

### Exploring Arabic Script Behavior Multi-script Database Series #6

Prepared by: Frank Oberle (فرينك بيرل) and J. N. Hummel (🎵🎵🎵🎵🎵🎵)

Some examples of Arabic, Devanagari, Hebrew, Korean and Thai Scripts were introduced in Design Note #2 (Exploring Complex Text Layout), but Arabic Script is a bit more complex. While other Scripts consist of characters that may swap positions with their neighbors (e.g. Devanagari), be written right-to-left (e.g. Hebrew), be grouped into blocks (e.g. Korean), or placed above or below their neighbors (e.g. Thai), characters nevertheless remain identifiable as those that were entered.

Arabic Characters, on the other hand, often seem to change as each subsequent character is typed. This can make it difficult for many developers – who may be unfamiliar with such languages – to determine whether the data for which they are responsible is being stored, retrieved, and displayed correctly.

This more detailed look at Arabic Script will help database and application developers make sense of what they are seeing during development and testing of systems required to support any of the more than thirty languages written in Arabic Script, and serve as a basic reference.

Revised for public distribution: 7 September 2018

See page 20 for information on other material from Antikythera Publications.



Copyright © 2015, 2018 by the Authors

Permission is granted to distribute unaltered copies of this document, so long as this is not done for recompense or similar commercial purposes.



[www.AntikytheraPubs.com](http://www.AntikytheraPubs.com)

## Database Design Note Series on Multi-Language/Multi-Script Databases

1. Exploring Alphabets
2. Exploring Complex Text Layout
3. Exploring UTF-8
4. Evaluating Fonts for use in Multi-Lingual Documents
5. Evaluating Bidirectional Text Entry
6. **Exploring Arabic Script Behavior**

# Database Design Note Series – Exploring Arabic Script Behavior

## INTRODUCTION

Arabic Script has a strong tradition of decorative and often quite elaborate calligraphy. As a result, its individual characters are often joined (as with English cursive) to prioritize the aesthetics of whole words over individual characters. Based on their position within a word – Initial, Medial (between two others), or Final – characters may take on shapes that differ from their “normal” (Isolate) forms.<sup>1</sup> While current font and rendering technologies are not yet supported widely enough in typical business or personal computing environments to fully reproduce Arabic Script<sup>2</sup>, technology practitioners should have a basic understanding of Arabic Script behavior in order to deal with data in the variety of disparate languages that utilize it.

Arabic alphabets are used to write Arwi (in Sri Lanka and Southern India), Azerbaijani (in Iran), Balochi (بلوچی), Balti, Belarusian (among the Tatars there), Bosniaks (in Bosnia), Brahui (in Pakistan), Central Kurdish (in Iraq and Iran), Chinese (in some areas), Dari (in Pakistan), Fulfulde-Pulaar/Fulani (in Senegal), Hausa, Kashmiri, Kazakh, Kyrgyz (in Central Asia), Luri, some Mandinka dialects (in West Africa), Malay, Mozarabic, Ottoman Turkish, Pashto/Pukhto (پښتو in Afghanistan and Pakistan), Persian/ Farsi (فارسی in Iran), Punjabi (پنجابی in Pakistan – ਧੰਨਾਬੀ in India), Sindhi, Swahili (in East Africa), Urdu (اردو, in Pakistan and India), Uyghur (in China and Central Asia), and Uzbek (in Central Asia) – this is an incomplete list, and some of these are written in multiple Scripts.

For practical reasons, however, this overview will be limited to Modern Standard Arabic, Farsi/Persian, and Urdu. These widely used, yet distinct and unrelated languages<sup>3</sup> should provide a wide enough sampling to demonstrate the basic concepts behind Arabic Script behavior and, more importantly, evaluate the capability of your infrastructure, data handling, and application portfolios to support any of the aforementioned languages.

The keyboard layouts provided in this reference are only a small sampling of those that may be encountered in practice, even for these three languages, but are the ones used for the Arabic typing examples in this series.

The table of glyph forms shows the Isolate form of each character, along with its Initial, Medial and Final forms.

---

## Contents

Arabic (Modern Standard) Keyboard Layout.....	4
Persian (Farsi) Keyboard Layout.....	5
Urdu (Pakistani, CRULP) Keyboard Layout.....	6
U.S. English Key Assignments/Codes for Arabic IME Keyboard Layouts.....	7
Glyph forms for different Character Positions in Arabic Script.....	9
Character Alteration when Joining Adjacent Arabic Characters.....	18
References for Further Exploration.....	19

\*

- 
- 1 It is very important to note however that *only the Isolate forms are stored on disc or transmitted* – usually in UTF-8 form – and it is only during the presentation of these characters on paper or screen that the positional variants are used. Following the preliminaries, a detailed example of how variations of one character are handled is presented on page 18.
  - 2 Due to long standing issues in LibreOffice, for instance, the capability to rotate right-to-left text has simply been removed in current (6.x) releases – and in spite of Unicode, UTF-8, and OpenType, some applications have no support for right-to-left Scripts at all. See Limitations in Current Arabic Script Rendering on page 19 for commentary on even more subtle issues.
  - 3 Modern Standard Arabic is recognized across the Arabic-speaking world in much the same manner as BBC English is understood by English speakers in Boston, New York, Baltimore, Mobile, and Dallas (to say nothing of Toronto, Canberra and Liverpool) whose speech characteristics (and even vocabulary) are quite different. Farsi (Persian) is an Indo-European Language, and closely related to Dari and Tajik, although the latter language is written using Cyrillic Script. Similarly, Urdu (written with Arabic Script) and Hindi (written with Devanagari Script) speakers have few problems talking to each other since these languages have similar origins, differing primarily in their written forms.



# تخطيط لوحة المفاتيح العربية

## Arabic (Modern Standard) Keyboard Layout

(pronunciation: takhtit lawhat almafatih alearabia) (keystrokes: jo'd' g,pm hglthjdp hguvfdm)



Unicode Arabic Script Planes: 0x0600-0x06ff, 0x0750-0x077f 0x08a0-0x08ff 0xfb50-0xfdff 0xfe70-0xfeff FreeSerif-11 ←RTL

Similar Glyphs / Letter Forms in the Right-to-Left Arabic Abjad

Numeric Characters are laid out in left-to-right order.

ي ئ ي ي	ث ت ب	خ ج ح	ض ص ش س	ظ ط	ع غ	ن ل ك ق ف	ة ه	لآ لا
0649 0620 0626 064a	0628 062a 062b	062d 062c 062e	0633 0634 0635 0636	0637 0638	0639 063a	0641 0642 0643 0644 0646	0647 0629	fefb fef5 fef7
n z d	f j e	p [ o	s a w q	' /	u y	t r ; g k	i m	b B G

Note the paired delimiter reversals on the (, ), D, F, C, and V keys.

(pronunciation: cheedamon chataykedee farsee) (keystrokes: ]dnlhk wtpi ;gdn thvsd)



Unicode Arabic Script Planes: 0x0600-0x06ff, 0x0750-0x077f 0x08a0-0x08ff 0xfb50-0xfdff 0xfe70-0xfeff FreeSerif-12 ←RTL

Similar Glyphs / Letter Forms in the Right-to-Left Persian/Farsi/Arabic Abjad

Numeric Characters are laid out in left-to-right order.

ی ئ ی ی	ث ت ب	خ ج ح	ض ص ش س	ظ ط	غ ع	ن ل ك ق ف	ه	لا لا
0649 0620 0626 064a	0628 062a 062b	062d 062c 062e	0633 0634 0635 0636	0637 0638	0639 063a	0641 0642 0643 0644 0646	0647 0629	fefb fef5 fef7
S	f j e	p [ o	s a w q	x z	u y	t r Z g k	i J	

Note the paired delimiter reversals on the **[ (, ) ]**, **[ {, } ]**, **[ (P) ]**, and **[ (O) ]** keys; Farsi uses French style Guillemets (« and », located on the K and L keys) in place of parentheses. What looks like a stylized comma located on the English “comma key” is actually a an actual alphabetic character Also note that the question mark is reversed.

# اردو زبان کی بورڈ ترتیب

## Urdu (Pakistani, CRULP) Keyboard Layout

(pronunciation: oardoo zabawn see ?? hekseem) (keystrokes: ardw zban ki brwD trtib)



Unicode Arabic Script Planes: 0x0600-0x06ff, 0x0750-0x077f 0x08a0-0x08ff 0xfb50-0xfdff 0xfe70-0xfeff FreeSerif-11 ←RTL

Similar Glyphs / Letter Forms in the Right-to-Left Urdu Arabic Abjad

Numeric Characters are laid out in left-to-right order.

ی ی ی ی	ت ت ب	ح ج ح	ض ص ش س	ظ ط	غ ع	ن ل ك ق ف	ہ	لا لا لا
0649 0620 0626 064a	0628 062a 062b	062d 062c 062e	0633 0634 0635 0636	0637 0638	0639 063a	0641 0642 0643 0644 0646	0647 0629	fefb fef5 fef7
U	b t C	h j K	s x S J	y V	e G	f q	l n	

The Urdu Language used widely in India and Pakistan is Indic, but written right-to-left using Arabic Script. This phonetic keyboard layout is that defined by the Center for Research in Urdu Language Processing and seems to be the standard layout now used throughout Pakistani government, business, and academic institutions. As Sahar Afshar<sup>4</sup> notes, the Arabic Script used to write Urdu is often spaced more widely than it is when used for other Languages, but this stylistic difference is not implemented by the font itself.

Note the paired delimiter reversals on the `{`, `}`, `[`, and `]` keys; Urdu doesn't use "regular/Western" parentheses. Also note that the question mark is reversed.

Actual Arabic forms of the numbers (١٢٣٤٥٦٧٨٩٠) are given precedence over the "Western" forms but, as with numbers in other right-to-left scripts, are written left-to-right.

4 Several of her lectures and writings are included in References for Further Exploration on page 19.

## U.S. English Key Assignments/Codes for Arabic IME Keyboard Layouts

Standard Arabic (qwerty), Persian (Farsi/Iranian), and Urdu (Pakistan, CRULP<sup>5</sup>)

Where the English key presses result in the same Unicode Character in any row, those identical characters are shaded in yellow.

KEY	Arabic Standard	Persian/Farsi	Urdu/Pakistani
A	ا u+0650	ؤ u+0624	آ u+0622
B	لا u+feef5	ZWNJ u+200c <sup>6</sup>	. u+002e
C	} u+007d	ژ u+0698	ث u+062b
D	] u+005d	ي u+064a	ڈ u+0688
E	ة u+064f	ۀ u+064d	ۀ u+0670
F	[ u+005b	إ u+0625	ۀ u+0651**
G	أ u+feef7	أ u+0623	غ u+063a
H	أ u+0623	آ u+0622	ه u+06be
I	÷ u+00f7	ۀ u+0651	ا u+0650
J	- u+0640	ة u+0629	ض u+0636
K	، u+060c <sup>4</sup>	» u+00bb	خ u+062e
L	/ u+002f	« u+00ab	ۀ u+0654
M	' u+0027	ء u+0621	ۀ u+0658
N	آ u+0622	ۀ u+0654	ن u+06ba
O	x u+00d7	] u+005d	ة u+06c3
P	؛ u+061b	[ u+005b	ة u+064f
Q	ة u+064e	ۀ u+0652	ۀ u+0652
R	ة u+064c	ۀ u+064b	ر u+0691
S	ة u+064d	ئ u+0626	س u+0635
T	لا u+feef9	ة u+064f	ط u+0679
U	` u+0060	ة u+064e	ئ u+0626
V	{ u+007b	ۀ u+0670	ظ u+0638
W	ۀ u+064b	ۀ u+064c	ۀ u+0651**
X	ۀ u+0652	ۀ u+0653	ژ u+0698
Y	إ u+0625	ا u+0650	ۀ u+064e
Z	~ u+007e	ك u+0643	ذ u+0630

Key	Arabic Standard	Persian/Farsi	Urdu/Pakistani
a	ش u+0634	ش u+0634	ا u+0627
b	لا u+feefb	ذ u+0630	ب u+0628
c	ؤ u+0624	ز u+0632	چ u+0686
d	ي u+064a	ی u+06cc	د u+062f
e	ث u+062b	ث u+062b	ع u+0639
f	ب u+0628	ب u+0628	ف u+0641
g	ل u+0644	ل u+0644	گ u+06af
h	ا u+0627	ا u+0627	ح u+062d
i	ه u+0647	ه u+0647	ی u+06cc
j	تا u+062a	تا u+062a	ج u+062c
k	ن u+0646	ن u+0646	ک u+06a9
l	م u+0645	م u+0645	ل u+0644
m	ة u+0629	پ u+067e	م u+0645
n	ی u+0649	د u+062f	ن u+0646
o	خ u+062e	خ u+062e	، u+06c1
p	ح u+062d	ح u+062d	پ u+067e
q	ض u+0636	ض u+0636	ق u+0642
r	ق u+0642	ق u+0642	ر u+0631
s	س u+0633	س u+0633	س u+0633
t	فا u+0641	فا u+0641	تا u+062a
u	ع u+0639	ع u+0639	ء u+0621
v	ر u+0631	ر u+0631	ط u+0637
w	س u+0635	س u+0635	و u+0648
x	ء u+0621	ط u+0637	ش u+0634
y	غ u+063a	غ u+063a	ع u+06d2
z	ئ u+0626	ظ u+0638	ز u+0632

5 Center for Research in Urdu Language Processing

6 Zero Width Non Joiner Mark; overrides automatic character/glyph joining behavior.

KEY	Arabic Standard	Persian/Farsi	Urdu/Pakistani
◌	◌ u+0651	÷ u+00f7	◌ u+064b
!	! u+0021	! u+0021	1 u+0031
@	@ u+0040	’ u+066c <sup>8</sup>	2 u+0032
#	# u+0023	, u+066b	3 u+0033
\$	\$ u+0024	ريال u+fdfc <sup>9</sup>	4 u+0034
%	% u+0025	% u+066a	5 u+0035
^	^ u+005e	× u+00d7	6 u+0036
&	& u+0026	، u+060c	7 u+0037
*	* u+002a	* u+002a	8 u+0038
)	) u+0029	) u+0029	9 u+0039
(	( u+0028	( u+0028	0 u+0030
_	_ u+005f	- u+0640	- u+005f
+	+ u+002b	+ u+002b	+ u+002b
<	< u+003c	} u+007d	} u+007d
>	> u+003e	{ u+007b	{ u+007b
	u+007c	u+007c	u+007c
:	: u+003a	: u+003a	: u+003a
"	" u+0022	؛ u+061b	" u+0022
,	, u+002c	> u+003e	< u+003c
.	. u+002e	< u+003c	، u+066b <sup>11</sup>
?	? u+003f	? u+061f	? u+061f
	ريال u+fdfc +others	ريال u+fdfc	Rs u+20a8
B		ZWNJ u+200c <sup>2</sup>	

Key	Arabic Standard	Persian/Farsi	Urdu/Pakistani
◌	◌ u+0630	ZWJ u+200d <sup>7</sup>	◌ u+007e
1	1 u+0031	۱ u+06f1	۱ u+06f1
2	2 u+0032	۲ u+06f2	۲ u+06f2
3	3 u+0033	۳ u+06f3	۳ u+06f3
4	4 u+0034	۴ u+06f4	۴ u+06f4
5	5 u+0035	۵ u+06f5	۵ u+06f5
6	6 u+0036	۶ u+06f6	۶ u+06f6
7	7 u+0037	۷ u+06f7	۷ u+06f7
8	8 u+0038	۸ u+06f8	۸ u+06f8
9	9 u+0039	۹ u+06f9	۹ u+06f9
0	0 u+0030	۰ u+06f0	۰ u+06f0
-	- u+002d	- u+002d	- u+002d
=	= u+003d	= u+003d	= u+003d
ج	ج u+062c	ج u+062c	] u+005d
د	د u+062f	د u+0686	[ u+005b
\	\ u+005c	\ u+005c	\ u+005c
ذ	ذ u+0643	ذ u+06a9	؛ u+061b
ط	ط u+0637	ط u+06af	' u+0027
و	و u+0648	و u+0648	، u+060c <sup>10</sup>
ز	ز u+0632	. u+002e	- u+06d4
ظ	ظ u+0638	/ u+002f	/ u+002f
	◌ u+20e3 <sup>12</sup>		
		ZWJ u+200d <sup>3</sup>	

### Use of zero-width joiner characters in Persian/Farsi (using Latin script example)

ⓕ Ⓛ Ⓨ fly or fly? – latter contains u+200c (zwnj) between “f” and “l” to override automatic ligature generation.

7 Zero Width Joiner Mark

8 Arabic Thousands Separator (used with Persian/Farsi)

9 Arabic ligature: Raa’ (ر u+0631) with Farsi Yaa’ (ي u+06cc), Alif (ا u+0627), and Laam (ل u+0644). The word “Riyal”.

10 Arabic Comma (used with Standard Arabic and Urdu)

11 Arabic Decimal Point (used with Urdu)




12 Post-fix Key Cap character used in this table.



## Glyph forms for different Character Positions in Arabic Script

Many characters in the Arabic Script, used not only for those Languages illustrated here, but a wide variety of others, may take differing forms depending on their position within a word (i.e. “contextual letter forms”). Arabic script is cursive in nature, and this permits the characters to be connected in a way that mimics – albeit incompletely – traditional Arabic calligraphic practices. The following charts will assist in identifying the characteristics of these forms.

**Legend for Arabic Script character form descriptions shown in the charts on the following pages.**

Example for the Baa' (ب) character:	Baa'	⇐ Arabic Character Name <sup>13</sup> (transliterated)
Initial Form (ب) of the Baa' (ب) character ⇒ (glyph truncated when connecting to a following character)	 u+0628 B U+D8A8	⇐ Unicode Hex Value for the Baa' (ب) character in its “normal” Isolate Form, the common English transliteration <sup>14</sup> and its UTF-8 Encoding, also in Hex.
Unicode Hex Value for this Form ⇒	u+fe91	⇐ If this cell has a full yellow background, it is a “Sun” letter/character. <sup>15</sup>
Medial Form (ب) of the Baa' (ب) character ⇒ (arm added when connecting between two adjacent characters)	 u+fe92	⇐ Isolate form of this character as it might appear on a keyboard (actual glyph will depend on the keyboard design) and stored on disk and in memory, or transmitted.
Unicode Hex Value for this Form ⇒	u+fe90	⇐ Key press to type this character when using the included IME keyboard layouts for standard Arabic, Farsi (Persian), or Urdu (Pakistani) respectively (see pages 4, 5, and 6).
Final Form (ب) of the Baa' (ب) character ⇒ (arm added to Isolate form when used as a final character)	 Ar: [ف] Fa: [ف] Ur: [ب]	⇐ Key press to type this character when using the included IME keyboard layouts for standard Arabic, Farsi (Persian), or Urdu (Pakistani) respectively (see pages 4, 5, and 6).
Unicode Hex Value for this Form ⇒	u+fe90	⇐ Pronunciation w/English example
See detailed example of Form usage on page 18.	B Bird	⇐

A user only “types” the normal form of the character, i.e. that which is shown on the top of each physical key cap or virtual keyboard map. Because of the multiple forms a character may take depending on its position within a word, displayed text may change dynamically as each additional character is typed. Not all characters have all forms. Actual “alphabetic characters”<sup>16</sup> are indicated by a double border around the cell containing the Isolate form of the character. The appearance of the characters and glyphs shown in this document may obviously differ depending on the font in use.

In all forms, the number of dot markings above or below the character (in this case “one below”) will remain constant, providing a clue to its “real” identity.

13 The names of the characters may differ in Farsi/Persian and Urdu since they are different languages; the Modern Standard Arabic name is used here for convenience.

14 Note that a character with an underscore represents a different sound. See the characters ح (u+062d) and ه (u+0647), for example, pronounced as H and H respectively.

15 Other characters are standard “Moon” characters. The distinction simply – although not *entirely* correctly – means that pronunciation of the definite article “Al” used with nouns is replaced with something more akin to an elision if it is used with a word beginning with a “Moon” character. “The book” (الكتاب) for example is pronounced “al Kitab”, while “the sun” (الشمس) is pronounced “as shum” rather than “al shum.” Note that both articles are written as ال; only the pronunciation differs.

16 The Arabic “alphabets” for most languages are actually Abjads (having only consonants), and “impure Abjads” at that, but such distinctions aren’t important here.

	Yaa' w/Hamza above	'Alif w/Hamza below	Waaw w/Hamza above	'Alif w/Hamza above	'Alif w/Madda above	Cutting Hamza	← Character Name
Initial Form	U+0626 U+D8A6	U+0625 U+D8A5	U+0624 U+D8A4	U+0623 U+D8A3	U+0622 U+D8A2	U+0621 U+D8A1	← Unicode HexCode
Medial Form	ئ	ا	و	آ	آ	ء	← Transliteration
Final Form	Ar: [Z] Fa: [S] Ur: [U]	Ar: [Y] Fa: [F] Ur: —	Ar: [C] Fa: [A] Ur: —	Ar: [H] Fa: [G] Ur: —	Ar: [N] Fa: [H] Ur: [A]	Ar: [X] Fa: [M] Ur: [U]	← UTF-8 Encoding
	Yaa' is u+06cc	'Alif is u+0627	Waaw is u+0648	'Alif is u+0627	Double 'Alif (u+0627)	Glottal Stop <sup>17</sup>	<b>BEGIN</b> ← This See final row for Arabic-specific punctuation characters. ← Arabic Keyboard ← Farsi Keyboard ← Urdu Keyboard ← Pronunciation

	Haa' (cf u+0647)	Jeem (Geem) (Zheem)	Thaa' (Theh')	Taa' (Teh') Maftuhah	Taa' (Teh') Marbuta	Baa' or Beh'	'Alif
Initial Form	ح u+fea3	ج u+fe9f	ث u+fe9b	ت u+fe97	ة u+D8A9	ب u+fe91	ا u+D8A7
Medial Form	ح u+fea4	ج u+fea0	ث u+fe9c	ت u+fe98	ة u+D8A9	ب u+fe92	ا u+D8A7
Final Form	ح u+fea2	ج u+fe9e	ث u+fe9a	ت u+fe96	ة u+fe94	ب u+fe90	ا u+fe8e
	Ar: [D] Fa: [D] Ur: [H]	Ar: [J] Fa: [J] Ur: [J]	Ar: [E] Fa: [E] Ur: [C]	Ar: [T] Fa: [T] Ur: [T]	Ar: [M] Fa: [J] Ur: —	Ar: [F] Fa: [F] Ur: [B]	Ar: [H] Fa: [H] Ur: [A]
	H Hit <sup>18</sup>	ZH meaSure	THing <sup>19</sup>	T Tap <sup>20</sup>		B Bird	A Apple (cf u+0649)

17 The cutting Hamza character ء (u+0621) acts as a consonant, but isn't considered an actual alphabetic character. The Cockney pronunciation of t' in bot'tle is similar.

18 The characters ح (u+062d) and ه (u+0647), both forms of the "h" sound, are pronounced as in the emphatic h in **hit** and the softer h in **he** respectively.

19 The characters ث (u+062b), ذ (u+0630), and ظ (u+0638), all forms of the "th" sound, are pronounced as in **thing**, **that**, and **those** respectively.

20 The characters ت (u+062a) and ط (u+0637), both forms of the "t" sound, are pronounced as in **tap** and **taught** (more emphatic) respectively.

	Sheen (Shin)	Seen (Sin)	Zaa' (Zay')	Raa' (Reh')	Dhaal (Zaal)	Daal	Khaa'
Initial Form	ش u+feb7 U+0634 SH U+D8B4	س u+feb3 U+0633 S U+D8B3	ز u+feb0 U+0632 Z U+D8B2	ر u+feae U+0631 R U+D8B1	ذ u+feac U+0630 DH U+D8B0	د u+feaa U+062F D U+D8AF	خ u+fea7 U+062E KH U+D8AE
Medial Form	ش u+feb8	س u+feb4	ز u+feb0	ر u+feae	ذ u+feac	د u+feaa	خ u+fea8
Final Form	ش u+feb6 Ar: [a] Fa: [a] Ur: [x]	س u+feb2 Ar: [s] Fa: [s] Ur: [s]	ز u+feb0 Ar: [z] Fa: [z] Ur: [z]	ر u+feae Ar: [r] Fa: [r] Ur: [r]	ذ u+feac Ar: [d] Fa: [b] Ur: [z]	د u+feaa Ar: [d] Fa: [n] Ur: [d]	خ u+fea6 Ar: [k] Fa: [k] Ur: [k]
	SH <sub>ee</sub> p	S <sub>i</sub> p	Z <sub>oo</sub>	rolled R	TH <sub>at</sub> <sup>19</sup>	Di <sub>d</sub> <sup>21</sup>	J.S. BaCH (b.1685)

		Ghayn (Ghain)	Ayn (Ayin)	Dhaa' (Zah')	Taa' (Tah')	Daad (Dad)	Saad (Sad)
Initial Form	U+063B TO U+063F	غ u+fecf U+063A GH U+D8BA	ع u+fecb U+0639 c U+D9B9	ظ u+fec7 U+0638 DH U+D8B8	ط u+fec3 U+0637 T U+D8B7	ض u+fecb U+0636 D U+D8B6	ص u+febb U+0635 S U+D8B5
Medial Form	Glyphs from u+63b to u+63f not shown as they are only used for ancient Persian and Azerbaijani	غ u+fed0	ع u+fec0	ظ u+fec8	ط u+fec4	ض u+fec0	ص u+fec0
Final Form		غ u+fece Ar: [v] Fa: [v] Ur: [G]	ع u+fec0 Ar: [u] Fa: [u] Ur: [e]	ظ u+fec6 Ar: [z] Fa: [z] Ur: [v]	ط u+fec2 Ar: [t] Fa: [x] Ur: [v]	ض u+febe Ar: [d] Fa: [q] Ur: [J]	ص u+feba Ar: [w] Fa: [w] Ur: [S]
		French R, as in Rue	Glottal Stop, ≈ u <sub>g</sub> k	Th <sub>ose</sub> <sup>19</sup>	T <sub>ought</sub> ; cf t earlier	Da <sub>ughter</sub> <sup>21</sup>	S <sub>w</sub> ord (open S)

21 The characters د (Daal u+062f) and ض (Daad u+0636), both forms of the “d” sound, are pronounced like the “d” in **Did** and **Daughter** (a bit more forceful) respectively.

	Noon (Nun)	Meem (Mim)	Laam	Kaaf	Qaaf	Faa' (Feh')	Tatweel (Kashida)
Initial Form	ز u+fee7 U+0646 N U+D986	م u+fee3 U+0645 M U+D985	ل u+fedf U+0644 L U+D984	ك u+fedb 22 U+0643 K U+D983	ق u+fed7 U+0642 Q U+D982	ف u+fed3 U+0641 F U+D981	U+0640 U+D980
Medial Form	ن u+fee8	م u+fee4	ل u+fee0	ك u+fedc	ق u+fed8	ف u+fed4	-
Final Form	ن u+fee6 Ar: [k] Fa: [k] Ur: [n]	م u+fee2 Ar: [l] Fa: [l] Ur: [m]	ل u+fede Ar: [q] Fa: [q] Ur: [l]	ك u+feda Ar: [k] Fa: [k] Ur: —	ق u+fed6 Ar: [r] Fa: [r] Ur: [q]	ف u+fed2 Ar: [f] Fa: [f] Ur: [f]	Under-score ⇒ Ar: [ʔ] Fa: [ʔ] Ur: —
	<b>N</b> Need	<b>M</b> Many	<b>L</b> Lift (light L)	<b>K</b> King	<b>C</b> Caught (≈ gk)	<b>F</b> Fish	(glyph extension)

	Kasratan	Dammatan	Fathatan (double FatHa)	Yaa' (Yeh')	'Alif Maksura	Waaw	Haa' (Heh')
Initial Form	U+064D U+D98D	U+064C U+D98C	U+064B U+D98B	ي u+fef3 U+064A Y   EE U+D98A	ا u+fbe8 U+0649 A U+D989	و u+fbeb U+0648 W   OO U+D988	ه u+feeb U+0647 H U+D987
Medial Form	◌ِ	◌ِ	◌ِ	ي u+fef4	ا u+fbe9	و u+feec	ه u+feec
Final Form	Ar: [s] Fa: Ur:	Ar: [r] Fa: [w] Ur: —	Ar: [w] Fa: [r] Ur: [ˁ]	ي u+fef2 Ar: [d] Fa: [d] Ur: —	ا u+fef0 Ar: [n] Fa: Ur: —	و u+feee Ar: [w] Fa: [w] Ur: [w]	ه u+feea Ar: [h] Fa: [h] Ur:
	indicates genitive case		Single FatHa is u+064e	Yes / ChEEse	Also see u+06cc	Wonder / FOOD	h he; cf h note earlier

22 Note the significant difference between the Isolate form of the Kaaf character and its initial and medial forms.

23 The Laam ل takes on a special shape when followed by an Alif ا alone, with a Madda, or with upper or lower Hamzas. (See u+fef5, u+fef7, u+fef9, and u+fefb).

24 The 'Alif Maksura ( ا u+0649) character, used as a word ending 'Alif in certain languages, has all three contextual variants, but there are versions used in other languages which have none, e.g. the Farsi (also Urdu) Yaa' at code point u+06cc. Unicode code point u+feef ( ا ) is a similar glyph not shown in this document.

25 A FatHa ( َ u+064e) is placed over the Waaw ( و u+0648) character ( و ) when the latter is used as a consonant.








	Hamza above	'Alif Maddah (above)	Sukoon	Shadda	Vowel: Kasra	Vowel: Damma	Vowel: FatHa
Initial Form	U+0654 U+D994	U+0653 U+D993	U+0652 U+D992	U+0651 U+D991	U+0650 I <sup>26</sup> U+D990	U+064F U <sup>26</sup> U+D98F	U+064E A <sup>26</sup> U+D98E
Medial Form							
Final Form	Ar: — Fa: [N] Ur: [L]	Ar: — Fa: [X] Ur: — u+feef	Ar: [X] Fa: [Q] Ur: [Q]	Ar: [~] Fa: [I] Ur: [W] Ur: [E]	Ar: [A] Fa: [Y] Ur: [I]	Ar: [E] Fa: [T] Ur: [P]	Ar: [Q] Fa: [U] Ur: [Y]
		Doubles an 'Alif	Cancel implied vowel	Doubles a Consonant	Long I, e.g. Bit	Long U, e.g. PUt	Long A, e.g. PAt
















  

	Arabic digit One	Arabic digit Zero	Zwarakay	Noon Ghunna	Inverted Damma	Subscript 'Alif	Hamza below
Initial Form	U+0661 1 U+D9A1	U+0660 0 U+D9A0	U+0659 U+D999	U+0658 U+D998	U+0657 U+D997	U+0656 U+D996	U+0655 U+D995
Medial Form							
Final Form	Ar: — Fa: — Ur: —	Ar: — Fa: — Ur: —	Ar: — Fa: — Ur: —	Ar: — Fa: — Ur: [M]	Ar: — Fa: — Ur: —	Ar: — Fa: — Ur: —	Ar: — Fa: — Ur: —







26 Corresponding short vowel sounds are seldom written or printed; when the FatHa (َ), Damma (ُ), or Kasra (ِ) are present, they invariably represent “long” vowel sounds.







27 Two sets of Arabic-Indic digits shown in this list. The first runs from u+0660-0669, and the second from u+06f0-06f9. Most are similar, but the 6, 5, and 4 differ enough (٦ ٥ versus ٦ ٥ ٤) that a separate sequence, known as Eastern-Arabic-Indic, is used for Farsi and Urdu key mappings shown here. Some sources show the Urdu 4 and 6 differently, though. Note that, unlike Arabic text, all forms of Arabic digits are written from left-to-right when forming number sequences.

	Arabic digit Eight	Arabic digit Seven	Arabic digit Six	Arabic digit Five	Arabic digit Four	Arabic digit Three	Arabic digit Two
Initial Form	U+0668 8 U+D9A8	U+0667 7 U+D9A7	U+0666 6 U+D9A6	U+0665 5 U+D9A5	U+0664 4 U+D9A4	U+0663 3 U+D9A3	U+0662 2 U+D9A2
Medial Form	 27	 27	 27	 27	 27	 27	 27
Final Form	Ar: — Fa: — Ur: —	Ar: — Fa: — Ur: —	Ar: — Fa: — Ur: —	Ar: — Fa: — Ur: —	Ar: — Fa: — Ur: —	Ar: — Fa: — Ur: —	Ar: — Fa: — Ur: —








	Paa'	Tteh	Superscript 'Alif	Arabic Thousands Mark	Arabic Decimal Mark	Arabic Percent Sign	Arabic Digit 9
Initial Form	U+067E U+D9BE	U+0679 U+D9B9	U+0670 U+D9B0	U+066C U+D9AC	U+066B U+D9AB	U+066A U+D9AA	U+0669 9 U+D9A9
Medial Form				 28	 28	 28	 27
Final Form	Ar: — Fa:  Ur: 	Ar: — Fa: — Ur: 	Ar: — Fa:  Ur: 	Ar: — Fa: — Ur: —	Ar: — Fa:  Ur: 	Ar: — Fa:  Ur: —	Ar: — Fa: — Ur: —
	<b>P Peel</b>			Numeric Punctuation	Numeric Punctuation	Numeric Punctuation	








28 Some Arabic languages use “lower ASCII” versions (comma, period and % of these numeric punctuation characters depending on context).

		Gaf	Keheh	Jeh	Rreh	Ddal	Tcheem
Initial Form		U+06AF U+DAAF	U+06A9 U+DAA9	U+0698 U+DA98	U+0691 U+DA91	U+0688 U+DA88	U+0686 TCH U+DA86
Medial Form		u+fecf 					
Final Form		Ar: — Fa: [ ] Ur: [G]	Ar: — Fa: [ ] Ur: [K]	Ar: — Fa: [C] Ur: [X]	Ar: — Fa: — Ur: [R]	Ar: — Fa: — Ur: [D]	Ar: — Fa: [J] Ur: [C]
		u+fecf					Ch Chip

	Full Stop (period)	Yaa' Baree	Farsi Yaa'	Taa' Marbuta Goal	Heh Goal	Heh Doachashmee	Noon (Nun) Ghunna
Initial Form	U+06D4 U+DB94	U+06D2 U+DB92	U+06CC A U+DB8C	U+06C3 U+DB83	U+06C1 U+DB81	U+06BE U+DABE	U+06BA U+DABA
Medial Form	-			 <sup>29</sup>			
Final Form	Ar: — Fa: — Ur: [ ]	Ar: — Fa: — Ur: [V]	Ar: — Fa: [d] Ur: [i]	Ar: — Fa: — Ur: [O] u+fe94	Ar: — Fa: — Ur: [O]	Ar: — Fa: — Ur: [H]	Ar: — Fa: — Ur: [N]
	Don't confuse w/u+640		cf Footnote 24 (u+0649)				

29 The Taa' Marbuta ۞ (u+06c3) character used in Urdu is always preceded by a Fatah َ (u+064e) character. It is only found at the end of a word.

	Arabic extended digit 6	Arabic extended digit 5	Arabic extended digit 4	Arabic extended digit 3	Arabic extended digit 2	Arabic extended digit 1	Arabic extended digit 0
Initial Form	U+06F6 6 U+DBB6	U+06F5 5 U+DBB5	U+06F4 4 U+DBB4	U+06F3 3 U+DBB3	U+06F2 2 U+DBB2	U+06F1 1 U+DBB1	U+06F0 0 U+DBB0
Medial Form	 27 30	 27	 27 30	 27	 27	 27	 27
Final Form	Ar: — Fa: ⑥ Ur: ⑥	Ar: — Fa: ⑤ Ur: ⑤	Ar: — Fa: ④ Ur: ④	Ar: — Fa: ③ Ur: ③	Ar: — Fa: ② Ur: ②	Ar: — Fa: ① Ur: ①	Ar: — Fa: ① Ur: ①
	cf Digits note above	cf Digits note above	cf Digits note above	cf Digits note above	cf Digits note above	cf Digits note above	cf Digits note above

	Laam + 'Alif	Laam + Hamza below	Laam + Hamza above	Laam + Madda above	Arabic extended digit 9	Arabic extended digit 8	Arabic extended digit 7
Initial Form	U+FEFB EFBBBB	U+FEF9 EFBBB9	U+FEF7 EFBBB7	U+FEF5 EFBBB5 <sup>31</sup>	U+06F9 9 U+DBB9	U+06F8 8 U+DBB8	U+06F7 7 U+DBB7
Medial Form					 20	 20	 20
Final Form	Ar: ③ Fa: — Ur: —	Ar: ① Fa: — Ur: —	Ar: ② Fa: — Ur: —	Ar: ④ Fa: — Ur: —	Ar: — Fa: ⑨ Ur: ⑨	Ar: — Fa: ⑧ Ur: ⑧	Ar: — Fa: ⑦ Ur: ⑦
	Lay ③ + ④	Ligature	Ligature	Ligature	cf Digits note above	cf Digits note above	cf Digits note above

30 Some Urdu keyboards map the “4” key to a character I can’t identify in Unicode; some map the “6” to u+0666. Those shown above seem to be the most common.

31 Character values in the lesser-used higher Arabic Unicode Planes (0xfb50-0xfdf and 0xfe70-0xfeff) require three bytes in UTF-8 format.



		Sheen + 'Alif Maksura	Zero Width Joiner	Zero Width Non-Joiner	Arabic Question Mark	Arabic Semicolon	Arabic Comma
Initial Form		U+FDFC EFB7BC	u+200d zwj E2808D	U+200C zwnj E2808C	U+061F ? U+D89F	U+061B ; U+D89B	U+060C , U+D88C
Medial Form		ريال			؟	؛	،
Final Form		Ar: — Fa: ﷲ Ur: —	Ar: — Fa: ﷲ Ur: —	Ar: — Fa: ﷲ Ur: —	Ar: ﷲ Fa: ﷲ Ur: ﷲ	Ar: ﷲ Fa: ﷲ Ur: ﷲ	Ar: ﷲ Fa: ﷲ Ur: ﷲ
		Ligature for Riyal <sup>32</sup>	Forces character joins	Prevents character joins	Punctuation	Punctuation	Punctuation

	Post-Fix Key Cap						
Initial Form	U+20E3 □ E283A3		U+20A8 R\$ E282A8				
Medial Form	□		Rs				
Final Form	u+20e3 is a postfix KeyCap used in this document.		Ar: — Fa: — Ur: —				
			Pakistani Rupee				

<sup>32</sup> This slightly stylized ligature used to represent the currencies used in Iran, Oman, Qatar, Saudi Arabia, and Yemen, (all called Riyal and spelled ريال, u+631, u+6cc, u+644) only appears on the Farsi/Persian keyboard shown in this document. Many of these countries often use western abbreviations (e.g. SR for the Saudi Riyal) as well.



## References for Further Exploration

### Arabic Script and Related Topics

*Arabic Alphabet General Reference* <sup>33</sup>

[https://en.wikipedia.org/wiki/Arabic\\_alphabet](https://en.wikipedia.org/wiki/Arabic_alphabet)

*Overview of Arabic Calligraphic Styles – Sahar Afshar*

<https://www.rosettatype.com/blog/2016/05/24/Arabic-calligraphic-styles>

*World Currency Symbols*

<https://www.xe.com/symbols.php>

*Unicode Charts*

<http://unicode.org/charts/>

### University of Reading Masters Program: Recommended Reading Material

<http://typefacedesign.net/resources/preparation-for-incoming-matd-students/>

### Limitations in Current Arabic Script Rendering

*Unicode Arabic lacks concept of contextually neutral, “inline” characters – Thomas Milo*

<https://unicode.org/L2/L2014/14109-inline-chars.pdf>

*Considering the Old; Designing the New (Video 17:29)*

“methods of designing Arabic typefaces with a new voice by use of case studies. It will do so by analyzing innovative practice in various aspects of font development such as design, technology, and solutions for harmonization, etc.” – Sahar Afshar

<https://www.youtube.com/watch?v=SBCpp-vxMBs>

*Typo Labs 2018 Lecture Series (includes Making it Fit video below)*

<https://www.typonetwork.com/news/article/typo-labs-2018-how-far-can-we-go>

*On Expanding Connections (aka Making it Fit, from TypoLabs 2018) (Video 45:41)*

“historic context and use of the Arabic extension, ... how, by using variable font technology, the long standing issue of organic extensions can be resolved to make possible a more dynamic typesetting for the Arabic script.” – Sahar Afshar and José Miguel Solé Bruning

[http://showclipaz.com/at/typo-labs-2018-sahar-afshar-jose-miguel-s\\_eMYLtwQWjTW4](http://showclipaz.com/at/typo-labs-2018-sahar-afshar-jose-miguel-s_eMYLtwQWjTW4)

*Related Formal Open Type Proposal from Sahar Afshar and José Miguel Solé Bruning*

[https://github.com/Microsoft/OpenTypeDesignVariationAxisTags/blob/master/Proposals/Glyph\\_Extension\\_Axis/ProposalSummary.md](https://github.com/Microsoft/OpenTypeDesignVariationAxisTags/blob/master/Proposals/Glyph_Extension_Axis/ProposalSummary.md)

*Possibilities offered by variable font technology for Arabic script (Sahar Afshar Tweet)*

<https://twitter.com/sahafshar/status/896389627065577472>

### Language Issues related to Arabic Scripts

*The Persian Language (and what makes it fascinating) (Video 8:49)*

<https://www.youtube.com/watch?v=tZlIDNcbeE8>

*Comparing Urdu and Hindi Languages 11:52*

It may seem that the Hindi (written left-to-right in Devanagari Script) and Urdu (written right-to-left in Arabic Script) Languages are quite different. They are not.

[https://www.youtube.com/watch?v=vxSd7pli\\_TA](https://www.youtube.com/watch?v=vxSd7pli_TA)

*Urdu character names:*

[http://scriptsource.org/cms/scripts/page.php?item\\_id=entry\\_detail&uid=ygnh6uc3nb](http://scriptsource.org/cms/scripts/page.php?item_id=entry_detail&uid=ygnh6uc3nb)

Names of Farsi/Persian characters are the same as those used in Modern Standard Arabic and used in this document, although they are pronounced differently.

### Typography Issues related to Arabic Scripts (and others)

*Liberating Digital Type from the Metal Rectangle (Panel Discussion; Video 57:02)*

– Just van Rossum, Bianca Berning, John Hudson, Toshi Omagari, Victor Gaultney, Rob McKaughan, Sahar Afshar

<https://www.youtube.com/watch?v=P2JA5EgoDC0>

Summary:

<https://www.typtalks.com/news/2017/04/08/typo-labs-liberating-digital-type-from-the-metal-rectangle/>

### Coordinating Styles across differing Scripts – Some Useful Commentary

When combining multiple Scripts in a single document, many of the font matching characteristics used in Latin-only text (terms such as Serif, Sans Serif, Italic, x-height, etc.) have no meaning in other Scripts with which they may be displayed; this is particularly true when mixing Latin and Arabic Scripts. Also see Multi-script Database Design Note #4. Evaluating Fonts for use in Multi-Lingual Documents.

*Review of Diodrum Arabic Font – Sahar Afshar on 5 July 2017*




Sahar’s review concentrates on the stylistic aspects of matching Latin and Arabic Glyphs

<https://typographica.org/typeface-reviews/diodrum-arabic/>

*Review of Suisse Int’l Arabic Font – Huda Smitshuijzen AbiFarès on 5 July 2017*

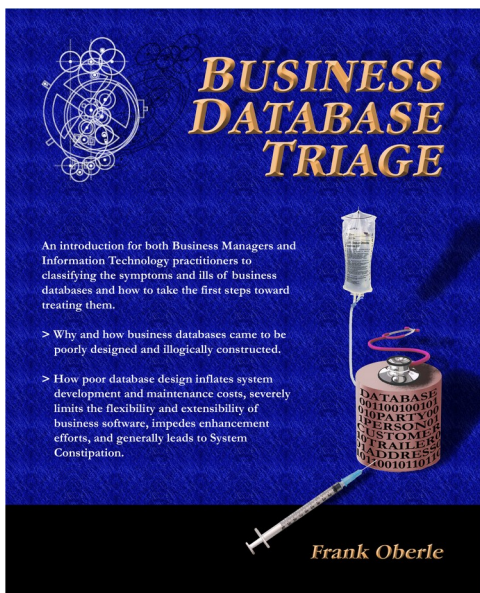
See particularly the sizing illustrations with overlays of Latin and Arabic Glyphs

<https://typographica.org/typeface-reviews/suisse-intl-arabic/>

قف Arabic:    Stop

\*

33 A useful reference, although the title should probably be “Arabic Script” rather than “Arabic Alphabet”



More information and sample pages at:  
[www.AntikytheraPubs.com](http://www.AntikytheraPubs.com)

In addition to an ongoing series of Database Design Notes, Antikythera Publications recently released the book “*Business Database Triage*” (ISBN-10: 0615916937) that demonstrates how commonly encountered business database designs often cause significant, although largely unrecognized, difficulties with the development and maintenance of application software. Examples in the book illustrate how some typical database designs impede the ability of software developers to respond to new business opportunities – a key requirement of most businesses.

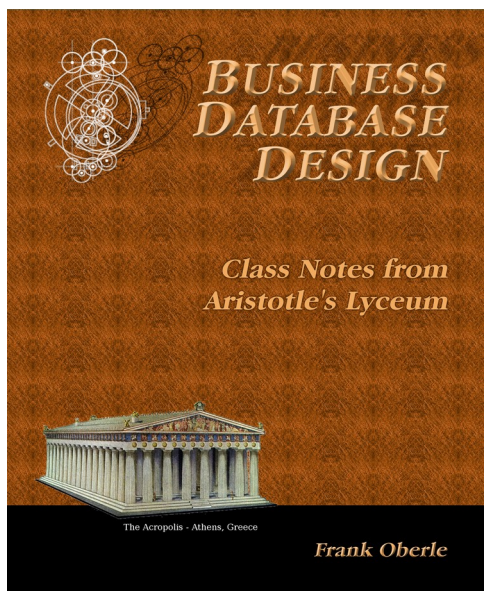
A number of examples of solutions to curing business system constipation are presented. Urban legends, such as the so-called object-relational impedance mismatch, are debunked – shown to be based mostly on illogical database (and sometimes object) designs.

“*Business Database Triage*” is available through major book retailers in most countries, or from the following on-line vendors, each of which has a full description of the book on their site:

CreateSpace: <https://www.createspace.com/4513537>

Amazon:

[www.amazon.com/Business-Database-Triage-Frank-Oberle/dp/0615916937](http://www.amazon.com/Business-Database-Triage-Frank-Oberle/dp/0615916937)



A follow-up book, “*Business Database Design – Class Notes from Aristotle’s Lyceum*” is due to be available in 2015.

“*Business Database Design*” leads the reader through the logical design and analysis techniques of data organization in more detail than the earlier work – which concentrated more on understanding and identifying problems caused by illogical database design rather than their solutions.

These logical approaches to data organization, espoused by Aristotle and an “A-List” of his successors, have formed the basis for scientific discovery over more than 2,400 years, and directly led to the technology we deal with today, notably including both relational and object theory.

“*Business Database Triage*” explained the reasons why these principles were virtually impossible to apply during the early years of our transition to the use of computers in business, but since the technology is now sufficiently mature that such compromises can no longer be justified, the time has come to relearn logical data organization techniques and apply them to our businesses.



The complete collection of publicly available Database Design Notes, including this series on handling Multi-Language/Multi-Script Databases, is available for free downloading from <http://www.antikytherapubs.com/>