# Confusing Character Shapes for Arabic Script Block

Arabic Script IDN Working Group (ASIWG)
Version 1.02
26 May 2010

Revision History:

| Date | Version | Changes | Changed by |
| --- | --- | --- | --- |
| 28 Aug 2008 | 1.0 | Created | Arfan Mansoor |
| 02 Sep 2008 | 1.01 | Review and Formatting | Arfan Mansoor and Atif Gulzar |
| 26 May 2010 | 1.02 | Reviewed | Rabia Sirhindi and Sarmad Hussain |

# 1. Objective

Arabic script is cursive in nature and each letter can have four different shapes: initial, medial, final and isolated according to its context in ligature. There are many characters in Arabic block of Unicode (U+06XX and U+0750 to U+077F) that are similar in one or more of these forms and these visual similarities between characters can cause confusions and may be major source of spoofing. The purpose of this document is to identify these confusing character shapes at different positions of ligature.

# 2. Categories

These confusing characters are divided into three different categories:

i)     Characters with same shape not distinct in any language

ii)    Confusable characters with similar shapes for at least one language but distinct in at least  one language

iii)   Confusable characters with similar shapes not distinct in any language

## 2.1  Characters with same shapes not distinct in any language

These characters are similar in shape at one or more positions and therefore they are confusable for any language using these characters. For example U+0649 (ى) and U+06CC (ی) have same isolated forms, whereas U+0643 (ك) and U+06A9 (ک) have same initial and medial forms although they have different isolated forms. All of these characters are discussed below.

## 2.2.1  'ک' Based Characters

ARABIC LETTER KAF (U+0643) and ARABIC LETTER KEHEH (U+06A9) which is used in Persian and Urdu are confusing at their initial and medial positions.

Table 1:  'ک' Based Characters

| Unicode | Initial Form | Medial Form | Final Form | Isolated Form |
|---|---|---|---|---|
| **U+06A9 (ک)** | ک | ک | ک | ک |

| U+0643 (ك) | ک | ک | كـ | ك |
|---|---|---|---|---|

## 2.2.2 'ه' Based Characters

HEH based characters are very confusing in Arabic script. ARABIC LETTER HEH (U+0647) and ARABIC LETTER HEH DOACHASHMEE (U+06BE) for Urdu are confusing at their initial, medial and isolated positions, while ARABIC LETTER HEH (U+0647), ARABIC LETTER HEH GOAL (U+06C1) for Urdu, and ARABIC LETTER AE (U+06D5) are confusing at their isolated positions.

Table 2: 'ه' Based Characters

| Unicode | Initial Form | Medial Form | Final Form | Isolated Form |
|---|---|---|---|---|
| U+0647 (ه) | ﻫ | ﻬ | ﻪ | ه |
| U+06BE (ﮬ) | ﮨ | ﮭ | ﮫ | ﮪ |
| U+06C1 (ہ) | ﮨ | ﮩ | ﮧ | ہ |
| U+06D5 (ە) | – | – | ﻪ | ە |

## 2.2.3 'ي' Based Characters

Although ARABIC LETTER FARSI YEH (U+06CC) does not have dots in isolated form, but it has two dots below at initial and medial position and confuse with ARABIC LETTER YEH (U+064A) at these positions. Furthermore, its isolated form is confusing with ARABIC LETTER ALEF MAKSURA (U+0649).

Table 3: 'ي' Based Characters

| Unicode | Initial Form | Medial Form | Final Form | Isolated Form |
|---|---|---|---|---|
| U+064A (ي) | ﻳ | ﻴ | ﻲ | ي |
| U+06CC (ی) | ﻳ | ﻴ | ﻰ | ی |

| U+0649 (ى) | - | - | ى | ک |
|---|---|---|---|---|

## 2.2.4 'ف' Based Characters[1]

ARABIC LETTER QAF WITH DOT ABOVE (U+06A7) and ARABIC LETTER FARSI FEH (U+0641) confuse at all (initial, medial, final, isolated) positions.

Table 4: 'ف' Based Characters

| Unicode | Initial Form | Medial Form | Final Form | Isolated Form |
|---|---|---|---|---|
| U+06A7 (ڧ) | ڧ | ڧ | ڧ | ڧ |
| U+0641 (ف) | ڧ | ڧ | ف | ف |

## 2.2.5 'ة' Based Characters

ARABIC LETTER TEH MARBUTA (U+0629) and ARABIC LETTER TEH MARBUTA GOAL (U+06C3) are confusing at their isolated positions.

Table 5: 'ة' Based Characters

| Unicode | Initial Form | Medial Form | Final Form | Isolated Form |
|---|---|---|---|---|
| U+0629 (ة) | - | - | ة | ة |
| U+06C3 (ۃ) | - | - | ـۃ | ۃ |

## 2.2.6 'ۀ' Based Characters

ARABIC LETTER HEH WITH YEH ABOVE (U+06C0) and ARABIC LETTER HEH GOAL WITH HAMZA ABOVE (U+06C2) are confusing at their isolated positions.

---

[1] In Arial Unicode MS, character under Numbering 2.2.4, 2.2.7, and 2.2.8 do not combine at their medial positions, but in other fonts like Tahoma (as in our examples) and Times New Roman they combine at their medial positions and create confusions.

Table 6: 'ۀ' Based Characters

| Unicode | Initial Form | Medial Form | Final Form | Isolated Form |
|---------|--------------|-------------|------------|---------------|
| U+06C0 (ۀ) | - | - | ـۀ | ۀ |
| U+06C2 (ۂ) | - | - | ـۂ | ۂ |

## 2.2.7  'ٹ' Based Character

ARABIC LETTER RNOON (U+06BB) and ARABIC LETTER TTEH for Urdu (U+0679) are confusing at their initial, final, and isolated positions.

Table 7: 'ٹ' Based Character

| Unicode | Initial Form | Medial Form | Final Form | Isolated Form |
|---------|--------------|-------------|------------|---------------|
| U+06BB (ڻ) | ڻ | ـڻ | ـڻ | ڻ |
| U+0679 (ٹ) | ٹ | ـٹـ | ـٹ | ٹ |

## 2.2.8  'ث' Based Characters

ARABIC LETTER NOON WITH THREE DOTS ABOVE (U+06BD) and ARABIC LETTER THEH (U+062B) are confusing at their all (initial, medial, final, isolated) positions.

Table 8: 'ث' Based Characters

| Unicode | Initial Form | Medial Form | Final Form | Isolated Form |
|---------|--------------|-------------|------------|---------------|
| U+06BD (ڽ) | ڽ | ـڽـ | ـڽ | ڽ |
| U+062B (ث) | ثـ | ـثـ | ـث | ث |

## 2.2.9  Confusing digits

Unicode provides two sets of digits for Arabic script. One from U+0660 to U+0669 is for Arabic digits and second from U+06F0 to U+06F9 is used for Eastern Arabic-Indic digits for languages of Iran, Pakistan and India (Persian,

Sindhi, Urdu etc.). Confusing character with same shapes from these two blocks are listed below.

Table 9:  Confusing digits (Arabic-Indic and Eastern Arabic-Indic)

| Arabic-Indic | Eastern Arabic-Indic |
|---|---|
| ٠(U+0660) | �۰(U+06F0) |
| ١ (U+0661) | ۱ (U+06F1) |
| ٢ (U+0662) | ۲ (U+06F2) |
| ٣ (U+0663) | ۳ (U+06F3) |
| ٥ (U+0665) | ۵ (U+06F5) |
| ٧ (U+0667) | ۷ (U+06F7) |
| ٨ (U+0668) | ۸ (U+06F8) |
| ٩ (U+0669) | ۹ (U+06F9) |

In case of digits we have one more problem. The codepoints (06F0 to 06F9) are defined for Persian, Sindhi and Urdu, but some of these characters have different shapes e.g. U+06F4, U+06F6, and U+06F7 have different glyphs in Persian than Sindhi and Urdu while U+06F4 has different glyphs in Sindhi, and Urdu.

Table 10:  Confusing digits (Eastern Arabic-Indic)

| Urdu and Sindhi | Persian |
|---|---|
| ۰(U+06F0) | ۰(U+06F0) |
| ١ (U+06F1) | ۱ (U+06F1) |
| ٢٢ (U+06F2) | ۲ (U+06F2) |
| ٣ (U+06F3) | ۳ (U+06F3) |
| ۴ (U+06F4) (Urdu) | ۴ (U+06F4) |
| ۴ (U+06F4) (Sindhi) | |

| | |
|---|---|
| ۵ (U+06F5) | ۵ (U+06F5) |
| ۶ (U+06F6) | ۶ (U+06F6) |
| ٧ (U+06F7) | ٧ (U+06F7) |
| ٨ (U+06F8) | ٨ (U+06F8) |
| ٩ (U+06F9) | ٩ (U+06F9) |

## 2.2  Confusable characters with similar shapes for at least one language but distinct in at least one language

There are characters that distinct in particular language but they may be confusable for users of other languages. For example U+06A9 (ک) and U+06AA (ڪ) are distinct Sindhi letters but confusing for Urdu speakers. Similarly U+064A (ي) and U+06CC (ی) are distinct Pashto letter but are confusing for Urdu speakers.  These letters are discussed below:

### 2.2.1  'ک' Based Character

ARABIC LETTER KEHEH (U+06A9) and ARABIC LETTER SWASH KAF (U+06AA) and ARABIC LETTER KAF (U+0643) are shape variants in Arabic. First two characters are confusing in Urdu but in Sindhi these are distinct characters.

Table 11: 'ک' Based Character

| Unicode | Initial Form | Medial Form | Final Form | Isolated Form |
|---|---|---|---|---|
| U+06A9 (ک) | ک | ک | ک | ک |
| U+06AA (ڪ) | ڪ | ڪ | ڪ | ڪ |
| U+0643 (ك) | ك | ك | ك | ك |

### 2.2.2  'ی' Based Character

ARABIC LETTER FRASI YEH (U+06CC) for Persian and Urdu is confusing with ARABIC LETTER YEH (U+064A) at its initial and medial positions and with ARABIC LETTER ALEF MAKSURA (U+0649) at its final and isolated positions.

Table 12: 'ى' Based Character

| Unicode | Initial Form | Medial Form | Final Form | Isolated Form |
|---|---|---|---|---|
| **U+064A (ي)** | يـ | ـيـ | ـي | ي |
| **U+06CC (ی)** | یـ | ـیـ | ـی | ی |
| **U+0649 (ى)** | – | – | ـى | ى |

## 2.2.3  ALEF with HAMZA ABOVE letters

ARABIC LETTER ALEF WITH HAMZA ABOVE (U+0623) and ARABIC LETTER ALEF WITH WAVY HAMZA ABOVE (U+0672) for Baluchi and Kashmiri are confusing in their final and isolated forms.

Table 13: ALEF with HAMZA ABOVE letters

| Unicode | Initial Form | Medial Form | Final Form | Isolated Form |
|---|---|---|---|---|
| **U+0623 (أ)** | – | – | ـأ | أ |
| **U+0672 (ٵ)** | – | – | ـٵ | ٵ |

## 2.2.4  ALEF with HAMZA BELOW letters

ARABIC LETTER 'ALEF WITH HAMZA BELOW' (U+0625) and ARABIC LETTER 'ALEF WITH WAVY HAMZA BELOW' for Baluchi and Kashmiri (U+0673) are confusing at their final and isolated positions.

Table 14: ALEF with HAMZA BELOW letters

| Unicode | Initial Form | Medial Form | Final Form | Isolated Form |
|---|---|---|---|---|
| U+0625 (إ) | - | - | إ | إ |
| U+0673 (ٳ) | - | - | ٳ | ٳ |

### 2.2.5 Shape variant confusion

In addition to these categories there may be confusable characters within some language. For example U+06D2 (ے) and U+064A (ي) are confusable for Arabic speakers even though these are not similar in shape, in Arabic these are considered stylistic variant of same character U+064A.

## 2.3 Confusable characters with similar shapes not distinct in any language

There are characters in different languages with slight variations e.g. vertical dots vs. horizontal dots, three dots pointing downwards vs. three dots pointing upward below etc. These characters can also create confusions for users. These characters are listed below.

Table 15: Confusing Characters with Similar Shapes

| | Unicode | Characters | | Remarks |
|---|---|---|---|---|
| i)<br>ii) | U+062A<br>U+067A | i)    ت<br>ii)    ٺ | i)<br>ii) | ARABIC LETTER THE<br>ARABIC LETTER TTEHEEH |
| i)<br>ii) | U+062B<br>U+067D | i)    ث<br>ii)    ٽ | i)<br><br>ii) | ARABIC LETTER THEH<br>ARABIC LETTER THE WITH THREE DOTS (Sindhi) |
| i)<br>ii) | U+063C<br>U+0764 | i)    ڜ<br>ii)    ݤ | i)<br><br><br>ii) | ARABIC LETTER KEHEH WITH THREE DOTS BELOW<br>ARABIC LETTER |

| | | | | |
|---|---|---|---|---|
| | | | | KEHEH WITH THREE DOTS BELOW POINTING UPWARDS BELOW |
| i)<br>ii) | U+064A<br>U+06D0 | i) ي<br>ii) ې | i)<br><br>ii) | ARABIC LETTER YEH<br>ARABIC LETTER E (Pashto and used as letter bbeh in Sindhi) |
| i)<br>ii) | U+067E<br>U+0752 | i) پ<br>i) ݒ | i)<br><br>ii) | ARABIC LETTER PEH (Urdu, Persian)<br>ARABIC LETTER BEH WITH THREE DOTS PONITNG UPWARDS BELOW |
| i)<br>ii) | U+0683<br>U+0684 | i) ڃ<br>ii) ڄ | i)<br><br>ii) | ARABIC LETTER NYEH (Sindhi)<br>ARABIC LETTER DYEH (Sindhi) |
| i)<br>ii) | U+0686<br>U+0758 | ii) چ<br>iii) ژ | i)<br><br><br>ii) | ARABIC LETTER TCHEH (Urdu, Persian)<br>ARABIC LETTER HAH WITH THREE DOTS POINTING UPWARDS BELOW |
| i)<br>ii) | U+068E<br>U+068F | i) ڎ<br>ii) ڏ | i)<br><br>ii) | ARABIC LETTER DUL (Burushaski)<br>ARABIC LETTER DAL WITH THREE DOTS ABOVE DOWNWORDS (Sindhi) |
| i)<br>ii) | U+06A0<br>U+075F | i) ڠ<br>iv) ݟ | i)<br><br><br>ii) | ARABIC LETTER AIN WITH THREE DOTS ABOVE<br>ARABIC LETTER AIN WITH THREE POINTING DOWNWORDS ABOVE |
| i)<br>ii) | U+06B2<br>U+06B3 | i) ڲ<br>ii) ڳ | i)<br><br><br>ii) | ARABIC LETTER GAF WITH TWO DOTS BELOW<br>ARABIC LETTER GUEH (Sindhi) |

| i) U+075D<br>ii) U+075E | i) غّ<br>ii) غ٘ | i) ARABIC LETTER AIN WITH TWO DOTS ABOVE<br>ii) ARABIC LETTER AIN WITH TWO DOTS VERTICALLY ABOVE |
| --- | --- | --- |
| i) U+0697<br>ii) U+076B | i) ژ<br>ii) ګ | i) ARABIC LETTER REH WITH TWO DOTS ABOVE (Dragwa)<br>ii) ARABIC LETTER REH WITH TWO DOTS VERTICALLY ABOVE (Torwali, Ormuri) |
| i) U+0649<br>ii) U+06CD<br>iii) U+06CC | i) ى<br>ii) ۍ<br>iii) ی | i) ARABIC LETTER ALEF MAKSURA (YEH shaped letter with no dots at any position)<br>ii) ARABIC LETTER YEH WITH TAIL (Pashto, Sindhi)<br>iii) ARABIC LETTER FARSI YEH (Arabic, Persian, Urdu) |

## 3. Conclusion

In Arabic script characters visually appear similar at different positions (initial, medial, final, and isolated). Presence of such characters causes confusions and they can be potential source of spoofing. Therefore confusing character's shape problem has to be solved for secure implementation of Internationalized Domain Names (IDNs) and for Email Address Internationalization (EAI). In this document we identified confusing characters in Arabic script and we have divided these confusing characters into three categories. First category consist of those confusing characters that are at script level i.e. not distinct in any language. In second category those characters are

identified that are confusing for at least one language but distinct in at least one language. Third category consists of those characters with similar shapes not distinct in any language.