

Inputting of Arabic inline characters in Unicode

Problem Statement

Unicode currently lacks a consistent and complete specification for handling Arabic inline characters. In this document we will consider two such characters:

1. *Hamzah*
2. *ʿAlif kharjiyyah* “Dagger Alef”

Hamzah

The case for inline *hamzah* has been discussed in some detail but, it appears, without much official resolution. <https://unicode.org/L2/L2014/14109-inline-chars.pdf> <https://unicode.org/L2/L2013/13226-koran-ortho.pdf> Section IV

We will attempt to summarize the problem here.

Being a later addition to the Arabic script, *hamzah* is written as a diacritic on the basic skeletal text. This *hamzah* can be written in two different ways:

1. Above a base letter, known as the “seat of the *hamzah*”: In writing the Arabic language *hamzah* may be seated above a vowel letter: أ و ي. It may also be sometimes written below vowel letters, depending on the writing style and vowel mark: إ ي. In Persian language writing, it may also appear above a heh: ه.
2. After (to the left of) a letter in the basic skeletal text. This has two sub-cases:
 - a. Standalone *hamzah*: After a non-joining base letter (ا د ذ ر ز و) or at the end of a word. In this case, *hamzah* is written at the baseline. Examples: سوء, دعاء.
 - b. Inline *hamzah*: After a joining base letter in the middle of a word. In this case *hamzah* is written between the two joining base letters, above the baseline, and without affecting the connection of the two base letters. Examples: شيءٌ, خطيئة, برئين

For a set of rules to determine the seat of the *hamzah* and whether *hamzah* is to be written inline, see this article.

For case 1, Unicode provides U+0654 ARABIC HAMZA ABOVE and U+0655 ARABIC HAMZA BELOW as combining diacritics.

For case 2a, Unicode provides U+0621 ARABIC LETTER HAMZA.

In order to display *hamzah* according to case 2b, most users are accustomed to approximating it with ّ which between two joining letters looks similar to an inline *hamzah* ّ. In fact, due to the prevalence of this approximation, many users will consider it the correct way of writing, even in handwritten text.

This solution is not ideal because it introduces a character in the basic skeletal text (tatweel) which is not part of the word's spelling. This can also adversely affect computerized search and indexing of text encoded in this manner.

[illegible]

However, this solution is problematic for some of the same reasons as for *hamzaʿi*: It introduces characters in the text which are not part of the word’s spelling: tatweel and NBSP/NNBSP. The user has to choose between between tatweel and NBSP/NNBSP simply on the basis of whether the preceding character is joining or non-joining. The computer’s shaping engine already handles joining of base letters. Therefore, this choice left to the user is non-semantic and unnecessary. It may adversely affect computerized search and indexing.

Arabic has the need to display inline characters as diacritics that: May generally be either above or after (to the left of) a base character. The two (above and after) are to be differentiated as they may affect pronunciation. Must not affect the joining of base letters. Solution We propose specifying the use of U+034F Combining Grapheme Joiner to determine whether a diacritic appears above or after a base letter.

Also that U+0654 Arabic *hamza* above be used for the diacritic inline *hamza* and U+0670 arabic superscript alef continue to be used for the diacritic dagger alef.

If the diacritic is to appear above the base letter then it shall be input after the letter without intermediate CGJ.

If the diacritic is to appear after the base letter then it shall be input after the letter with an intermediate CGJ.

CGJ and/or the diacritic shall not affect the joining of base letters.

We recommend to font developers and typographers that CGJ + *hamza* above after a non-joining letter be rendered identical to U+0621.

Here are some examples of the proposed input sequences and the displayed rendering: