# Appendix A: Character Shaping Behavior

Characters from the Arabic and Indic scripts can combine or change shape depending on their context. A character's appearance is affected by its ordering with respect to other characters, the font used to render the character, and the application or system environment. These variables can cause the appearance of such characters to be different from their nominal glyphs (such as those shown in the Unicode charts).

In the Indic scripts, the presence of certain characters causes a change in the display order of other characters. This reordering is not commonly seen in non-Indic scripts. It also occurs independently of any bidirectional character reordering that might be in effect.

The following section describes minimal requirements for legibly rendering Arabic, Devanagari, and Tamil scripts as part of a plain text sequence. It describes the mapping between Unicode values and the glyphs in a font used to represent these characters. This section also describes the combining and ordering of those glyphs. The reader should be familiar with the General Principles of the Unicode Standard (Volume 1, Chapter 2), and especially with the distinction between characters and glyphs.

The algorithms provide minimal requirements for legibly rendering interchanged Unicode text with characters from these scripts. As with any script, a more complex algorithm can add rendering characteristics, depending on the font and application.

It is important to emphasize that in a font that is capable of rendering Arabic and Indic scripts, the set of glyphs is much greater than the number of Unicode characters. If a particular form is not found in the Unicode charts, implementors should be aware that it may be a different form of a Unicode character, or the result of a ligature of two or more Unicode characters.

For brevity, the Indic scripts are represented in some detail with Devanagari and Tamil, since they illustrate the range of types of behavior found in the other Indic scripts. The behavior of the other scripts can be derived from a close examination of these examples.

## **Arabic Character Shaping**

Arabic is always written in "cursive" or "handwritten" form, where the characters link together as they are written. All cursive scripts require linking rules for proper rendering. The linking rules for Arabic may be easier to understand if compared to Latin cursive handwriting.

There are also many optional ligatures which may be found in different fonts, as well as special forms of characters such as the *snake-kaf* or *extended noon* which may be used for justification or æsthetic reasons.

The minimal shaping algorithm is described in terms of the display (visual) order. In other words, the position of characters in the included examples is presented as it would appear on the screen after the bidirectional algorithm has determined character order in a line of text. Wherever the rules do not specify a change in form, characters are presented with their nominal forms, as they are presented in the code charts.

### Linking Classes

For brevity, the words ARABIC LETTER have been omitted from the names of all Arabic characters. For the purposes of describing Arabic linking, Unicode characters fall into the following classes. A complete list follows at the end of this section.

Dual-linking:

such as BAA, TAA, THAA, JEEM...

Right-linking:

ALEF, DAL, THAL, RA, ZAIN...

Link-causing:

ZERO WIDTH JOINER, TATWEEL (kashida)

Non-linking:

ZERO WIDTH NON-JOINER, and all spacing characters (other than those above) are non-linking, including: HAMZAH, HIGH HAMZA,

spaces, digits, punctuation, non-Arabic letters, and so on.

Transparent:

All combining marks and format marks are transparent, including:

FATHATAN, DAMMATAN, FATHAH, DAMMAH, KASRAH, SHADDAH,

SUKUN, ALEF ABOVE, RIGHT-LEFT MARK and so on.

For convenience in the linking rules, two supersets of some of the preceding classes are defined as follows: A *right link-causing character* is either a right-linking, dual-linking or link-causing character. A *left link-causing character* is either a left-linking, dual-linking or link-causing character.

#### Notation

If X refers to a character, then various glyph types representing that character are referred to by adding a subscript:

X <sub>n</sub>	Nominal form of a character, as it appears in the code charts (subject to stylistic variation). Unless otherwise stated, characters will adopt a nominal form.
$X_r$	Form of the character which joins to the right (right-linking and dual-linking characters may have this form).
× <sub>l</sub>	Form of the character which joins to the left (left-linking and dual-linking characters may have this form).
X <sub>m</sub>	Form of the character which joins to both the right and the left (dual-linking characters may have this form).
$X.Y_n$	Ligature form representing a combination of a form of $X_n$ and a form of $Y_n$ .

### Cursive Linking Rules

J1. Transparent glyphs do not affect the linking behavior of base (spacing) characters. For example:

J2. A right-linking glyph  $X_n$  which has a right link-causing glyph on the right will adopt the form  $X_r$ . For example:

- I3. A left-linking glyph  $X_n$  which has a left link-causing glyph on the left will adopt the form  $X_l$ .
- J4. A dual-linking glyph  $X_n$  which has a right link-causing glyph on the right and a left link-causing glyph on the left will adopt the form  $X_m$ . For example:

TATWEEL
$$_n$$
 + MEEM $_n$  + TATWEEL $_n$  → TATWEEL $_n$  + MEEM $_m$  + TATWEEL $_n$ 

J5. A dual-linking glyph  $X_n$  which has a right link-causing glyph on the right and no left link-causing glyph on the left will adopt the form  $X_r$ . For example:

$$\mathsf{MEEM}_n + \mathsf{TATWEEL}_n \to \mathsf{MEEM}_r + \mathsf{TATWEEL}_n$$

$$+ - \to + - \to +$$

J6. A dual-linking glyph which has a left link-causing glyph on the left and no right-link causing glyph on the right will adopt the form  $X_1$ . For example:

As noted in Section 3.1 of Volume 1, "Arabic," the zero width non-joiner can be used to prevent linking, as in the Persian plural suffix or Ottoman Turkish vowels.

## Ligature Classes

A certain class of ligatures are obligatory regardless of font design: they must be used in well-rendered Arabic text. There are many other ligatures that are optional, depending on font design. Since they are optional, those ligatures are not covered by the Unicode standard.

For the purpose of describing the obligatory Arabic ligatures, Unicode characters fall into the following classes (for a complete list, see the end of this section):

## **Obligatory Ligature Rules**

The following rules describe the formation of ligatures. They are applied after the preceding linking rules.

L1. Transparent glyphs do not affect the ligating behavior of base glyph. For example:

$$\mathsf{ALEF}_r + \mathsf{FATHAH}_n + \mathsf{LAM}_l \rightarrow \mathsf{LAM}.\mathsf{ALEF}_n + \mathsf{FATHAH}_n$$

(However, placement of a combining mark on a resulting ligature will depend on which of the original base characters the mark modifies).

L2. Any sequence with  $ALEF_r$  on the left and  $LAM_m$  on the right will form the ligature  $LAM.ALEF_l$ .

$$l + l \rightarrow \lambda$$
 (not u)

L3. Any sequence with  $ALEF_r$  on the left and  $LAM_l$  on the right will form the ligature  $LAM.ALEF_n$ .

$$l + J \rightarrow Y$$
 (not U)

### **Optional Features**

There are many other ligatures and contextual forms that are optional—depending on the font and application—such as the following:

In addition, the context-sensitive placement of non-spacing vowels such as fatha can greatly improve the appearance of the text.

### Character Types

The following is a detailed list of the Arabic characters that are either right-linking or dual-linking. All other Arabic characters (aside from TATWEEL) are non-linking, including ARABIC LETTER AE. For brevity in the names, the words Two, THREE and FOUR are abbreviated.

Most of the extended Arabic characters are merely variations on the basic Arabic shapes, with additional or different marks. For compatibility, many precomposed forms are included.

In some cases there are characters that only occur at the end of words in correct spelling, called trailing characters. Examples include teh marbuta, alef maqsurah and dammatan. When trailing characters are joining (such as teh marbuta), they are classed as right-joining (even when similarly-shaped characters are dual-joining). When trailing characters do not join or cause joining (such as dammatan), they are classed as transparent, which treats all combining marks similarly.

Note: In the case of HA, HA<sub>l</sub> is also shown in the code chart box. This is often done to reduce the chance of misidentifying HA as ARABIC INDIC DIGIT 5, which has a very similar shape. The nominal form of HA is the isolate form, which looks like ARABIC LETTER AE. In the case of HA GOAL, the nominal form is not even listed, and also resembles ARABIC LETTER AE.

The characters in the following charts are grouped by shape, and not by standard Arabic alphabetical order.

# Dual-linking

Group	CHAR <sub>n</sub>	CHAR <sub>r</sub>	CHAR m	CHAR <sub>1</sub>	Characters with Similar Shaping Behavior	
BAA	, Marketa A	, than al	stor grill	eg-e-man	TAA, THAA, TAA WITH SMALL TAH, TAA WITH 2 DOTS VERTICAL ABOVE, BAA WITH 2 DOTS VERTICAL BELOW, TAA WITH RING, TAA WITH 3 DOTS ABOVE DOWNWARD, TAA WITH 3 DOTS BELOW, TAA WITH 4 DOTS ABOVE, BAA WITH 4 DOTS BELOW	
наа	۲	ح	×	en effet tra	JEEM, KHAA, HAMZAH ON HAA, HAA WITH 2 DOTS VERTICAL ABOVE, HAA WITH MIDDLE 2 DOTS, HAA WITH MIDDLE 2 DOTS VERTICA. HAA WITH 3 DOTS ABOVE, HAA WITH MIDDLE 3 DOTS DOWNWARD, HAA WITH MIDDLE 4 DOTS	
SEEN	س	ا ۱۸۹۸ ا	tr più na	erit ura w	SHEEN, SEEN WITH DOT BELOW AND DOT ABOVE, SEEN WITH 3 DOTS BELOW, SEEN WITH 3 DOTS BELOW AND 3 DOTS ABOVE	
SAD	ص	ص	-2	ص	DAD, SAD WITH 2 DOTS BELOW, SAD WITH 3 DOTS ABOVE	
TAH	ط	ط	<u>ط</u>	ط	dhah, tah with 3 dots above	
AIN	٤	ځ		ء	GHAIN, AIN WITH 3 DOTS ABOVE	
FA	ف	ن	ف	ۏ	DOTLESS FA, FA WITH DOT MOVED BELOW, FA WITH DOT BELOW, FA WITH 3 DOTS ABOVE, FA WITH 4 DOTS ABOVE	
QAF	ق	ق	ä	و ق ماللا	QAF WITH DOT ABOVE, QAF WITH 3 DOTS ABOVE	
CAF	ك	십	ج	5	CAF WITH DOT ABOVE, CAF WITH 3 DOTS ABOVE, CAF WITH 3 DOTS BELOW	
LAM	J	ل	ı	j	LAM WITH SMALL V, LAM WITH DOT ABOVE, LAM WITH 3 DOTS ABOVE	
МЕЕМ	م	۴	•	م		
NOON	ن	ڹ	÷	3	DOTLESS NOON, DOTLESS NOON WITH SMALL TAH, NOON WITH RING, NOON WITH 3 DOTS ABOVE	
на	ه	٨	4  r	A		

# Dual-linking (continued)

YA	ي	ي	KI + DEE	ANTE ME	HIGH HAMZAH YA, HAMZAH ON YA, DOTLESS YEH, YA WITH SMALL V, YA WITH 2 DOTS VERTICAL BELOW, YA WITH 3 DOTS BELOW
SWASH CAF	5	ڪ	ڪ	5	WIND WORTH AS A CONTROL OF THE CONTR
GAF	گ	گ	٤	\$	OPEN CAF, CAF WITH RING, GAF WITH RING, GAF WITH 2 DOTS ABOVE, GAF WITH 2 DOTS BELOW, GAF WITH 3 DOTS ABOVE
KNOTTED HA	ھ	a	a	Ą	A SECURIAL STATES STATE
HA GOAL	A	~	Ŧ	ţ	HAMZAH ON HA GOAL

# Right-linking

Group	CHARn	CHAR <sub>T</sub>	Characters with Similar Shaping Behavior
ALEF	1		ALEF, HAMZAH ON ALEF, MADDAH ON ALF, HAMZAH UNDER ALEF, WAVY HAMZAH ON ALEF, HIGH HAMZAH ALEF
DAL	3	٨	THAL, DAL WITH SMALL TAH, DALL WITH RING, DALL WITH DOT BELOW, DAL WITH DOT BELOW AND SMALL TAH, DAL WITH 2 DOTS, DAL WITH 2 DOTS BELOW, DALL WITH 3 DOTS ABOVE, DALL WITH 3 DOTS ABOVE DOWNWARD, DAL WITH 4 DOTS ABOVE
RA	ر	in twell	ZAIN, RA WITH SMALL TAH, RA WITH SMALL V, RA WITH RING, RA WITH DOT BELOW, RA WITH SAMM V BELOW, RA WITH DOT BELOW AND DOT ABOVE, RA WITH 2 DOTS ABOVE, RA WITH 3 DOTS ABOVE, RA WITH 4 DOTS ABOVE
WAW	و		HAMZAH ON WAW, HIGH HAMZAH WAW, HIGH HANZAH WAW WITH DAMMAH, WAW WITH RING, WAW WITH BAR, WAW WITH SMALL V, WAW WITH DAMMAH, WAW WITH ALEF ABOVE, WAW WITH INVERTED SMALL V, WAW WITH 2 DOTS ABOVE, WAW WITH 3 DOTS ABOVE
ALEF MAQSURAH	ی	ی	YA WITH TAIL
TAA MARBUTAH	ä	ä	
TAA MARBUTAH GOAL	ة	ベ	HAMZAH ON HA
YA BAREE	_	_	HAMZAH ON YA BARREE

#### Other Classes

Link-causing	ZERO WIDTH JOINER, TATWEEL
Non-linking	ZERO WIDTH NON-JOINER, and all spacing characters (other than those mentioned above) are non-linking, including: HAMZAH, HIGH HAMZAH, HAMZAT WASL ON ALEF, spaces, digits, punctuation, non-Arabic letters and so on.
Transparent	All combining marks and format marks are transparent, including: fathatan, dammatan, fathah, dammah, kasrah, shaddah, sukun, alef above, rightleft mark and so on.

		•					Blanca.
Ind	ox	hv	112	1100	de	Val	110
Titte		$\nu_{y}$	$\mathbf{c}$	***	ne	1 000	ne

Unic.	Link	Link Group	Name	
0622		ALEE	1887-181 (1886)(18 1.1-4.5) (1948)	
0622	R	ALEF	MADDAH ON ALEF	
0623	R	ALEF	HAMZAH ON ALEF	
0624	R	WAW	HAMZAH ON WAW	
0625	R	ALEF	HAMZAH UNDER ALEF	
0626	D	YA Jin Assentian s	HAMZAH ON YA	
0627	R	ALEF	ALEF AND MADE AND MAD	
0628	D	BAA	BAA JIII III III III III III III III III	
0629	R	TAA MARBUTAH	TAA MARBUTAH	
062A	D	BAA	TAA	
062B	D	BAA	THAA	
)62C	D	HAA	JEEM	
062D	D	HAA	HAA	
062E	D	HAA	KHAA	
062F	R	DAL	DAL	
0630	R	DAL	THAL	
0631	R	RA	RA	
0632	R	RA	ZAIN	
0633	D	SEEN	SEEN	
0634	D	SEEN	SHEEN	
0635	D	SAD	SAD	
0636	D	SAD	DAD	
637	D	TAH	TAH	
0638	D	TAH	DHAH	
0639	D	AIN	AIN	
063A	D	AIN	GHAIN	
0640	C	<no shaping=""></no>	TATWEEL	
0641	D	FA	FA	
0642	D	QAF	QAF	
0643	D	CAF	CAF	
)644	D	LAM	LAM	
645	D	MEEM	MEEM	
)646	D	NOON	NOON	
)647	D	НА	НА	
0648	R	WAW	WAW	
0649	R	ALEF MAQSURAH	ALEF MAQSURAH	
)64A	D	YA	YA	
0671	U	<no shaping=""></no>	HAMZAT WASL ON ALEF	
0672	R	ALEF	WAVY HAMZAH ON ALEF	
0673	R	ALEF	WAVY HAMZAH UNDER ALEF	
674	U	<no shaping=""></no>	HIGH HAMZAH	
675	R	ALEF	HIGH HAMZAH ALEF	
0676	R	WAW	HIGH HAMZAH WAW	
0677	R	WAW	HIGH HAMZAH WAW WITH DAMMA	Н
0678	D	YA	HIGH HAMZAH YA	M.17
0679	D	BAA	TAA WITH SMALL TAH	
067A	D	BAA	TAA WITH 3 MALE TAN TAA WITH 2 DOTS VERTICAL ABOVE	
067B	11-01-0			
J0/B	D	BAA	BAA WITH 2 DOTS VERTICAL BELOW	

067C	D	BAA Would a croduct	TAA WITH RING
067D	D	BAA	TAA WITH 3 DOTS ABOVE DOWNWARD
067E	D	BAA	TAA WITH 3 DOTS BELOW
067F	D	BAA SAADAA BARATA H	TAA WITH 4 DOTS ABOVE
0680	D	BAA WOUNDERFOORER	BAA WITH 4 DOTS BELOW
0681	D	HAA alla suo malla madis el	HAMZAH ON HAA
0682	D	HAA JVORA STOCK H	HAA WITH 2 DOTS VERTICAL ABOVE
0683	D	HAA	HAA WITH MIDDLE 2 DOTS
0684	D	HAA	HAA WITH MIDDLE 2 DOTS VERTICAL
0685	D	HAA BUDDA TO DE SO	HAA WITH 3 DOTS ABOVE
0686	D	HAA	HAA WITH MIDDLE 3 DOTS DOWNWARD
0687	D	HAA IAS IIIAME HYW MOON	HAA WITH MIDDLE 4 DOTS
0688	R	DAL and and	DAL WITH SMALL TAH
0689	R	DAL avena atom entitle	DAL WITH RING
068A	R	DAL	DAL WITH DOT BELOW
068B	R	DAL	DAL WITH DOT BELOW AND SMALL TAH
068C	R	DAL	DAL WITH 2 DOTS ABOVE
068D	R	DAL JACK AS MOS	DAL WITH 2 DOTS BELOW
068E	R	DAL JACO MATURA	DAL WITH 3 DOTS ABOVE
068F	R	DAL	DAL WITH 3 DOTS ABOVE DOWNWARD
0690	R	DAL	DAL WITH 4 DOTS ABOVE
0691	R	RA - VIII WHENT	RA WITH SMALL TAH
0692	R	RA	RA WITH SMALL V
0693	R	RA	RA WITH RING
0694	R	RA V LIMEN COMPANY INC.	RA WITH DOT BELOW
0695	R	RA WORL STOCK HT	RA WITH SMALL V BELOW
0696	R	RA COMPANIES OF HED	RA WITH DOT BELOW AND DOT ABOVE
0697	R	RA	RA WITH 2 DOTS ABOVE
0698	R	RA	RA WITH 3 DOTS ABOVE
0699	R	RA V MANY 2	RA WITH 4 DOTS ABOVE
069A	D	SEEN	SEEN WITH DOT BELOW AND DOT ABOVE
069B	D	SEEN	SEEN WITH 3 DOTS BELOW
069C	D	SEEN	SEEN WITH 3 DOTS BELOW AND 3 DOTS ABOVE
069D	D	SAD	SAD WITH 2 DOTS BELOW
069E	D	SAD	SAD WITH 3 DOTS ABOVE
069F	D	TAH	TAH WITH 3 DOTS ABOVE
06A0	D	AIN	AIN WITH 3 DOTS ABOVE
06A1	D	FA	DOTLESS FA
06A2	D	FA	FA WITH DOT MOVED BELOW
06A3	D	FA	FA WITH DOT BELOW
06A4	D	FA	FA WITH 3 DOTS ABOVE
06A5	D	FA	FA WITH 3 DOTS BELOW
06A6	D	FA	FA WITH 4 DOTS ABOVE
06A7	D	QAF	QAF WITH DOT ABOVE
06A8	D	QAF	QAF WITH 3 DOTS ABOVE
06A9	D	GAF	OPEN CAF
06AA	D	SWASH CAF	SWASH CAF
06AB	D	GAF	CAF WITH RING
06AC	D	CAF	CAF WITH DOT ABOVE
06AD	D	CAF	CAF WITH 3 DOTS ABOVE

06AE	D	CAF	CAF WITH 3 DOTS BELOW	
06AF	D	GAF	GAF	
06B0	D	GAF	GAF WITH RING	
06B1	D	GAF	GAF WITH 2 DOTS ABOVE	
06B2	D	GAF	GAF WITH 2 DOTS BELOW	
06B3	D	GAF	GAF WITH 2 DOTS VERTICAL BELOW	1
06B4	D	GAF	GAF WITH 3 DOTS ABOVE	
06B5	D	LAM	LAM WITH SMALL V	
06B6	D	LAM	LAM WITH DOT ABOVE	
06B7	D	LAM	LAM WITH 3 DOTS ABOVE	
06BA	D	NOON	DOTLESS NOON	
06BB	D	NOON	DOTLESS NOON WITH SMALL TAH	
06BC	D	NOON	NOON WITH RING	
06BD	D	NOON	NOON WITH 3 DOTS ABOVE	
06BE	D	KNOTTED HA	KNOTTED HA	
06C0	R	TAA MARBUTAH	HAMZAH ON HA	
06C1	D	HA GOAL	HA GOAL	
06C2	R	HAMZAH ON HA GOAL	HAMZAH ON HA GOAL	
06C3	R	HAMZAH ON HA GOAL	TAA MARBUTAH GOAL	
06C4	R	WAW	WAW WITH RING	
06C5	R	WAW	WAW WITH BAR	
06C6	R	WAW	WAW WITH SMALL V	
06C7	R	WAW	WAW WITH DAMMAH	
06C8	R	WAW	WAW WITH ALEF ABOVE	
06C9	R	WAW	WAW WITH INVERTED SMALL V	
06CA	R	WAW	WAW WITH 2 DOTS ABOVE	
06CB	R	WAW	WAW WITH 3 DOTS ABOVE	
06CC	D	YA	DOTLESS YA	
06CD	R	ALEF MAQSURAH	YA WITH TAIL	
06CE	D	YA	YA WITH SMALL V	
06D0	D	YA	YA WITH 2 DOTS VERTICAL BELOW	
06D1	D	YA	YA WITH 3 DOTS BELOW	
06D2	R	YA BARREE	YA BARREE	
06D3	R	YA BARREE	HAMZAH ON YA BARREE	
06D5	U	<no shaping=""></no>	AE	

## Devanagari Character Shaping

Devanagari characters, like characters from many other scripts, can combine or change shape depending on their context. A character's appearance is affected by its ordering with respect to other characters, the font used to render the character, and the application or system environment. These variables can cause the appearance of Devanagari characters to be different from their nominal glyphs (such as those shown in the Unicode charts).

Additionally, a few Devanagari characters cause a change in the order of the displayed characters. This reordering is not commonly seen in non-Indic scripts and occurs independently of any bidirectional character reordering that might be required.

The following algorithm describes minimal rendering of Devanagari as part of a plain text sequence. It describes the mapping between Unicode values and the glyphs in a Devanagari font. It also describes the combining and ordering of those glyphs.

This algorithm provides minimal requirements for legibly rendering interchanged Unicode Devanagari text. As with any script, a more complex algorithm can add rendering characteristics, depending on the font and application.

It is important to emphasize that in a font that is capable of rendering Devanagari, the set of glyphs is greater than the number of Unicode Devanagari characters.

#### Notation

For brevity, the words "DEVANAGARI LETTER" and "DEVANAGARI VOWEL SIGN" are omitted from the names of Unicode characters in this section where such omission is not ambiguous. If X refers to a character, then various glyph types representing that character are referred to by adding a subscript.

$X_m$	Non-spacing mark form.
$x_l$	Left side form.
X <sub>n</sub>	Nominal form of a character, as it appears in the code charts (subject to stylistic variation). Unless otherwise stated, characters will adopt a nominal form.
$X_d$	Form of a consonant representing a "dead" form, where the implicit vowel is not pronounced.
$X_h$	Form of a dead consonant where the <i>virama</i> is absorbed, and a different shape is taken (see below). This is also known as a "half consonant."
$X.Y_n$	Ligature of a dead-consonant $XA_n$ with a consonant $Y_n$ . Such ligatures are also called <i>conjunct consonants</i> .

Vowel signs are indicated by superscripting with vs.

XVS

Vowel sign (not to be confused with the independent vowels).

AA is not the same as  $AA^{VS}$ .

#### Character Classes

The principal classes of Devanagari characters are the following (a complete list follows at the end of this section):

Consonants: KA, KHA, GA, GHA...

Vowels:

A. AA...

Vowel Signs:

AAVS, IVS...

### Background

Independent vowels can be followed by various members of the class of non-spacing marks. The marks are applied to the vowel in a manner similar to marks on Latin letters. Independent vowels generally have no other interaction with other characters.

Other well-formed Devanagari is organized as a sequence of syllable clusters. Phonetically, a syllable cluster is a sequence of consonants (usually one or two, occasionally three, rarely four), followed by a vowel, with perhaps some modification of the vowel (such as nasalization).

Other than the application of non-spacing marks, digits and free-standing signs have no interaction with other characters.

The *virama* indicates a null vowel. As such, it is used to join together the sequence of consonants in a syllable cluster. Additionally, it may be used at the end of the cluster if the cluster's vowel is null.

#### **Dead Consonants**

A dead consonant is defined as a sequence consisting of a consonant followed by a VIRAMA<sub>n</sub>.

The minimal rendering for a dead consonant is to position the *virama* as a non-spacing mark bound to the character. (See, however, Optional Forms in the following section.)

For example:

$$TA_n + VIRAMA_n \rightarrow TA_d$$

#### Consonant Clusters

A *consonant cluster* is defined as a consonant or independent vowel, optionally preceded by a sequence of one or more dead consonants, and optionally followed by a vowel sign, as in the following rules.

(Note that for the purposes of rendering, the independent vowels behave very much like consonants, even though they are rarely used this way. For example,  $\frac{1}{2}$  or  $\frac{3}{2}$  can be produced by an application of rules discussed below. See "Independent Vowels" for more information.)

For example:

C1. When  $RA_n$  occurs to the right of a VIRAMA<sub>n</sub>, then it combines with the VIRAMA<sub>n</sub> to form a low non-spacing mark  $(RA_m)$ , which attaches to the consonant on the left.

VIRAMA<sub>n</sub> + RA<sub>n</sub> → RA<sub>m</sub>

$$\bigcirc + ₹ → \bigcirc$$

For example:

THA<sub>n</sub> + VIRAMA<sub>n</sub> + RA<sub>n</sub> 
$$\rightarrow$$
 TTA<sub>n</sub> + RA<sub>m</sub>  $\rightarrow$  +  $\bigcirc$  +  $\bigcirc$  +  $\bigcirc$   $\rightarrow$   $\bigcirc$ 

C2. If the non-spacing  $RA_m$  is attached to a consonant having a vertical bar on its right, then it loses its right side, and is positioned attached to the vertical bar:

$$PHA_n + RA \rightarrow PHA_n + RA$$

C3. The half consonant form of  $RA_n$  behaves as a non-spacing mark when followed by a consonant. The vowel sign is then reordered (as in rule O1 in a following section). The ordering of these rules allows the production of ligatures such as 87.

C4. There are three forms that produce significantly different shapes when ligating. These forms are mandatory ligatures:

$$KA_d + SSHA_n \rightarrow K.SSHA_n$$
 $\overline{A}_{Q} + \overline{Q} \rightarrow G$ 
 $JA_d + NYA_n \rightarrow J.NYA_n$ 
 $\overline{Q} + \overline{A} \rightarrow \overline{Q}$ 
 $TA_d + RA_n \rightarrow T.RA_n$ 
 $\overline{Q} + \overline{Q} \rightarrow \overline{Q}$ 

C5. Non-spacing marks modifying a consonant are placed right after that consonant in the backing store, and are attached to that consonant in rendering. If a virama is included, it should be placed after the consonant-modifying non-spacing marks in the backing store.

$$KA_n + NUKTA_n + AA_n$$
  
क  $+ \circ + \mathbf{1} \rightarrow$ का

C6. Non-spacing marks modifying a cluster are placed right after the final element in the cluster in the backing store, and are attached to the rightmost free-standing element of the cluster. Note that the bindus follow vowel signs, but precede svaras in the backing store. The relative placement of these non-spacing marks is horizontal rather than vertical: the horizontal rendering order may vary according to typographical concerns.

$$KA_n + AA_n + CANDRABINDU_n$$
  
क  $+ T + \mathring{\Box} \rightarrow \mathring{a}\mathring{1}$ 

# **Optional Forms**

In most cases, the following ligatures are the preferred forms, although the unligated form can be used and understood. There are generally two shape changes that happen in forming ligatures of consonant clusters.

L1. If a preceding consonant contains a vertical bar on the right, then that bar is removed, the virama disappears, and the characters are joined:

$$GA_d + GHA_n \rightarrow G.GHA_n$$
 $\overline{\eta} + \overline{\mathbf{u}} \rightarrow \overline{\mathbf{v}}\overline{\mathbf{u}}$ 

In implementation, this form can be rendered either by making a single composed ligature glyph to represent the cluster, or just by replacing the dead consonant by a special "half-consonant" form, such as:

$$GA_d \rightarrow GA_h$$
 $T \rightarrow T$ 
 $NA_d \rightarrow NA_h$ 
 $T \rightarrow \overline{T}$ 

The choice of whether to use the half-consonant forms or to use a composite glyph is implementation-dependent. Depending on the characters involved, both can produce acceptable display. The choice depends on the extent to which the two consonants undergo further shape modifications in order to link properly.

Character shaping can be explicitly controlled by using the ZERO WIDTH NON-JOINER and ZERO WIDTH JOINER. As described in Volume 1, the ZERO WIDTH NON-JOINER and ZERO WIDTH JOINER can be used to request that the *virama* be explicitly rendered to prevent any automatic formation of a ligature or half-consonant. The ZERO WIDTH JOINER can be used to request that the virama be absorbed into the half consonant form, and prevent any further automatic formation of a ligature.

For example:

$$KA_n + VIRAMA_n + SSHA_n \rightarrow K.SSHA_n$$
  
 $\overline{\Phi}_1 + \overline{\mathbb{Q}}_1 + \overline{\mathbb{Q}}_1 \rightarrow \overline{\mathbb{Q}}_1$   
 $KA_n + VIRAMA_n + ZWNJ + SSHA_n \rightarrow KA_d + SSHA_n$   
 $\overline{\Phi}_1 + \overline{\mathbb{Q}}_1 + \overline{\mathbb{Q}}_1 \rightarrow \overline{\mathbb{Q}}_1$   
 $KA_n + ZWJ + VIRAMA_n + SSHA_n \rightarrow KA_n + SSHA_n$   
 $\overline{\Phi}_1 + \overline{\mathbb{Q}}_1 + \overline{\mathbb{Q}}_1 \rightarrow \overline{\mathbb{Q}}_1$   
 $\overline{\Phi}_1 + \overline{\mathbb{Q}}_1 + \overline{\mathbb{Q}}_1 \rightarrow \overline{\mathbb{Q}}_1$ 

L2. Where the vertical bar does not exist or cannot be detached, or where the distinctive portion of the following character does not provide a linkage point, the letters are joined vertically, and then generally reduced in size.

$$z + z \to z$$

L3. In some cases, non-spacing marks will also combine with a base consonant, either attaching at a non-standard location, or changing shape.

In minimal rendering there are only two cases,  $RA_n$  with  $U_{vs}^n$  and  $UU_{vs}^n$ .

$$RA_n + UU_n \to RUU_n$$

$$\overline{\zeta} + Q \to \overline{\xi}$$

L4. Whenever a ligature is formed, it can then combine with other glyphs, ligatures, and half-forms in an even more complicated ligature:

## Character Ordering

There are cases in Devanagari where characters are rendered in a different order from the backing store. The reordering rules affect two characters.

O1. Cluster 
$$+I_n \rightarrow I_n + Cluster$$

Reordering depends upon whether consonants within the cluster ligate. For example:

$$GA_d + GHA_n + I_n \rightarrow I_n + G.GHA_n$$
 $\eta + \mathbf{E} + \mathbf{\hat{\cap}} \rightarrow \overline{\eta}\mathbf{E}$ 
 $G.GHA_n + I_n \rightarrow I_n + G.GHA_n$ 
 $\overline{\mathbf{1}}\mathbf{E} + \mathbf{\hat{\cap}} \rightarrow \overline{\mathbf{1}}\mathbf{E}$ 

O2. 
$$RA_h + Cluster \rightarrow Cluster + RA_d$$

Example:

$$RA_h + G.GHA_n \rightarrow G.GHA_n + RA_d$$
  $^{\circ}$  +  $^{\circ}$   $^{\circ}$  +  $^{\circ}$   $^{\circ}$   $^{\circ}$   $^{\circ}$ 

Once conjunct consonants have been formed and reordering has taken place, more advanced contextual modifications of the character shapes can be used to create a better presentation. For example, the hook in the  $I^{vs}$  can be extended to intersect the last vertical bar in the following (after reordering) consonant cluster, or the  $RA_d$  can be aligned with the vertical bar.

These modifications are optional, but produce a better appearance.

## Independent Vowels

There is a one-to-one correspondence between the independent vowels and the dependent vowel signs. Except for A, independent vowels are sometimes represented by a form of the dependent vowels. This representation is sometimes found in schoolbooks. Independent forms can be represented by a sequence of the A followed by the dependent vowel. This does not imply that the corresponding independent and dependent vowels are equivalent backing stores: they are two separate (but related) backing stores and renderings.

The dependent vowels should not be represented by a combination of *virama* plus independent vowel in interchange:

#### Character Classes

The following character types are found in written Devanagari:

#### Digits

0966 **DIGIT ZERO** 

096F DIGIT NINE

## Various Free-Standing Signs

0903 SIGN VISARGA

093D SIGN AVAGRAHA

0950 OM

0964 DANDA

0965 DOUBLE DANDA

0970 ABBREVIATION SIGN

## Independent Vowels

0905 Α

0914 ΑU

**VOCALIC RR** 0960

0961 **VOCALIC LL** 

# Cluster-Modifying Non-Spacing Marks

#### Bindus

0901 SIGN CANDRABINDU

0902 SIGN ANUSVARA

### Svaras

0951 STRESS SIGN UDATTA

0952 STRESS SIGN ANUDATTA

0953 **GRAVE ACCENT** 

0954 ACUTE ACCENT

# Consonant-Modifying Non-Spacing Marks

093C SIGN NUKTA

## Diacritical Vowel Signs

0941 **VOWEL SIGN U** 

0948 **VOWEL SIGN AI** 

0962 VOWEL SIGN VOCALIC L

0963 **VOWEL SIGN VOCALIC LL** 

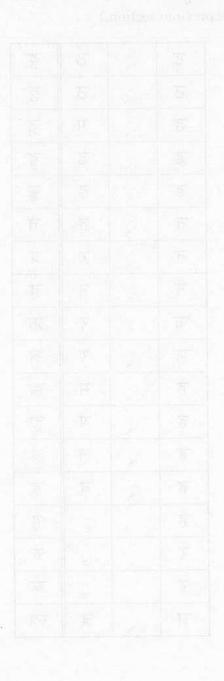
### Free-Standing Vowel Signs

093E **VOWEL SIGN AA** 

0940 **VOWEL SIGN II**  094A VOWEL SIGN SHORT O
094B VOWEL SIGN O
094C VOWEL SIGN AU

Preceding Type
093F VOWEL SIGN I

Virama 094D VIRAMA





# Sample Ligatures

The following charts show examples of ligatures that can be used in Devanagari. These ligatures are glyphs, not characters. (See L4 in the previous section.)

क	্	क	霰
क	্	त	क्त
क	্	र	क्र
क	্	ष	क्ष
ङ	্	क	ङ
ङ	্	ख	इं
ভ	Q	ग	ङ
ভ	Q	घ	ন্ধ
স	્	ज	ञ्ज
ज	Q	স	<b>ज्ञ</b>
द	Q	घ	द्ध
द	Q	द	द
द	্	ध	द्ध
द	્	ब	द्व
द	্	भ	ब्द
द	्	म	द्म
द	Q	य	द्य
द	্	व	द्व
ट	Q	ट	ट्ट

ट	્	ਰ	छ
ਰ	্	ਰ	ड
ड	Q	ग	ত্ত্
ड	্	ड	கு
ड	୍ଦ	ढ	क्रि
त	୍	त	त्त
त	oʻ	र	त्र
न	্	न	न्न
फ	Q	र	फ्र
श	્	र	श्र
虑	Q	म	展
ह	્	य	ह्य
ह	Q	ल	震
ह	্	ਕ	慮
ह		្វ	ह
र		ુ	रु
र		ૂ	रू
स	Q	त्र	स्त्र

Sample Half-Forms

क	Q	व	- ,न	Q.s	F
ख	Q	ख	प	Q	τ
ग	Q	nate for	फ	-Q	먁
घ	्	3	ब	্	9
च	୍ଦ	च	भ	a Q	£
ज	Q	ত	म	্	Ŧ
झ	Q	इ	य	Q	ट
স	Q	19	ल	Q	₹
ण	Q	σ	а	Q	5
त	Q	7	श	Q	য়
थ	Q	દ	ष	Q	7
ध	Q	3	स	a o	Æ

As explained in the text, ligatures may also have half-forms. The following shows examples of this:

Ligature Half-Forms

क	્	ष	Q	82
ज	্	ञ	Q	₹
त	्	त	্	₹
त	Q	र	Q	5
श	્	र	्	8

# Tamil Shaping Behavior

The South Indic scripts function in much the same way as Devanagari, with the additional feature of two-part vowels. The example of Tamil will be used to illustrate this feature. The following is a summary of the Tamil letters. As in the Devanagari example, the words "TAMIL LETTER" and "TAMIL VOWEL SIGN" will be omitted where this does not cause ambiguity. The latter are indicated by superscripting with 1/5.

It is important to emphasize that in a font that is capable of rendering Tamil, the set of glyphs is greater than the number of Unicode Tamil characters.

										وعظاهم	
க	囮	ச	ஜ	ஞ	<b>L</b>	ண	த	ந	ன	= LJ	1
KA	NGA	CA	JA	NYA	TTA	NNA	TA	NA	NNNA	PA	
ம	ш	Ţ	ற	ಉ	ள	ழ	ഖ	ஷ	സ	ஹ	
MA	YA	RA	RRA	LA	LLA	LLLA	VA	SSA	SA	НА	
அ	ஆ	<b>@</b>	ानः	உ	<u>গ্রু</u>	எ	ஏ	නු	99	ஓ	ஒள
Α	AA	Já Lá	11	U	UU	E	EE	Al	0	00	AU
	п	٦	కి	ղ <del>իս</del>	ු <del> </del> ම	େ	ෙ	ത	ொ	ோ	ெள
Α	AA	E	- 11	U	υυ	E	EE	AI	О	00	AU
	តា	80.8				11	-				4
VIRAMA	AU LENGTH		70								V5-0211

# Independent versus Dependent Vowels

As with Devanagari, the dependent vowel signs are not equivalent to a sequence of *virama* + independent vowel. For example:

As in the case of Devanagari, a consonant cluster is any sequence of one or more consonants separated by *viramas*, possibly terminated with a *virama*.

#### Two-Part Vowels

Certain Indic vowels consist of two discontiguous elements. As in other cases of discontiguous elements, there are two sequences of Unicode values which can be used to express equivalent spellings. This is similar to the case of letters such as "â," which can either be spelled with "a" followed by non-spacing "^", or spelled with a single Unicode character "â".

Note that the 6 I in the third example is not the LLA vowel; it is the AU LENGTH MARK.

If the precomposed forms are used in the backing store instead of the separate characters, then a similar transformation occurs in the rendering process. The precomposed form on the left is transformed into the two separate forms equivalent to those on the right, which are then subject to vowel reordering, as below. Thus in rendering:

## Vowel Reordering

The following vowels are always reordered in front of the previous consonant cluster, similar to the rendering behavior of the DEVANAGARI VOWEL SIGN I.

For example:

Backing sto	ore management		Display
க	<b>6</b>	$\rightarrow$	கெ
க	<b>6</b> - 6	$\rightarrow$ 700	கே
க	െ	100) 🗪 ji 🗀	கை

The same effect occurs with the results of vowel splitting, as follows:

Backi	ng store				Display
а	Б	ொ		+ V3.00) → 1	கொ
а	Б	ଚ	452 <b>U</b>	TORRE TRANS	கொ
а	Б	ோ		$\rightarrow$	கோ
д	5	<b>©</b>	H <sub>t</sub> = 1 may 1.1	* * 1.55 (11 to 1)	கோ
Э	5	ெள		itz galilaed <u>la</u> li n	கௌ
д	5 ed ou	බ	តា	$\rightarrow$	கௌ

In both cases, the ordering of the elements is *unambiguous*: the consonant (cluster) occurs *first* in the backing store. The vowel of also has two discontinuous parts, and can also be composed using the AU LENGTH MARK.

#### Ligatures

The following examples illustrate the range of ligatures available in Tamil. These changes take place after vowel reordering and vowel splitting. Unlike Devanagari, there are very few conjunct consonants; most ligatures are located between a vowel and a neighboring consonant.

1. Conjunct consonants.

As with Devanagari, vowel reordering occurs around conjunct consonants. For example:

2. The vowel ∏ optionally ligates with 600T, 60T, or ဤ on its left:

Since this process takes place after reordering and splitting, the following ligatures may also occur:

Precomposed Vowels  $mathref{m} \operatorname{Precomposed} \operatorname{Vowels}$   $m + \operatorname{Old} \operatorname{Ol$ 

3. The vowel signs and form ligatures with ∟ on their left.

$$\Gamma + \begin{subarray}{c} \to \begin{subarray}$$

These vowels often change shape or position slightly in order to link up with the appropriate shape of the consonant on their left:

ର + 
$$^{\circ}$$
 → ଚୌ

4. The vowel signs and typically change form or ligate:

×	x + ी	× + ೃ •೨
க	கு	<del>መ</del>
ங	ାହା	ஙூ
Ф	퓩	<b></b>
ஞ	து	ஞூ
L	டு	G
ண	ഞ്ച	ணூ
த	ਭੀ	தூ
ந	நு	நூ
ன	னு	னூ

x	х + ्П	x + ৣ •Ð
	4	ГP
О	மு	மூ
ш	Щ	T)
J	ரு	ரூ
Д	Д	றூ
ಉ	லு	லூ
ள	ளு	ளூ
மி	ழ	ழ
ഖ	ഖ	ഖ്യ

To the right of g, 知, 和, or 歌, these forms have a spacing form. For example:

5. The vowel sign 600 changes to 20 to the left of 6001, 601, 60, or 611.

Remember that this change takes place after the vowel reordering: in the first example, the vowel 600 follows 6001 in the backing store. After vowel reordering, it is on the left of 6001, and thus changes form. The complete process is:

$$6000 + 6000 \rightarrow 6000 + 60000 \rightarrow 26000$$

6. The consonant  $\Pi$  changes shape to  $\Pi$ .

This occurs when the  $\Pi$  would not be confused with the vowel  $\Pi$ . That is, when  $\Pi$  is combined with (600),  $\Omega$ , or  $\mathcal{C}$ .