# Normalization Character List for Arabic Script Block

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

Revision History:

| Date | Version | Changes | Changed by |
|---|---|---|---|
| 04 Sep 2008 | 1.01 | Created | Arfan Mansoor |
| 26 May 2010 | 1.02 | Reviewed | Rabia Sirhindi and Sarmad Hussain |

## Objective

In Arabic script block we have some characters that can be written in two forms: i) composed form and ii) decomposed form. In composed form characters occur as a single entity in Unicode block i.e. single codepoint is used to for that character e.g. آ (U+0622). In decomposed form the same character can be written by the combination two or more codepoints e.g. آ (U+0622) can be written by the combination of ا (U+0627) and ٓ (U+0653). The formation of same character by more than one way can cause security risks because it is potential source of phishing problem in Internationalized Domain Names (IDNs).

The purpose of this document is to identify all such characters in Arabic script that can be formed in more than one way. Below is the list of these characters that exits in composed and decomposed forms in Arabic block. All such cases have to be normalized to make IDNs more secure.

| Combining Mark | Composed Form | Decomposed Form | Unicode Normalized Form |
|---|---|---|---|
| ٓ<br>U+0653 | آ<br>U+0622 | ٓ ا<br>U+0627 U+0653 | Defined |
| ٔ<br><br><br><br>U+0654 | أ<br>U+0623 | ٔ ا<br>U+0627 U+0654 | Defined |
| | ؤ<br>U+0624 | ٔ و<br>U+0648 U+0654 | Defined |
| | | ٔ ي<br>U+064A U+0654 | Defined |
| | | ٔ ى | Not Defined |

| Combining Mark | Composed Form | Decomposed Form | Unicode Normalized Form |
|---|---|---|---|
| | ئ<br>U+0626 | U+0649 U+0654 | |
| | | ◌ی<br>U+06CC U+0654 | Not Defined |
| | ۀ<br>U+06C0 | ◌ە<br>U+06D5 U+0654 | Defined |
| | | ◌ه<br>U+0647 U+0654 | Not Defined |
| | ۂ<br>U+06C2 | ◌ہ<br>U+06C1 U+0654 | Defined |
| | | ◌ه<br>U+0647 U+0654 | Not Defined |
| | ۓ<br>U+06D3 | ◌ے<br>U+06D2 U+0654 | Defined |
| | ځ<br>U+0681 | ◌ح<br>U+062D U+0654 | Not Defined |
| | ݬ<br>U+076C | ◌ر<br>U+0631 U+0654 | Not Defined |
| ◌ٕ<br>U+0655 | إ<br>U+0625 | ◌ٕا<br>U+0627 U+0655 | Defined |

| Combining Mark | Composed Form | Decomposed Form | Unicode Normalized Form |
|---|---|---|---|
| ـُ<br>U+064F | ۇ<br>U+06C7 | ۇ ّ<br>U+0648 U+064F | Not Defined |
| | | ۇ ّ<br>U+0648 U+0619 | Not Defined |
| ـٰ<br>U+0670 | ۈ<br>U+06C8 | ۈ ّ<br>U+0648 U+0670 | Not Defined |
| ـٰ<br>U+06EC | ۏ<br>U+06CF | ۏ ّ<br>U+0648 U+06EC | Not Defined |
| | غ<br>U+063A | غ ّ<br>U+0639 U+06EC | Not Defined |
| | ض<br>U+0636 | ض ّ<br>U+0635 U+06EC | Not Defined |
| | خ<br>U+062E | خ ّ<br>U+062D U+06EC | Not Defined |
| | ﭺ<br>U+06BF | ﭺ ّ<br>U+0686 U+06EC | Not Defined |
| | ذ<br>U+0630 | ذ ّ<br>U+062F U+06EC | Not Defined |
| | ز<br>U+0632 | ز ّ<br>U+0631 U+06EC | Not Defined |
| | ڶ<br>U+06B6 | ڶ ّ<br>U+0644 U+06EC | Not Defined |
| | ف<br>| ف ّ<br>| Not Defined |

| Combining Mark | Composed Form | Decomposed Form | Unicode Normalized Form |
|---|---|---|---|
| | U+06A7 | U+066F U+06EC | |
| | ڧ<br>U+0641 | ۅ<br>U+06A1 U+06EC | Not Defined |
| | ن<br>U+0646 | ۑ<br>U+06BA U+06EC | Not Defined |
| | ڬ<br>U+06AC | ك<br>U+0643 U+06EC | Not Defined |
| | ڴ<br>U+0762 | ک<br>U+06A9 U+06EC | Not Defined |
| | ݥ<br>U+0765 | م<br>U+0645 U+06EC | Not Defined |
| U+0615 | ڂ<br>U+0772 | ح<br>U+062D U+0615 | Not Defined |
| | ٹ<br>U+0679 | ٮ<br>U+066E U+0615 | Not Defined |
| | ڑ<br>U+0691 | ر<br>U+0631 U+0615 | Not Defined |
| | ڈ<br>U+0688 | د<br>U+062F U+0615 | Not Defined |
| | ڗ<br>U+0771 | ز<br>U+0697 U+0615 | Not Defined |
| | ن<br> | ن<br> | Not Defined |

| Combining Mark | Composed Form | Decomposed Form | Unicode Normalized Form |
|---|---|---|---|
| | U+0768 | U+0646 U+0615 | |
| | ڋ<br>U+068B | ڊ ٕ<br>U+068A U+0615 | Not Defined |
| | ڻ<br>U+06BB | ں ٕ<br>U+06BA U+0615 | Not Defined |
| ٛ<br>U+065B | ئ<br>U+063D | ی ٛ<br>U+06CC U+065B | Not Defined |
| | ۉ<br>U+06C9 | و ٛ<br>U+0648 U+065B | Not Defined |
| | ݾ<br>U+077E | س ٛ<br>U+0633 U+065B | Not Defined |
| | ڮ<br>U+06EE | د ٛ<br>U+062F U+065B | Not Defined |
| | ۯ<br>U+06EF | ر ٛ<br>U+0631 U+065B | Not Defined |
| | ۿ<br>U+06FF | ھ ٛ<br>U+06BE U+065B | Not Defined |
| | | ه ٛ<br>U+0647 U+065B | Not Defined |
| ۛ<br>U+06DB | ۓ<br>U+063F | ی ۛ<br>U+06CC U+06DB | Not Defined |
| | | ی ۛ | Not Defined |

| Combining Mark | Composed Form | Decomposed Form | Unicode Normalized Form |
|---|---|---|---|
| | | U+0649 U+06DB | |
| | ش<br>U+0634 | س<br>U+0633 U+06DB | Not Defined |
| | ݜ<br>U+069C | ݛ<br>U+069B U+06DB | Not Defined |
| | ث<br>U+062B | ٮ<br>U+066E U+06DB | Not Defined |
| | څ<br>U+0685 | ح<br>U+062D U+06DB | Not Defined |
| | ژ<br>U+0698 | ر<br>U+0631 U+06DB | Not Defined |
| | ڎ<br>U+068E | د<br>U+062F U+06DB | Not Defined |
| | ڠ<br>U+06A0 | ع<br>U+0639 U+06DB | Not Defined |
| | ڤ<br>U+06A4 | ڡ<br>U+06A1 U+06DB | Not Defined |
| | ڨ<br>U+06A8 | ٯ<br>U+066F U+06DB | Not Defined |
| | ڭ<br>U+06AD | ك<br>U+0643 U+06DB | Not Defined |
| | ڴ<br>U+06B4 | گ<br>U+06AF U+06DB | Not Defined |
| | ڷ | ل | Not Defined |

| Combining Mark | Composed Form | Decomposed Form | Unicode Normalized Form |
|---|---|---|---|
| | U+06B7 | U+0644 U+06DB | |
| | ث<br>U+06BD | ں ٛ<br>U+06BA U+06DB | Not Defined |
| | ݣ<br>U+0763 | ک ٛ<br>U+06A9 U+06DB | Not Defined |
| ◌ٜ<br>U+065C | ب<br>U+0628 | ٮ ٜ<br>U+066E U+065C | Not Defined |
| | ڊ<br>U+068A | د ٜ<br>U+062F U+065C | Not Defined |
| | ڋ<br>U+068B | ڈ ٜ<br>U+0688 U+065C | Not Defined |
| | ڔ<br>U+0694 | ر ٜ<br>U+0631 U+065C | Not Defined |
| | ڣ<br>U+06A3 | ف ٜ<br>U+0641 U+065C | Not Defined |
| | ڹ<br>U+06B9 | ن ٜ<br>U+0646 U+065C | Not Defined |
| | ػ<br>U+06FC | غ ٜ<br>U+063A U+065C | Not Defined |
| | ڬ<br>U+06FB | ض ٜ<br>U+0636 U+065C | Not Defined |
| | �̪ب<br>U+0751 | ث ٜ<br>U+062B U+065C | Not Defined |

| Combining Mark | Composed Form | Decomposed Form | Unicode Normalized Form |
|:---:|:---:|:---:|:---:|
|  | ﭦ<br>U+0766 | م ◌<br>U+0645 U+065C | Not Defined |
| ◌<br>U+065A | ڵ<br>U+06B5 | ل ◌<br>U+0644 U+065A | Not Defined |
|  | ۆ<br>U+06C6 | و ◌<br>U+0648 U+065A | Not Defined |
|  | ئ<br>U+06CE | ى ◌<br>U+06CC U+065A | Not Defined |
|  |  | ى ◌<br>U+0649 U+065A | Not Defined |
|  | ݖ<br>U+0756 | ٮ ◌<br>U+066E U+065A | Not Defined |
|  | ݩ<br>U+0769 | ن ◌<br>U+0646 U+065A | Not Defined |

In addition to these characters there are some characters that are formed by more than two characters combinations. These characters are listed in table given below.

| | | | |
|:---:|:---:|:---:|:---:|
| ښ<br>U+069A | س<br>U+0633 | ◌<br>U+065C | ◌<br>U+06EC |
| ڣ<br>U+06A3 | ڡ<br>U+06A1 | ◌<br>U+065C | ◌<br>U+06EC |
| ۺ<br>U+06FA | س<br>U+0633 | ◌<br>U+06DB | ◌<br>U+065C |
| ض | ص | ◌ | ◌ |

| U+06FB | U+0635   U+065C   U+06EC |
|--------|--------------------------|
| ڻ U+06FC | ع U+0639   U+065C   U+06EC |

## Conclusion

Occurrence of characters in more than one form in Unicode Arabic script block can be a major security risk and these needs to be normalized to make IDNs more secure. In this document we have identified all those characters for Arabic script block that can occur in more than one form.