

The Challenges and Pitfalls of Arabic Romanization and Arabization

Jack Halpern

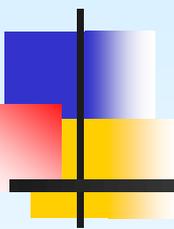
CEO

The CJK Dictionary Institute, Inc.

المؤسسة المعجمية للغات الشرقية

日中韓辭典研究所

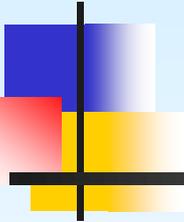
漢



Automatic Romanizer of Arabic Names

ARAN

الناقل اللغوي الآلي للأسماء العربي



Non-Arabic Name Arabizer

نفع NANA

نقل عربي نقل عربي

Transcription and Transliteration

☺ *Never the Twain Shall Meet* ☺

Transliteration: representing the source script letters (graphemes *not* phonemes) with the characters of another script. محمد > /mHmd/

Transcription: representing the source script of a language in the target script in a manner that reflects pronunciation. This includes:

1. Phonetic transcription represents the actual speechsounds.

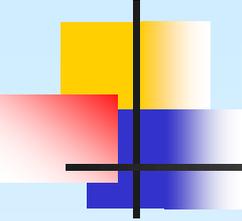
محمد > [muħëimmëd] (IPA)

2. Phonemic transcription represents the phonemes of the source.

محمد > /muHammad/

3. Popular transcription roughly represents pronunciation.

محمد > Mohammed, Muhammad, Moohammad, Moohamad...+200



NANA Conversion Modes

- Latin Clinton → کلینتون
- Japanese 埼玉(Saitama) → سائیتاما.
- Chinese 杨海洋 (Yang Haiyang) → هاییانغ یانغ
- Korean 부산 (Busan) → بوسان.

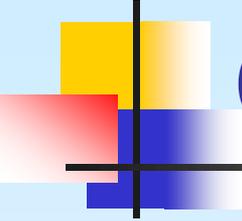
Diphthong Ambiguity for

福井 /fu-ku-i/

No.	Arabic	Google hits	Buckwalter
1	ف و ك و ئ ي	468	fwkw}y
2	ف و ك و ئ	9	fwkw}
3	ف و ك و ي	1950	Fwkwy
4	ف و ك و ي ي	335	Fwkwy y

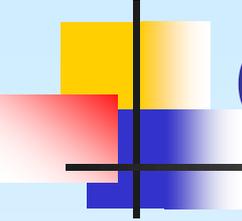
Long and Short Vowels

No.	Kanji	Kana	Phonemic	Arab1	Arab2	Arab3
1	太田	おおた	oota	أوتا		
2	風馬	ふうま	fuuma	فوما		
3	敬子	けいこ	keiko	ك ي ك و ك		
4	空野	くうの	kuuno	كوزو		
5	久野	くの	kuno	كوزو		
6	日枝	ひえだ	hieda	ه ي د ا	ه ي ئ د ا	ه ي ئ د ا
7	芳江	よしえ	yoshie	وش ي ي ي	وش ي ئ ي	وش ي ئ ي ي



Orthographic Ambiguity

- Short vowel omission (كاتب / kAtb /)
- Short vowel Representation (جامعة /jaami` a/)
- Multiple long /aa/
 - '*alif Tawiila* (ا) سوريا
 - '*alif maduuda* (آ) آسيا
 - '*alif maqSuura* (ى) آسيا الوسطى ا.
- Long vowel omission (هدا /haadha/)
- Long /aa/ ambiguity (شكرا, انا) (/an/, /a/)



Orthographic Ambiguity

- Otiose alif is silent كتبوا > /katabuu/
- Omission of *shadda* (مُحَمَّد < محمد) /Muhammad/
- Omission of *tanwiin* شكراً ¥\$ukrAF¥ (شُكْرًا)
- Complex hamza rules (فوكوأوكا vs فوكوؤوكا)
- Hamza omission (سائيتاما < سايتاما)
- Phonological alternation (لرجل الطويل)
'alrajulu alTawiilu/> /'arrajulu-Ttawiilu >
- Shortening long vowels (القاهرة في)
'fii-lqaahira/ > /fi-lqaahira/

Output from ARAN modules

Unvocalized (input)	Vocalized (ADAN)	Phonemic (ATAN)	Graphemic (AXAN)	Phonetic (APAN)	Popular (AVAN)*
محمد	مُحَمَّدٌ	muHammad	mHmd	muħɛimmɛd	Muhammad
قَابُوس	قَابُوسٌ	qaabuus	qAbws	qɑ:bu:s	Qaboos
جمال	جَمَالٌ	jamaal	jmAl	dʒɛimɛ:l	Jamal
مَكَّة	مَكَّةٌ	makka	mkp	mɛkkɛ	Mecca

*Only one popular variant is shown, but in reality there could be dozens. For example, for قَابُوس AVAN generates Qabuus, Qabus, Qabous, Qabooss, ... and many more.

ARAN Processing of قابوس /qAbws/

Conversion process	ARAN module	Input	Output	Remarks
Phonemic Transcription	ATAN	ق ا ب و س	/qaabuus/	linguistic representation of phonemes
English Transcription	ATAN	ق ا ب و س	Qaboos	Standard English spelling
Popular Transcriptions	ATAN	ق ا ب و س	Qabuus, Qabus, Qabous, Qabooss, Qaaboos, Kaboos, Kabuus, Gabous...	some of the many popular variants

ARAN Processing of قابوس /qAbws/

Conversion process	ARAN module	Input	Output	Remarks
Phonetic Transcription	APAN	قَابُوس	[qɑːbuːs]	scientific transcription in IPA
Unvocalized Transliteration	AXAN	قَابُوس	/qAbws/	Buckwalter transliteration of unvocalized Arabic
Vocalized Transliteration	AXAN	قَابُوس	/qaAbuws/	Buckwalter transliteration of vocalized Arabic
Diacrticization	ADAN	قَابُوس	قَابُوس	adding vowels (vocalization) and diacrtics to unvocalized Arabic
Arabization	NANA	Qabuus, Qabus, etc.	قَابُوس	converting non-Arabic to Arabic script

Major Arabic Romanization Systems

Example: شولوخ

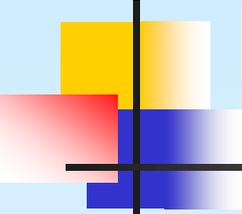
System	Example	Description
ALC-LC	shwllwkh	Romanization standard of the American Library Association-Library of Congress.
DIN	šūlūḥ	This refers to DIN 31635, the DIN standard for Arabic transliteration.
IPA	ʃuːluːx	International Phonetic Alphabet, a scientific system of uniquely and accurately representing speech sounds.
English	Shoulokh	One of many (at least 10) possible popular transcriptions.
Buckwalter	\$wllwx	A strict transliteration system widely used in information processing.

Sample output from ADAN module

Unvocalized	Vocalized	Transcription	English
محمد	مُحَمَّدٌ	muHammad	Muhammad
إِبْرَاهِيمَ	إِبْرَاهِيمَ - يَم	'ibrahiim	Abraham
إِسْحَاقَ	إِسْحَاقَ	isHaaq	Isaac
الرِّيَّاضَ	الرِّيَّاضَ	arriyaaD	Riyadh
مَكَّةَ	مَكَّةَ	makkah	Mecca
القَاهِرَةَ	القَاهِرَةَ	alqaahirah	Cairo

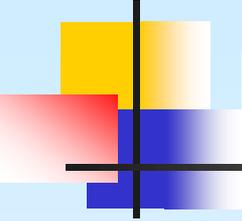
Variation in Arabic Names

Standard	Buckwalter	English	Variant	Error	Remarks
أبو ظبي	>bw Zby	Abu Dhabi	ظ بجاو و	ظ بجاو و	V: omit hamza E: alif maqsura replaces yaa'
سكندرية	Al<skndryp	Alexandria	سكندرية	سكندرية	V: omit hamza E: haa' replaces taa' marbuuTa
جدة	jdp	Jeddah	جدة	جده	V: explicit shadda E: haa' replaces taa' marbuuTa
الأردن	Al>rdn	Jordan	الاردن		V: omit hamza
بالو ألتو	bAlw >ltw	Palo Alto	ال تو ال و أ ال تو ال و		V1: omit hamza V2: madda replaces hamza
الرياض	AlryAD	Riyadh	الرياض		V: explicit shadda
طوكيو	Twkyw	Tokyo		توكيو	E: taa' replaces Taa'



Popular Transcriptions

Arabic	Buckwalter Transliteration	Popular Transcription
م عمر	mEmr	Moammar
م عمر	mEmr	Muammar
م عمر	mEmr	Mu'ammarr
م عمر	mEmr	Mu`ammarr
م عمر	mEmr	Mo'ammarr
م عمر	mEmr	Moammarr
م عمر	mEmr	Moamer
م عمر	mEmr	Moamar
م عمر	mEmr	Mohamar



MSA Flavors

Arabic	English	Phonemic	Phonetic Gulf	Phonetic Egyptian	Phonetic Levantine
قَابُوس	Qaboos	/qaabuus/	[qa:bu:s]	[ʔa:bu:s]	[qa:bu:s]
جمال	Jamal	/jamaal/	[dzëmɛ:l]	[gëmɛ:l]	[ʒɛmɛ:l]

Variation in Arabic Names

Standard	Buckwalter	English	Variant	Error	Remarks
ظ بي أبو و	>bw Zby	Abu Dhabi	ظ بي ابو و	ظ بي أبو و ظ بي ابو و	V: omit hamza E: 'alif maqsura replaces <u>yaa</u> '
سدك ندرية	Al<skndryp	Alexandria	سدك ندرية	سدك ندرية	V: omit hamza E: haa' replaces taa' marbuuTa
أل تو ب ال و	bAlw >ltw	Palo Alto	ال تو ب ال و آل تو ب ال و		V1: omit hamza V2: madda replaces hamza
طوك يو	Twkyw	Tokyo		توك يو	E: taa' replaces Taa'

English-to-Arabic Errors

NUM	English	Error	Google_Hits	Google_Arabic	CJKI_Arabic
1	Muhammad	NO	24,700,000	ﻣﻮﺣﺎﻣﺪ	ﻣﻮﺣﺎﻣﺪ
2	Moohammad	YES	2,720	ﻣﻮﺣﺎﻣﺪ	ﻣﻮﺣﺎﻣﺪ
3	Moohamad	YES	975	ﻣﻮﺣﺎﻣﺪ	ﻣﻮﺣﺎﻣﺪ
4	Mohammad	NO	15,100,000	ﻣﻮﺣﺎﻣﺪ	ﻣﻮﺣﺎﻣﺪ
5	Mohamad	NO	4,030,000	ﻣﻮﺣﺎﻣﺪ	ﻣﻮﺣﺎﻣﺪ
6	Muhamad	NO	795,000	ﻣﻮﺣﺎﻣﺪ	ﻣﻮﺣﺎﻣﺪ
7	Mohamed	NO	23,700,000	ﻣﻮﺣﺎﻣﺪ	ﻣﻮﺣﺎﻣﺪ
8	Mohammed	NO	26,700,000	ﻣﻮﺣﺎﻣﺪ	ﻣﻮﺣﺎﻣﺪ
9	Mohemmed	YES	940	ﻣﻮﺣﺎﻣﺪ	ﻣﻮﺣﺎﻣﺪ
10	Muhemmed	YES	12,800	ﻣﻮﺣﺎﻣﺪ	ﻣﻮﺣﺎﻣﺪ
11	Muhamed	NO	742,000	ﻣﻮﺣﺎﻣﺪ	ﻣﻮﺣﺎﻣﺪ
12	Muhammed	NO	7,800,000	ﻣﻮﺣﺎﻣﺪ	ﻣﻮﺣﺎﻣﺪ
13	Moohammed	YES	1,840	ﻣﻮﺣﺎﻣﺪ	ﻣﻮﺣﺎﻣﺪ
14	Mouhammed	NO	102,000	ﻣﻮﺣﺎﻣﺪ	ﻣﻮﺣﺎﻣﺪ

Variants and Errors

Rank	Type	ARABIC	Buckwalter	Googlits	Remarks
1	N	سدك ندرية	Askndryp	2,930,000	Normalized, no hamza
2	S	سدك ندرية	Askndryp	690,000	Standard form, with hamza
3	E	سدك ندرية	Askndryh	89,200	No hamza, taa marbuta replaced by haa
4	V	سدك ندرية	Askndry~p	954	Explicit shadda
5	E	سدك ندرية	Askndryh	897	taa marbuta replaced by haa
6	V	سدك ندرية	Askndry~p	245	no hamza, shadda explicit
7	E	سدك ندرية	AskndryA	80	hamza omitted, taa marbuuta replaced by alif
8	V	سدك ندرية	Asokanodary~ap	24	fully vocalized
9	E	سدك ندرية	Askndry~h	12	no hamza, shadda explicit, taa marbuta replaced by haa
10	E	سدك ندرية	AskndryA	7	taa marbuta replaced by alif tawiila
11	E	سدك ندرية	Askndry~h	5	taa marbuuta replaced by haa, shadda explicit