Notes on some Unicode Arabic characters: recommendations for usage

Jonathan Kew

Draft 2 — April 21, 2005

Contents

1	Introduction	2
2	KAF-based letters	2
	2.1 Arabic	2
	2.2 Persian	3
	2.3 Urdu	3
	2.4 Sindhi	4
	2.5 Jawi (Malay) gaf	4
	2.6 Moroccan Arabic gaf	4
	2.7 Uighur, Kirghiz and Kazakh eng	5
3	HEH-based letters	5
	3.1 Arabic	5
	3.2 Persian	6
	3.3 Urdu	6
	3.4 Sindhi	7
	3.5 Parkari	7
	3.6 Kurdish	7
4	YEH-based letters	8
	4.1 Arabic	8
	4.2 Persian	8
	4.3 Urdu	9
	4.4 Sindhi	9
	4.5 Kurdish	9
	4.6 Uighur	9
5	For font designers: summary of heh glyph variants	10

1 Introduction

In certain cases, the Unicode standard encodes separate characters for forms that would be considered glyph variants of a single character in Arabic. While this is sometimes necessary, in order to support writing systems where the shapes are used contrastively, it also raises sometimes raises questions of which character to use, among several possibilities. This document discusses some of these situations, and attempts to offer guidance for implementers and users of the Standard.

To an Arabic reader, the glyphs $\stackrel{\checkmark}{\smile}$, and $\stackrel{\checkmark}{\smile}$ are all clearly recognizable as forms of the same letter, kaf. The first, $\stackrel{\checkmark}{\smile}$, is typical of the designs seen in common text typefaces based on a simplified Naskh style of writing. $\stackrel{\checkmark}{\smile}$ is an alternate form that seems to be based on Nastaliq style, and $\stackrel{\checkmark}{\smile}$ is a swash form sometimes used, normally in initial or medial position, for stylistic effect or as part of line justification. Similarly, $\stackrel{\checkmark}{\smile}$ and $\stackrel{\checkmark}{\smile}$ are both yeh, the dots being optional.

However, as the Arabic script has been adopted and adapted for writing many other languages, these different shapes have sometimes been taken and used as distinct letters in such writing systems. Even where the alternate forms of a single Arabic letter are not used contrastively within a single writing system, the range of shapes that are recognized and accepted may be much more restricted than was the case with the original Arabic letter.

Note that this document does not discuss the "presentation forms" of Arabic letters. These are not recommended for encoding data; they exist only for legacy compatibility reasons. Thus, except where the context specifically refers to joining forms, references here to different "shapes", "forms", or "glyphs" for a given Unicode character are not referring to the initial, medial, and final linking forms, or to ligatures, but to different designs of the basic unjoined letter (and correspondingly different linked forms).

Not every character nor every language is discussed here (far from it); however, it is hoped that the principles used can be applied where similar encoding choices need to be made for other writing systems and additional letters.

Some of the recommendations given here are based in part on the presentation *Guidelines to Use* of *Arabic Characters* by Kamal Mansour at the 24th Internationalization and Unicode Conference (September 2003) in Atlanta, GA. Others are based on discussions with specialists studying various of the languages concerned, and on experience gained in implementing a variety of fonts and software systems.

2 KAF-based letters

Here, we consider the Unicode characters U+0643 $\stackrel{\checkmark}{\smile}$, U+06A9 $\stackrel{\checkmark}{\smile}$, and U+06AA $\stackrel{\hookleftarrow}{\smile}$, and other characters based on these forms. These are all forms of the Arabic letter *kaf*, written in different styles.

I am not aware of any language whose writing system uses both $\[Delta]$ and $\[Delta]$ contrastively; indeed, this seems highly unlikely, as in both initial and medial positions, their linked forms are the same: $\[Delta]$ joins as $\[Delta]$, while $\[Delta]$ joins as $\[Delta]$. On the other hand, $\[Delta]$ and $\[Delta]$ do occur together and must be distinguished; and in some writing systems, the default shape of U+0643 $\[Delta]$ is not considered correct for $\[Delta]$. Similarly, where the alphabet has been extended by the addition of dots or other marks to $\[Delta]$, this may apply only to one specific shape of the letter.

2.1 Arabic

The Arabic letter kaf is encoded as U+0643 $\stackrel{\ \ \ \ }{\sim}$. Depending on the type design, and possibly other stylistic factors, this character might be rendered with forms more like $\stackrel{\ \ \ \ }{\sim}$ or $\stackrel{\ \ \ \ }{\hookrightarrow}$, but kaf in Arabic

should nevertheless always be encoded with U+0643. The selection of alternate glyphs would occur as a result of typeface choice, formatting processes, and higher-level protocols, without altering the encoded text.

In the absence of specific reasons to use a different *kaf* character, U+0643 should also be considered the default choice to encode the corresponding /k/ letter in other languages where the Arabic script is used. However, if the script has been adopted not directly from Arabic, but from another source such as Persian or Sindhi, the practices of that more immediate source should generally be considered first.

- use U+0643 **4** for *kaf*
- U+06A9 \leq and U+06AA \leq should *not* be used for stylistic effect

2.2 Persian

In Persian (Farsi), the typical Arabic shape $\stackrel{\ \ }{}$ is not considered an acceptable form for *kaf*. The standard *Information Technology – Persian Information Interchange and Display Mechanism, using Unicode* (ISIRI 6219)¹ recommends the use of U+06A9 of or Persian *kaf*, permitting both Arabic and Persian forms to co-occur in plain text without needing markup or other higher-level protocols to distinguish the two.

While the recommendation is to use U+06A9 of for kaf when encoding Persian text in Unicode, users should be aware that there is likely to be a considerable amount of Persian text where U+0643 is used, making no distinction from Arabic kaf. In many cases, Arabic fonts have been "adapted" for Persian by simply changing the glyph at U+0643 (and its corresponding final form), to obtain the correct Persian appearance with software systems (keyboards, mappings from legacy codepages, etc.) that were designed for Arabic.

Therefore, while *producers* of Persian text should use $U+06A9 \checkmark$ for *kaf*, it may be advisable for *consumers* of Persian text data, especially if accepting input data from arbitrary sources, to recognize U+0643 as well, perhaps offering an option to remap this code to U+06A9 if appropriate.

- use U+06A9 **5** for *kaf*
- U+0643 \(\delta\) for kaf may be encountered in data

2.3 Urdu

Urdu tends to follow Persian writing conventions more closely than Arabic, and in particular the shape \leq is clearly the preferred kaf, with $\stackrel{\ \ \ \ \ \ \ \ \ }{\ \ \ \ \ \ \ }$ being viewed as Arabic and "foreign". This preference probably arises because Urdu is almost universally written in Nastaliq style script, where the form of kaf resembles $\stackrel{\ \ \ \ \ \ \ }{\ \ \ \ \ }$ (even when the language is Arabic); however, in Urdu the preference is so strongly established that $\stackrel{\ \ \ \ \ \ \ }{\ \ \ \ \ }$ would be considered incorrect even in non-Nastaliq styles, rather than being seen as dependent on the style in use. (The history is probably similar for Persian, which also has a long tradition of Nastaliq calligraphy, even though that style is less widely used now.)

The same encoding recommendation therefore applies for Urdu as for Persian:

- use U+06A9 \leq for *kaf*
- U+0643 4 for kaf may be encountered in data

¹See http://www.farsiweb.info/standard/; note that the document is in Persian.

2.4 Sindhi

The Sindhi language has a contrast between unaspirated and aspirated consonants. When the Arabic script was adopted and extended to write Sindhi, the form \checkmark was used to represent an aspirated velar consonant /kh/, while the form \hookrightarrow was used for the unaspirated /k/. The form \checkmark is not used in writing Sindhi.

To encode Sindhi, then, the two Unicode characters U+06AA \subseteq and U+06A9 \subset should be used for /k/ and /kh/ respectively. It is probably less likely that U+0643 will be found in Sindhi data than in Persian or Urdu, as Sindhi does not have the same history as Persian and Urdu of legacy implementations based on slightly-extended Arabic systems with a few glyph changes. If it does occur in Sindhi text, it will most likely be representing /kh/ (properly encoded as U+06A9), as in some positions these share similar glyph shapes.

(It may be interesting to note that the Unicode character name of U+06A9 С ARABIC LETTER КЕНЕН looks like an attempt to indicate in transcription the *aspirated kaf* sound of Sindhi. This supports the view that this character was encoded, perhaps originally in a legacy codepage, specifically for the contrastive Sindhi /kh/ usage where is not a recognized form.)

- use U+06A9 \leq for aspirated *kaf*/kh/
- use U+06AA $\stackrel{\checkmark}{=}$ for unaspirated *kaf*/k/
- U+0643 \(\delta\) should not occur, but probably represents /kh/ if encountered in data

2.5 Jawi (Malay) gaf

Malay written in Arabic script (known as Jawi) uses a kaf modified by the addition of a dot above to represent a voiced consonant /g/. This could be encoded using U+06AC $\dot{\omega}$, and indeed the Names List annotation found in Unicode versions up to 4.0 suggests this. However, old Malay sources consistently write this character as $\dot{\omega}$, using the Persian kaf as a base and not the Arabic kaf. This is true even where the Malay sources use $\dot{\omega}$ for kaf, and applies to both printed and hand-written materials. The form $\dot{\omega}$ does not appear to be a legitimate rendering of Jawi gaf.

The strength of the preference for the shape in rather than in may be gauged from the fact that some writers, faced with computer systems that only provided U+06AC in have used this character but added a *kashida* (extender) character after it in final or isolated position, in order to get a printed result such as in Although this is typographically quite unsatisfactory, it has been preferred over the is shape.

It is therefore recommended that Jawi *gaf* be encoded as U+0762 (newly added in Unicode version 4.1); the use of U+06AC is *not* recommended, though it may be found in some existing text data, especially in view of the fact that in Unicode versions prior to 4.1, U+0762 was not encoded. The character U+06AC should be used only for languages where its nominal form would be an acceptable, recognized way to write the relevant letter.

- use U+0643 <u>4</u> for *kaf*
- use U+0762 5 for gaf
- U+06AC 4 for gaf may be encountered in existing data

2.6 Moroccan Arabic gaf

Like Malay, Moroccan Arabic adds a *gaf* letter to the standard Arabic alphabet. In this case, it is written as a *kaf* with three dots above. However, like the Jawi (Malay) case, the base form used is consistently \leq and not \leq , even though the \leq shape is used for *kaf*. Just as with Malay, there are

sources that show authors deliberately using a final *kashida* to force a \angle -shaped character to take on the shape \angle , so strong is the feeling that \angle rather than \angle is the correct form for *gaf*.

Therefore, it is recommended that Moroccan *gaf* should be represented as U+0763 (new in Unicode version 4.1), and not as U+06AD. The latter should be used for languages where a based shape is accepted, as appears to be the case for Uighur *eng*, for example.

- use U+0643 **4** for *kaf*
- use U+0763 $\stackrel{*}{\smile}$ for gaf

2.7 Uighur, Kirghiz and Kazakh eng

In a number of Central Asian languages, the letter *eng* (a velar nasal) appears to be an example where an extended Arabic letter, constructed as a *kaf* with three dots above, has been based directly on the traditional shape. Although these writing systems use the Persian and Urdu form U+06AF for *gaf*, which is based on a shape, the *eng* is typically written as and should therefore be encoded as U+06AD.

- use U+0643 **4** for *kaf*
- use U+06AF \mathcal{S} for gaf
- use U+06AD 5 for eng

3 HEH-based letters

Like kaf, the letter heh also has several variant glyphs. The nominal character may be written as either \bullet or \bullet . In this case, there are also variants of the associated joining forms; thus, a document where isolated heh is written as \bullet might use either \bullet or \bullet as the medial form, and any of \bullet , \bullet or \bullet as the corresponding final form. In some cases, these variants may be freely interchangeable at the discretion of the writer or type designer, but in some languages and for some extended Arabic letters based on a heh form, there are specific requirements concerning the shapes used.

Unicode encodes arabic letter heh as U+0647 , with alternate possibilities including heh doachashmee at U+06BE , heh goal at U+06C1 , and letter ae at U+06D5 . In addition, there are various modified forms of some of these, with dots or other marks added.

To better distinguish among these *heh* letters, we also consider typical linked glyphs used for each character (remembering that there may be variation in these):

USV	Name	Isolate	Final	Medial	Initial
U+0647	ARABIC LETTER HEH	٥	4	8	ھ
U+06BE	ARABIC LETTER HEH DOACHASHMEE	ھ	B	8	ھ
U+06C1	ARABIC LETTER HEH GOAL	٥	~	1 ∼	ų
U+06D5	ARABIC LETTER AE	٥	4	,	

Note that U+06D5 • differs from the others in that it is right-linking, not dual-linking. It is thus similar to a sequence *<heh*, *zwnj>*, although there is no such formal equivalence defined via a canonical (or even compatibility) decomposition.

3.1 Arabic

Arabic has just one *heh* character, encoded with the basic Arabic block code U+0647 . This character is sometimes written as A, particularly when used as part of the *abjad-hawaz* numbering system using Arabic letters; and in medial position it may appear as either 4 or codepending on the style of

script. Both forms may occur even within the same line of text, at the discretion of the calligrapher. Nevertheless, these are not regarded as different characters, but merely as glyph variants of the same character, U+0647.

There is also the feminine ending written as \ddot{s} , a form based on *heh* despite its name TEH MARBUTA. This has its own Unicode value, U+0629. Although it could be considered to be a form of *teh*, this is not reflected in the encoding, as the written form is different, and Unicode encodes the script, not the language.

- use U+0647 for *heh*
- other *heh*-like characters should not occur

3.2 Persian

Persian uses the same letter *heh* as Arabic, and has no reason to encode it differently.

A uniquely Persian usage is the form 5, a letter *heh* with *hamza* above, used in compound words (*izafet*). There is a Unicode character U+06C0 5 that has this appearance. However, this is defined as canonically equivalent to the sequence <U+06D5 0, U+0654 5, and *not* to a sequence involving the normal *heh*, U+0647 0. As Persian is considered to have only one *heh*, the ISIRI 6219 standard specifies that U+06C0 5 is *not* to be used in Persian text; rather, *izafet* is encoded as <U+0647 0, U+0654 5, with a zero width non-joiner if necessary to prevent linking to a following letter.

- use U+0647 for *heh*
- use <U+0647 •, U+0654 •, U+200C zero width non-joiner> for izafet
- other *heh*-like characters should not occur
- in particular, do not use U+06C0 5 for izafet

3.3 Urdu

In Urdu writing, different forms of *heh* are used to show the distinction between an /h/ consonant and aspiration of a plosive. Aspiration is written with the double-looped *doachashmee* form (sometimes also referred to as a *butterfly* or *knotted heh*), while the separate /h/ consonant is most commonly written with non-looped forms such as (A) As this semantic distinction is required as part of the Urdu writing system, two separate *heh* characters must be encoded.

For aspiration, Unicode provides U+06BE , encoded for this exact purpose.

The /h/ consonant is more problematic. The Unicode character U+06C1 oprovides by default the typical shapes expected for Urdu heh, and was apparently added to the standard for this purpose. (The name HEH GOAL is a rough transcription of the Urdu expression (round heh), indicating its origin as a character intended for Urdu usage.) However, the Urdu heh is not really a different letter from the Arabic or Persian one; it is the same letter. It is merely written (usually) with a different choice among the possible glyph variants of that letter, the "Urdu shapes" being more typical of Nastaliq script, the calligraphic style normally used for Urdu. For consistency across the languages that use Arabic script, therefore, it is preferable to use the standard heh character, U+0647 o, and treat the preferred forms for Urdu as a language-specific font variant.

- use U+06BE ه to form aspirate digraphs (ده , یه , etc.)
- use U+0647 o for heh, with joining glyphs similar to those seen for U+06C1:
- do not use U+06C1 o, which is a glyph variant of U+0647
- treat U+06C1 as equivalent to U+0647 if found in Urdu data

3.4 Sindhi

Like Urdu, Sindhi also uses a *doachashmee* form of *heh* to write some aspirated consonants (specifically and specifically); there are unique letters for other aspirates). Unlike Urdu, the normal /h/ consonant is also commonly written with a form; however, some writers, at least, make a distinction between the letter *heh* /h/ and aspiration by writing the medial form of *heh* as while using for the aspirate digraphs. Thus, (initial aspirated /jh/) contrasts with (initial /j/ followed by medial /h/). Some writers of Sindhi script also prefer to use double-looped forms a rather than a 4 in isolate and (more rarely) final positions.

- use U+06BE م to form aspirate digraphs جم and على and
- use U+0647 o for heh, preferring double-looped glyphs in most cases

3.5 Parkari

The Parkari language (and other related languages in the southern Sindh) uses a writing system based on Sindhi (see above); however, it has an additional *heh*-based letter, encoded as U+06FF â. Font designers should note that the linked glyphs associated with U+06FF should all be based on a double-looped *heh*, thus: â â â. For consistency, it is common to use the forms a for the normal *heh*, U+0647.

Parkari also uses a final *heh* with the form , when writing certain vowel sounds. This can be encoded using U+06C1 • *heh goal*, which is an appropriate choice here as this shape needs to contrast with the normal *heh*.

- use U+06BE م to form aspirate digraphs جم and علي and
- use U+0647 for heh, preferring double-looped glyphs • a •
- use U+06C1 in final position for the vowel written as 2 < U+0626, U+06C1>
- Parkari also uses U+06FF 🏝

3.6 Kurdish

In Kurdish, there is a normal /h/ consonant written with the double-looped *heh* , and in addition the /e/ vowel is written with the *heh*-based forms 4.6 (similar to Parkari usage of the ... form). In Sorani, probably the most widely-written Kurdish dialect, /h/ does not occur in final position, and therefore there is no confusion between /h/ and this /e/ vowel, but in the Behdini dialect, it is important that final /h/ is distinguished from /e/. For this reason, writers use a double-looped final form of *heh* for the /h/ consonant.

Some implementations appear to have encouraged the use of U+06BE a for *heh* in Kurdish, because typical glyph shapes used for this character more nearly correspond to Kurdish preferences. However, this is not recommended, as it obscures the identity of the letter as the "normal" *heh* of Arabic script, and the shape of the final form of U+06BE (typically 4) is still not in accordance with Behdini usage. It is therefore recommended that Kurdish *heh* should be encoded as U+0647 4, just as for other languages, with the preferred glyph shapes being provided as a language-specific font variant.

- use U+06D5 of or the final /e/ vowel written as 45 < U+0626, U+06D5 > 0
- do not use U+06BE for *heh*, but recognize that some data may include it

4 YEH-based letters

The letter *yeh*, like other Arabic letters, was originally written without dots, ω . Over time, the dotted form ω developed, representing the same letter. There is also an alternate form ω seen in some styles of writing. (Note that although the use of dots under the isolated and final forms of *yeh* may be optional, they are not optional in modern writing under the initial ω and medial ω forms.

The yeh shape & is also used, optionally with a superscript alef (encoded as U+0670 $\dot{\circ}$), to write a final alef in some instances. This usage, alef maksura, omits the dots and is regarded as a different letter than normal yeh.

4.1 Arabic

Although it may be written in various ways, Arabic has just one *yeh* letter and should encode it consistently using U+064A . Although the representative glyph shown for this Unicode character has dots below, it may be rendered as a dotless shape in some styles and locales (this would be expected in a Nastaliq font, for example); it could even appear as a shape in appropriate calligraphic contexts.

- use U+064A \sim for yeh
- use U+0649 s for alef maksura (only occurs word finally)
- do not use U+06CC \searrow or U+06D2 \nearrow for yeh

4.2 Persian

In Persian, *yeh* is consistently written without dots, \mathcal{L} , in isolated and final positions. The dotted form that is the typical rendering of U+064A \mathcal{L} would be considered incorrect. To allow the Persian form \mathcal{L} to be reliably rendered, and to co-occur with the typical Arabic \mathcal{L} in the same text, Unicode provides a separate code U+06CC \mathcal{L} FARSI YEH. Although the Persian *yeh* is clearly related to the Arabic letter, and could have been considered a glyph variant of it, the difference in typical appearance is not an optional feature; it is a requirement for proper rendering of Persian text.

Note that it is possible to produce the appropriate *appearance* for Persian by the use of Arabic yeh U+064A in initial and medial positions, and alef maksura U+0649 in final and isolated positions. This represents an inconsistent encoding of the Persian letter, and should not be done; however, users working with systems designed for Arabic may have used such "hacks" to achieve the desired rendering.

There have also been cases where Arabic fonts have been altered to render U+064A with a dotless glyph, and U+0649 with a dotted one (in order to support the Persian preference for *yeh*, while still providing the possibility of rendering the Arabic form where needed). This swapping of glyphs represents a deviation from the Unicode standard, and leads to data interchange problems; it is not a correct way to encode Persian text.

When Arabic words spelled with *alef maksura* are used in Persian, it could be considered most logical to encode these with U+0649 \odot , for consistency with Arabic. However, as U+0649 \odot is indistinguishable from the Persian *yeh* U+06CC \odot in word-final position (the only position *alef maksura* should occur, it is unlikely that users will make such a distinction.

- use U+06CC \sim for yeh
- be aware that some Persian text may be encoded with U+064A and/or U+0649 for yeh

4.3 Urdu

In Urdu, a distinction is made between the form على, representing an /i/ vowel, and رمے, representing /e/. The two forms are known as بحولُ ہے small yeh and بری لے large yeh respectively, and are considered separate letters. Thus, unlike in Arabic, the form ملك must be encoded separately from على, not treated as an optional calligraphic glyph variant.

For $\frac{2}{2}$, $\frac{2}{2$

- use U+06CC of for small yeh/i/
- use U+06D2 \(\sum \) for large yeh /e/
- be aware that U+064A and/or U+0649 may also be found

4.4 Sindhi

Although Sindhi is a neighboring language to Urdu, it follows the Arabic convention of writing *yeh* with dots, \mathcal{L}_{s} , rather than in the Persian and Urdu way. A dotless form is only seen when writing Arabic-derived words with *alef maksura*, not for the normal *yeh* letter. It therefore uses the same character codes as Arabic for *yeh*.

- use U+064A \sim for yeh
- use U+0649 & for alef maksura in Arabic-derived words

4.5 Kurdish

Kurdish practice appears to be to write *yeh* in the Persian style, without dots in final/isolate position, indicating that U+06CC ω is the preferred character. In addition, there is a version of *yeh* with a 'small v' mark above, encoded as U+06CE ω . Like the regular *yeh*, this character also acquires two dots below when it links to a following letter, so the joining forms are ω .

It should be noted that the 'small v' here is an integral part of the letter, not a vowel mark akin to *fatha*, etc. The 'small v' is analogous to the dots and other marks used to create new letters as modifications of a basic Arabic letter shape. Therefore, it would be inappropriate to represent this Kurdish letter as U+06CC \swarrow *yeh* with the vowel mark U+065A $^{\sim}$ added.

I have also seen Kurdish texts that use a *yeh* with a horizontal bar above, \mathcal{Z} , but do not have adequate information concerning this at present.

- use U+06CC \sim for yeh
- use U+06CE & for yeh with small v
- do not use <U+06CC کی U+065A Č> as alternative to U+06CE

4.6 Uighur

The Uighur writing system includes several unusual innovations. In addition to a standard Arabic yeh, written with dots as ς , it uses a form with two dots placed vertically below ς , and a dotless form ς . Unlike alef maksura in Arabic, the dotless form also occurs within words (not only in final position); and unlike the Persian yeh, it does not acquire dots when linked (as it must be distinguished from the normal yeh). Thus, there are three distinct yeh-shaped letters, each with a full set of joining forms.

The dotless and vertical-dotted *yeh* letters are typically preceded with a *hamza* when they occur in initial or final position, or at a syllable break, and may be shown in this way in sources that discuss the writing system. For example, the vertical-dotted *yeh* may be shown as رئب ہو بنی ہی This is, however, a spelling convention and not an inherent aspect of the shaping of the letter.

- use U+0649 ك (joining forms (د مري) for the dotless yeh
- use U+064A ي (joining forms (پيټي) for the dotted yeh
- use U+06D0 ي (joining forms (ب به ي) for the yeh with vertical dots
- insert U+0626 في where a form with preceding *hamza* is required, e.g., ئىدى is encoded <U+0626 في U+0649 كى, U+0649 كى, U+0649 كى)
- do not use U+06CC (Persian yeh) for the dotless yeh

5 For font designers: summary of *heh* glyph variants

This document recommends that the "normal" Arabic *heh* character U+0647 • should be rendered with different glyph shapes to match the requirements and expectations of particular languages and locales. To provide a convenient reference for developers of general-purpose Unicode fonts that cover the Arabic block and are intended to be widely useful, some appropriate language-specific renderings are summarized here.

The glyphs shown here are a typical "simplified Naskh" style, which forms the basis for most general-purpose fonts. In other styles such as Nastaliq, Ruqa, Kufi, etc., there may be significant differences to consider. This summary does not claim to be a substitute for in-depth knowledge of the particular writing traditions of different user communities; it is merely a guidepost intended to promote better consistency in type design and encoding practices.

Typical default shapes for			Final	Medial	Initial
U+0647	ARABIC LETTER HEH	٥	4	or ~	ھ
U+06BE	DOACHASHMEE	ھ	8	8	ھ
U+06C1	GOAL	٥	^	4~	٦
U+06FF	WITH INVERTED SMALL V ABOVE	ۿ	ۿ	ۿ	ۿ
<i>Urdu</i> U+0647	ARABIC LETTER HEH	٥	^	4-	ન
Sindhi U+0647	ARABIC LETTER HEH	ھ	ه or ۵	▲ or ₄	ھ
<i>Parkari</i> U+0647	ARABIC LETTER HEH	æ	۵	ھ	ھ
Kurdish U+0647	ARABIC LETTER HEH	ھ	A	€	ھ