di.ki	wa.mu.Xla."m:a.di.ki
0109	وَمُحَمَّدِهِ	wamuḥammádihi	wa.mu.ḥa.'m:a.di.hi	wa.mu.Xla."m:a.di.hi

0110	وَمُحَمَّدٌ هَا	wamuḥammádiḥa	wa.mu.ḥa.'m:a.di.ha	wa.mu.Xla."m:a.di.ha
0111	وَمُحَمَّدٌ نَا	wamuḥammádinā	wa.mu.ḥa.'m:a.di.na	wa.mu.Xla."m:a.di.na
0112	وَمُحَمَّدٌ كُمْ	wamuḥammádikum	wa.mu.ḥa.'m:a.di.kum	wa.mu.Xla."m:a.di.kum
0113	وَمُحَمَّدٌ كُنْ	wamuḥammadikúnna	wa.mu.ḥa.m:a.di.'ku.n:a	wa.mu.Xla.m:a.di."ku.n:a
0114	وَمُحَمَّدٌ كُمَا	wamuḥammadíkumā	wa.mu.ḥa.m:a.'di.ku.ma	wa.mu.Xla.m:a."di.ku.ma
0115	وَمُحَمَّدٌ هُمْ	wamuḥammádihim	wa.mu.ḥa.'m:a.di.him	wa.mu.Xla."m:a.di.him
0116	وَمُحَمَّدٌ هِنْ	wamuḥammadihínna	wa.mu.ḥa.m:a.di.'hi.n:a	wa.mu.Xla.m:a.di."hi.n:a
0117	وَمُحَمَّدٌ هِمَا	wamuḥammadíhimā	wa.mu.ḥa.m:a.'di.hi.ma	wa.mu.Xla.m:a."di.hi.ma
0118	وَمُحَمَّدٌ هِمَا	wamuḥammadíhimā	wa.mu.ḥa.m:a.'di.hi.ma	wa.mu.Xla.m:a."di.hi.ma
0119	وَمُحَمَّدًا	wamuḥámmadan	wa.mu.'ḥa.m:a.dan	wa.mu."Xla.m:a.dan
0120	وَمُحَمَّدَ	wamuḥámmada	wa.mu.'ḥa.m:a.da	wa.mu."Xla.m:a.da