Indic Scripts in Unicode

Summary

¶3 This document is a *very rough draft* of a description of the encoding of the Indic scripts in Unicode. It probably contains many mistakes. It is not a publication of the Unicode Consortium, nor has it been endorsed in any way by the Unicode Technical Committee.

Contents

1 Introduction

^{¶4} Most of the scripts of South Asia, from north of the Himalayas to Sri Lanka in the south, from Pakistan in the west to the easternmost islands of Indonesia, are derived from the ancient Brahmi script. The oldest lengthy inscriptions of India, the edicts of Ashoka from the third century BCE, were written in two scripts, Kharoshthi and Brahmi. These are both ultimately of Semitic origin, probably deriving from Aramaic, which was an important administrative language of the Middle East at that time. Kharoshti, written from right to left, was supplanted by Brahmi and its derivatives. The descendants of Brahmi spread with myriad changes throughout the subcontinent and outlying islands. There are said to be some 200 different scripts deriving from it. By the eleventh century, the modern script known as Devanagari was in ascendancy in India proper as the major script of Sanskrit literature.

¶5 The North Indian branch of scripts was, like Brahmi itself, chiefly used to write Indo-European languages such as Pali and Sanskrit, and eventually the Hindi, Bengla, and Gujarati languages, though it was also the source for scripts for non-Indo-European languages such as Tibetan, Mongolian, and Lepcha.

¶6 The South Indian scripts are also derived from Brahmi and, therefore, share many structural characteristics. These scripts were first used to write Pali and Sanskrit but were later adapted for use in writing non-Indo-European languages—namely, the languages of the Dravidian family of southern India and Sri Lanka. Because of their use for Dravidian languages, the South Indian scripts developed many characteristics that distinguish them from the North Indian scripts. South Indian scripts were also exported to the southeast Asia and were the source of scripts such as Lanna and Myanmar, as well as the insular scripts of the Philippines and Indonesia.

¶7 The shapes of letters in the South Indian scripts took on a quite distinct look from the shapes of letters in the North Indian scripts. Some scholars suggest that this occurred because writing materials such as palm leaves encouraged changes in the way letters were written.

¶8 The common origin explains why those scripts share so many common features and why a side-by-side comparison of a few will often reveal structural similarities even in the modern letterforms. It is useful to exploit this situation when describing the scripts. Section <u>2</u>, <u>Common Features</u>, describes those common features. Section <u>3</u>, <u>Encoding Models</u>, describes the various encoding models which are used in Unicode and gives encoding rules which are shared by all the scripts. Each script, having evolved from this common origin, and being pressed into service to write languages with different phonological requirements, has developed unique particularities; those are discussed in section <u>6</u>, <u>South Asian Scripts</u>.

2 Common Features

¶9 In this section, we will present the characteristics which are broadly shared by the scripts derived from the Brahmi script. Most of the examples will use Devanagari as this is the most widely used script of that lineage.

2.1 Letters

¶10 From the first century BCE on, the evolution of the Indic scripts was heavily influenced by the phonetics of Sanskrit, as developed by the ancient pandits. Thus the repertoire and ordering of the basic letters shows a remarkable consistency. Even when the scripts have been adapted to write new languages, many of the original features remained intact, if only to accomodate the numerous Sanskrit loan words and to transcribe accurately religious texts.

¶11 This document contains many small fragments of text using the Indic scripts. Because a reader may not be familiar with all the shapes, we often include a transliteration using Latin characters, as shown in the tables below.

¶12 The traditional repertoire starts with the short and long vowels:

short long अ a आ ā इ i ई ī उ u ऊ ū ए ē ऐ ai ओ ō औ au

¶13 Nasalized vowels are written with the diacritic sign $\overset{\circ}{\sim}$ known as the *candrabindu* in Hindi.

¶14 Next come the vocalic letters:

short long ऋ r, ऋ r, ऌ l, ऌ l,

¶15 Next come consonant letters. The first group covers the occlusive consonants, organized by point of articulation and other phonetic features:

	Voice	eless	Voi	ced	Nasals
	unaspirated	l aspirated	l unaspirated	l aspirated	l
velar	क ka	ख kha	ग ga	ਬ gha	ৰু na
palatal	च ca	छ cha	ज ja	झ jha	স ña
retroflex	द ța	ਰ țha	ड d़a	ਫ d़ha	ण ṇa
dental	त ta	थ tha	द da	ध dha	न na
labial	Ч ра	फ pha	ब ba	भ bha	म ma

¶16 Next come the sonorants and fricatives:

Sonorants Fricatives palatal य ya श śa retroflex र ra ष ṣa dental ल la स sa labial व va

¶17 The remaining consonants are not classified by articulation:

ਫ ha

ळ la

¶18 Nasalized consonants are written with the diacritic sign is known as the anusvara in Hindi.

¶19 The unvoiced aspiration is written with the diacritic sign : known as the visarga in Hindi.

¶20 As in most writing systems, the traditional inventory of letters often needs to be extended to transcribe words of other languages. A common device in the Indic scripts is the addition of the diacritic sign \bigcirc known as the *mukta* in Hindi.

¶21 The Dravidian languages, spoken in the south of India, have additional letters to write sounds which are not found in the Indo-European languages of the north, illustrated here by their form in Tamil:

പ്പ ஒ o ഞ <u>n</u>a ற <u>r</u>a ழ <u>l</u>a

2.2 Digits and Numbers

¶22 Many of the Indic scripts have digits, which can often be used in decimal notation. In some cases, these script-specific digits are still in modern use; in other cases, the Latin digits are used instead in modern texts.

¶23 A number of scripts also have distinct signs which are assembled to write fractions.

2.3 Punctuation and Other Signs

¶24 Traditional punctuation is achieved by the sign |, called *danda* for a full stop, and by the sign ||, called *double danda*, to mark the end of a verse in traditional texts. The term *danda* is from Sanskrit, and those signs sometime have another name in the Indic languages (for example *viram* in Hindi). As with the letters, the appearance of those signs varies across scripts.

¶25 Many modern languages written in the Indic scripts intersperse punctuation derived from the Latin script. Thus U+002C ',' COMMA and U+002E '.' FULL STOP are freely used in writing Hindi, and the *danda* is usually restricted to more traditional texts. The *danda* may be preserved when such traditional texts are transliterated into the Latin script.

 $\[\] 26 The avagraha symbol S is used todo . It is not present in every script. \]$

¶27 The *om* symbol 35 'om is a religious todo . It is not present in every script.

¶28 The rupee sign \mathfrak{Z}_0 (Gujarati) is used todo. It is not present in every script.

2.4 Typographic Clusters

¶29 In many scripts, text is essentially written by the juxtaposition of letters, with limited typographic interaction between successive letters; furthermore, those interactions, such as an "fi" ligature, serve mostly an aesthetic purpose and are not used to indicate distinctions of meaning.

¶30 The Indic scripts, on the other hand, use typographic interactions to indicate differences of meaning. For example, the letter $\overline{\sigma}$ ka

followed by the letter $\overline{\langle}$ ra, without typographic interaction ($\overline{\sigma}\overline{\langle}$) is usually read /kara/; with typographic interaction ($\overline{\sigma}$), the same sequence of letters is usually read /kra/. We call *typographic cluster* a succession of letters which interact typographically. The visible result of the interaction is that the shapes of the component letters are modified or even merged and spatially rearranged, both horizontally and vertically.

¶31 We use the notation $l_1 \bullet l_2 \dots \bullet l_n$ to indicate that the letters $l_1, l_2 \dots l_n$ form a typographic cluster. Thus $\langle \overline{\sigma} | ka \bullet \overline{\langle} ra \rangle \rightarrow \overline{\sigma} kra$, while $\langle \overline{\sigma} | ka, \overline{\langle} ra \rangle \rightarrow \overline{\sigma} \overline{\langle} kara$.

¶32 A consonantal typographic cluster corresponds roughly to a sequence of consonant sounds followed by a vowel sound. Graphically, this translates in *core consonants, satellite consonants, vowel signs* and *other signs* for secondary phonetic features. The vast majority of consonantal typographic clusters are made of a single core consonant.

¶33 A vowel typographic cluster corresponds roughly to a vowel sound alone. Graphically, this translates in a *vowel letter, vowel signs* and *other signs* for secondary phonetic features. The vast majority of vowel typographic clusters are made of a single vowel letter.

¶34 *Core Consonants.* The core of a typographic cluster writes a sequence of consonant sounds. Its written appearance is often a shape unique to that sequence:

क ka • ष şa → क्ष kşa श śa • र ra → श्र śra द da • म ma → द्म dma

¶35 In other cases, it is more straightforward to recognize a combination of the basic letters; typically, one letter retains its normal shape, and the others take a reduced form:

क ka • र ra → क्र kra क ka • न na → क्न kna

¶36 Yet in other cases (rare in the scripts of the North and essentially the norm in Tamil), the core is rendered by the juxtaposition of the basic letters, with the addition of the diacritic sign \bigcirc known as *halant* in Hindi:

ङ na•ग ga → ङ्ग nga

¶37 Those examples illustrate more a continuum in the assembly of basic letters than strict categories in which every typographic cluster core falls neatly.

¶38 *Satellite Consonants.* Depending on the script, some consonants in initial or final position within a cluster may behave graphically more as satellites to the rest of the core of the typographic cluster. For example, in Devanagari, an initial \exists ra is depicted by a small hook, or *repha* at the extreme right of the core; in Myanmar, a final \exists ra will be depicted by an enclosing sign on the left:

र ra•क ka → र्क rka र ra•क ka•न na → र्क्न rkna の ka•ရ ra → ऴ kra

¶39 *Vowel signs.* Consonantal typographic clusters can include one or more *vowel sign* to denote their vowel sound. Those signs are diacritic forms of the vowel letters. We illustrate those on a few cores:

 क ka
 र ra
 क्न kna
 क rka

 आ ā
 का kā
 रा rā
 क्ना knā
 का rkā

 इ i
 कि ki
 रि ri
 क्न kni
 कि rki

 ई ī
 की kī
 री rī
 क्नी knī
 की rkī

 उ u
 कु ku
 रु ru
 क्नु knu
 कु rku

 ऊ ū
 कू kū
 रु rū
 क्नू knū
 कू rkū

 ए ē
 के kē
 रे rē
 क्ने knē
 के rkē

ऐ ai कै kai रै rai क्नै knai कैं rkai ओ ō को kō रो rō क्नो knō कीं rkō औ au कौ kau रौ rau क्नौ knau कीं rkau

¶40 The normal form of a vowel letter is called an *independent vowel* or a *vowel letter*, while its diacritic form is called a *dependent vowel* or *vowel sign*.

¶41 As can be seen in those examples, the vowel signs can be placed all around the core, even though they write a vowel sound pronounced after that core. In some cases, they actually fuse into the core, as in $\overline{\nabla}$ ru or $\overline{\nabla}$ rū.

¶42 In some scripts, a vowel sign can be rendered by fragments on multiple sides of the core. For example, the vowel sign of Bengla \mathfrak{F} au attaches to the simple core \mathfrak{F} ka on both sides: (को kau; the vowel sign of Oriya la au attaches to the simple core \mathfrak{F} ka on the left, top and right sides: (କ) kau. In Unicode terminology, those vowels are called *split vowels*.

¶43 In some cases, vowel signs are also present in vowel typographic clusters; they attach to the vowel letter that form the core of the typographic cluster.

144 *Clusters without vowels.* Sometimes clusters represent only a sequence of consonant sounds without a vowel sound. This is most common at the end of words of Sanskrit or foreign origin. This situation is sometimes reflected in writing, by the addition of the halant sign to the cluster.

Hindi ् todo need example Telugu ్ అప్ ap (water)

¶45 *Other Signs.* Finally, the signs $\check{\circ}$ \check{m} , $\dot{\circ}$ \check{m} and \circ : \dot{h} can be added to a cluster. Their position may vary across scripts. For example, the Bengla \circ ? \dot{m} attaches to the right of the core, rather than to the top.

¶46 *Vowel Clusters.* When a word starts with a vowel sound, there is no consonant to which to attach the diacritic form of the vowel letter. Instead, the non-diacritic form of the vowel acts as the core of the cluster, as in Hindi आगर agara. This core can be decorated with other signs, for example Hindi উঁন্যা ūmঁcā.

¶47 *Clusters and Syllables.* The typograhic clusters need not correspond exactly with phonological syllables, especially when a consonant cluster is involved. For example, Hindi पुत्ति purti is made of two typographic clusters y pu and ति rti, while the syllables are /pur/ and /ti/.

3 Encoding Models

3.1 Three Models

¶48 Three distinct models have been used in the Unicode standard for encoding the Indic scripts. They differ principally in the scheme used to indicate typographic clusters.

¶49 The *consonant linking model* (also known as the *virama model*) encodes one coded character for each consonant; one *consonant linker* coded character; one coded character for each vowel which can appear as a *vowel letter*; and one coded character for each vowel which can appear as a *vowel sign*. A typographic cluster is formed by placing the consonant linker between the consonants of the cluster.

¶50 The *subjoined consonant model* encodes two or more characters for each consonant, one non-combining and the others combining; one coded character for each vowel which can appear as a *vowel letter*; and one coded character for each vowel which can appear as a *vowel letter*; and one coded character for each vowel which can appear as a *vowel letter*; and one coded character for each vowel which can appear as a *vowel letter*; and one coded character for each vowel which can appear as a *vowel sign*. A typographic cluster is formed by using the non-combining coded character for the first consonant, and combining coded

characters for the remaining consonants.

¶51 In those two models, coded characters appear in logical order, that is in the same order as the sounds they write. This is particularly visible for vowels rendered on the left side of their typographic cluster, or the left part of a two-part vowel: the corresponding characters appear toward the end of the typographic cluster.

¶52 The *visual model* encodes a non-combining character for each spacing sign, and a combining character for each non-spacing sign. Furthermore, the coded characters appear in textual representation in visual order. This model was motivated by simple rendering technologies (such as typewriters).

¶53 For a given script, the choice of model in Unicode is guided by a number of factors, including legacy practices, closeness to other scripts, and perception of the script by its users.

¶54 In all cases, Unicode does not encode a character for each typographic cluster, but one or two for each letter. On the other hand, the glyph complement of fonts is more likely to include separate glyphs for many of the typographic clusters, or at least for their cores. Thus, the number of glyphs in an Indic font can far exceed the number of corresponding characters. It is up to the rendering software to select the appropriate glyph sequence for a given character sequence.

3.2 The Consonant Linking Model

¶55 The consonant linking model uses a *consonant linker* coded character to link together the consonants of a typographic cluster.

¶56 The consonant linking model is used for scripts covered by ISCII (which also uses that model) and closely related scripts: Devanagari, Bengla, Gurmukhi, Gujarati, Oriya, Tamil, Telugu, Kannada, Malayalam, Sinhala.

¶57 Coded Characters. In the consonant linking model, we have

• **§**58 one coded character for each consonant; for example:

क ka \rightarrow U+0915 क DEVANAGARI LETTER KA

₹ ra \rightarrow U+0930 ₹ DEVANAGARI LETTER RA

• ¶59 one coded character which acts as a *consonant linker*; for example:

U+094D ू DEVANAGARI SIGN VIRAMA

• **§**60 one coded character for each vowel which can appear as a *vowel letter*; for example:

 ${\ensuremath{\smuremath{\ens$

• **§**61 one coded character for each vowel which can appear as a *vowel sign*; for example:

 ${\ensuremath{\overline{s}}}\xspace$ i ${\rightarrow}$ U+093F \bigcirc Devanagari vowel sign i

• **§**62 one coded character for each diacritic sign; for example:

ὄṁ→U+0901 ὄ DEVANAGARI SIGN CANDRABINDU

¶63 A cluster such as $\overline{\sigma}$ ka • $\overline{\tau}$ ra • $\overline{\mathfrak{F}}$ i $\rightarrow \overline{\mathfrak{F}}$ kri is represented by the sequence of coded characters:

U+094D ् DEVANAGARI SIGN VIRAMA

U+0930 ₹ DEVANAGARI LETTER RA

U+093F ि DEVANAGARI VOWEL SIGN I

164 That is: each consonant is represented by the corresponding coded character; successive consonants in the same cluster are linked by the consonant linker coded character; vowels signs in the cluster are represented by the corresponding vowel sign coded character.

¶65 A cluster such as $\exists i \bullet \overset{\circ}{i} \overset{\circ}{m} \rightarrow \overset{\circ}{\exists} \overset{\circ}{i} \overset{\circ}{m}$ is represented by the sequence of coded characters:

U+0907 इ DEVANAGARI LETTER I U+0901 ँ DEVANAGARI SIGN CANDRABINDU

166 That is: each vowel letter is represented by the corresponding vowel letter coded character; diacritic signs are represented by the corresponding coded character.

¶67 The ordering of the coded characters matches the pronunciation order. In the cluster $\overline{\sigma}$ ka • $\overline{\varsigma}$ ra • $\overline{\varsigma}$ i $\rightarrow \overline{\beta}\overline{\rho}$ kri, the coded character for the vowel sign is last in the sequence because it is pronounced last, eventhough it is written first.

168 The representative glyph for the consonant linker is often the halant sign of the script: U+094D \bigcirc DEVANAGARI SIGN VIRAMA. However, the coded character needs to be understood more as a control character, which denotes the formation of typographic clusters, rather than as a graphic sign. The rendering of a cluster may result in the halant sign, but this caused only indirectly by the consonant linker. This situation is very similar to hyphenation opportunities which are denoted by U+00AD SOFT HYPHEN and can result in the - sign.

¶69 *Atomicity of Vowel Letters.* It is sometimes possible to analyze the visual appearance of vowel letters, and notice that they are formed by assembling the rendering of other vowel letters and signs. For example, the Devanagari letter $\Im \ \bar{o}$ is clearly composed graphically of an $\Im \ \bar{a}$ and a $\ \bar{o}$ -ē. This characteristic is more than a graphical accident, and leads to two possible approaches for the encoding of vowel letters. The approach chosen by Unicode is to encode each vowel letter separately, without any decomposition: $\Im \ \bar{a}$ is represented by the single coded character U+0913 $\Im \ DEVANAGARI LETTER O$ and not by the sequence $\langle U+0906 \ \Im \ DEVANAGARI LETTER AA, U+0947 \ DEVANAGARI VOWEL SIGN E>$.

¶70 Similarly, vowel signs can be sometimes be analyzed. For example, Bengla (1 is composed graphically of a and a 1. The approach chosen by Unicode is to encode the constituent parts: <math>(1 is represented by the sequence of coded characters < U+09C7)BENGALI VOWEL SIGN E, U+09BE 1 BENGALI VOWEL SIGN AA> rather than by a single coded character that has no decomposition. For compatibility with existing standards, some of these sequences are also encoded as a single code point (U+09CB (1 BENGALI VOWEL SIGN O in this case), but that code point canonically decomposes in the sequence.

¶71 In the descriptions of the individual scripts, we will list the vowel letters and signs which can be confused graphically.

¶72 Use of Joiners for Consonants. As we have seen before, there is generally a high degree of typographic interaction within a cluster, and the amount of interaction is better understood as a continuum. Nevertheless, it is possible to identify three broad forms for the rendering of the core of a cluster:

- ¶73 fully-joined forms, where the consonants merge into a shape where both original shapes are modified
- ¶74 conjoined forms, where one of the consonant retains its full shape, and the others take reduced forms that attach to the unmodified consonant
- ¶75 non-joined forms, where each consonant retains its full shape, and a vowel killer sign is added to all but the last consonant

¶76 The following illustrates the three forms:

fu	lly-joined	l conjoined 1	non-joined
क ka • ष ṣa	क्ष	क्ष	क्ष
क ka • र ra	क्र	कर	क्र
क ka • न na		क्न	क्न
ङ na • ग ga	ङ्ग		ङ्ग

¶77 As shown above, a given cluster can generally be rendered with more than one form, although not all forms need to exist. While all the forms are equivalent and can be used intercheangly, the context (point in time, nature of the text, etc) usually dictates one preferred form. For example, the forms \overline{qq} and \overline{qq} can be understood by all readers, but the form \overline{q} is pretty much the norm. On the other hand, the form

ক্ন is mostly found in older texts, and the form ভ্র্য is more commonly found in modern texts. The historical trend is to abandon the use of fully-joined forms for the rare clusters.

¶78 A rendering system can render a cluster with any of the possible forms, and should normally use what is considered the preferred form. Some degree of control is achieved by the selection of a font appropriate for the nature of the text, much like one may wish to select a font with "ct" and "st" ligatures for an 18th century text.

¶79 The author of a text can exert some degree of control over which form should be used by the rendering system, by inserting the characters U+200C ZERO WIDTH NON-JOINER and U+200D ZERO WIDTH JOINER appropriately. The general rules are as follows.

180 Without any joiner, no preference is expressed, and the rendering system should select the fully-joined form if it exists in the font.

क ka • ष sa <0915, 094D, 0937> क्ष preferred, then क्ष, then क्ष क ka • र ra <0915, 094D, 0930> क्र preferred, then कर, then क्र क ka • न na <0915, 094D, 0928> कन preferred, then क्न ङ na • ग ga <0919, 094D, 0917> ङ्ग preferred, then ङ्ग

¶81 With a ZWJ between a consonant and a consonant linker, the fully-joined form of the cluster is discouraged. The cluster should be rendered preferrably with a conjoined form, with the consonant before the ZWJ in full form and the consonant after the consonant linker in the a reduced form. If no such form is possible, then a non-joined form should be used.

क ka • ष sa <0915, 200D, 094D, 0937> क्ष

क ka • र ra <0915, 200D, 094D, 0930> क्र preferred, then क्र

क ka • न na <0915, 200D, 094D, 0928> क्न

জ na • ग ga <0919, 200D, 094D, 0917> জ্য

¶82 With a ZWJ between a consonant linker and a consonant, the fully-joined form of the cluster is discouraged. The cluster should be rendered preferrably with a conjoined form, with the consonant after the ZWJ in full form and the consonant before the consonant linker in the a reduced form. If no such form is possible, then a non-joined form should be used.

क ka • ष sa <0915, 094D, 200D, 0937> वष preferred, then क्ष क ka • र ra <0915, 094D, 200D, 0930> कर preferred, then क्र क ka • न na <0915, 094D, 200D, 0928> कन preferred, then क्न ङ na • ग ga <0919, 094D, 200D, 0917> ङ्ग

¶83 With a ZWNJ between a consonant linker and a consonant, the fully-joined form and the conjoined form of the cluster are discouraged. The cluster should be rendered with both consonants in full form, and a consonant killer sign.

क ka • ष şa <0915, 094D, 200C, 0937> क्ष क ka • र ra <0915, 094D, 200C, 0930> क्र क ka • न na <0915, 094D, 200C, 0928> क्न ङ ńa • ग ga <0919, 094D, 200C, 0917> ङ्ग

¶84 The rules above are generally applicable to all the scripts covered in this document with the notable exception of Sinhala.

185 Use of Joiners for Vowel Signs. A similar situation exists for the rendering of vowel signs in a cluster. Two broad forms can be identified:

- **§**86 fully-joined forms, where the vowel sign and the rest of the cluster merge into a unique shape
- **187** conjoined forms, where the general diacrictic form is used

188 Here are illustrations of the two forms:

fully-joined conjoined গga•উu গু গু

ক ka • উ u কু

¶89 Again, the joiners characers can be used to encourage one or another rendering:

¶90 Without any joiner, no preference is expressed, and the rendering system should select the most common form.

গ ga • উ u <0997, 09C1> গু fallback to গু ক ka • উ u <0995, 09C1> ক্

¶91 With a ZWJ between the consonant and the vowel sign, the cluster should be rendered in a fully-joined form.

গ ga • উ u <0997, 200D, 09C1> গু fallback to গু

ক ka • উ u <0995, 200D, 09C1> কু

¶92 With a ZWNJ between the consonant and the vowel sign, the cluster should be rendered in a conjoined form.

গ ga • উ u <0997, 200C, 09C1> গু fallback to গু ক ka • উ u <0995, 200C, 09C1> কু

¶93 Use of Joiners for Typographic Fragments. Unicode also supports the representation of typographic fragments, such as the half form $\overline{\sigma}$ of the letter $\overline{\sigma}$ ka. This can be useful in didactic materials. This support is typically limited to the half forms of consonant letters and is achieved using the joiners. In general, the approach is to use sequences of the form <consonant, consonant linker, ZWJ> for consonant which take a reduced form when first in a cluster, and <U+0020 SPACE, ZWJ, consonant linker, consonant> for consonants which take a reduced form when second in a cluster. The script descriptions below will give the precise list of sequences which are recognized by the Unicode standard.

¶94 *Other Uses of Joiners.* In a few situations, U+200C ZERO WIDTH NON-JOINER and U+200D ZERO WIDTH JOINER are used for other purposes than those listed above. Unlike the cases listed so far, some of those situations correspond to semantic differences.

3.3 The Subjoined Model

¶95 The subjoined model is used for Tibetan and to a large extent for Myanmar.

¶96 Most of the scripts using the subjoined model do not have vowel letters, and only vowel signs.

¶97 Coded Characters. In the subjoined model, we have:

• ¶98 two or more coded characters for each consonant; for example:

 $\P\ ka \longrightarrow U\!\!+\!\!0F40$ $\P\ TIBETAN\ LETTER\ KA$

 \rightarrow U+0F90 \Re TIBETAN SUBJOINED LETTER KA

 τ ra → U+0F62 τ TIBETAN LETTER RA

 \rightarrow U+0FB2 $\stackrel{\scriptscriptstyle \odot}{\scriptstyle \sim}$ TIBETAN SUBJOINED LETTER RA

• **199** one coded character for each vowel which can appear at a *vowel sign*; for example:

° i → U+0F72 ° TIBETAN VOWEL SIGN I

• ¶100 one coded character for each diacritic sign; for example:

^{*}→ U+0F83 ^{*} TIBETAN SIGN SNA LDAN

¶101 A cluster such as \mathbb{T} ka • \mathbb{T} ra • \mathbb{T} i \rightarrow kri is represented by the sequence of coded characters:

U+0F40 TIBETAN LETTER KA

U+0FB2 $\stackrel{\mbox{\tiny \Im}}{=}$ TIBETAN SUBJOINED LETTER RA

U+0F72 [°] TIBETAN VOWEL SIGN I

3.4 The Visual Model

¶102 The visual model is used where there is a well-established legacy practice: Thai, Lao.

4 Commonalities in Encoding

4.1 The ISCII Standard

¶103 A number of blocks in the Unicode Standard are based on ISCII-1988 (Indian Script Code for Information Interchange). The ISCII standard of 1988 differs from and is an update of earlier ISCII standards issued in 1983 and 1986.

¶104 The Unicode Standard encodes the characters in the same relative positions as those coded in positions A0-F4₁₆ in the ISCII-1988 standard. The same character code layout is followed for nine Indic scripts in the Unicode Standard: Devanagari, Bengla, Gurmukhi, Gujarati, Oriya, Tamil, Telugu, Kannada, and Malayalam. This parallel code layout emphasizes the structural similarities of the Brahmi

scripts and follows the stated intention of the Indian coding standards to enable one-to-one mappings between analogous coding positions in different scripts in the family. Sinhala, Tibetan, Thai, Lao, Khmer, Myanmar, and other scripts depart to a greater extent from the Devanagari structural pattern, so the Unicode Standard does not attempt to provide any direct mappings for these scripts to the Devanagari order.

¶105 In November 1991, at the time *The Unicode Standard, Version 1.0*, was published, the Bureau of Indian Standards published a new version of ISCII in Indian Standard (IS) 13194:1991. This new version partially modified the layout and repertoire of the ISCII-1988 standard. Because of these events, the Unicode Standard does not precisely follow the layout of the current version of ISCII. Nevertheless, the Unicode Standard remains a superset of the ISCII-1991 repertoire. Modern texts encoded with ISCII-1991 may be automatically converted to Unicode code points and back to their original encoding without loss of information.

4.2 Letters with Nukta

¶106 The scripts which use a *nukta* to create new letters encode their own nukta character.

¶107 Some, but not all, combinations of a base letter with a nukta which are actually used are also encoded as individual coded characters. All those coded characters have a canonical decomposition into the coded character for the base letter and the coded character for the nukta. Either the individual coded characters or their decompositions can be used.

4.3 Dandas

¶108 The *danda* and *double danda* are encoded only once at U+0964 | DEVANAGARI DANDA and U+0965 || DEVANAGARI DOUBLE DANDA. Those two coded characters, despite their location in the Devanagari block and the word DEVANAGARI in their names, are used regardless of the script.

4.4 Spoken Versus Written

¶109 While the basic inventory is strongly aligned with pronunciation, there are numerous cases of a disconnect between the spoken and the written. In Hindi, many words do not end with a final /a/ sound, yet they are written as if they did: राम rāma is pronounced /rām/; the same can occur inside words: सरकार sarakāra is pronounced /sarkār/.

¶110 The Unicode standard is first and foremost an encoding of the written signs rather than an encoding of the sounds. Therefore, सरकार is represented by < sa, \forall ra, $\overline{\sigma}$ ka, \bigcirc I - \overline{a} , \forall ra>, eventhough there is no /a/ sound after the first /r/. The more exact representation of the sounds, < sa, \forall ra, \bullet ka, \bigcirc I - \overline{a} , \forall ra> would actually look like सकरि sarkāra, but is not the common way to write the word.

¶111 Another manifestation of the predominance of the written can be seen with the nasal consonants and the anusvara. The Oriya word $\Im \mathfrak{T}$ and a an also be written with an anusvara: $\Im^{\circ} \mathfrak{T}$ and \mathbb{T} and \mathbb{T}

¶112 This is very similar to the way "color" and "colour" are two different spellings of the same English word, and have different representations matching their written form.

5?

5.1 Collation

¶113 todo

5.2 Security

¶114 todo

5.3 Identifiers

¶115 todo

6 South Asian Scripts

6.1 Devanagari

¶116 The Devanagari script is used for writing classical Sanskrit and its modern historical derivative, Hindi. Extensions to the Sanskrit repertoire are used to write other related languages of India (such as Marathi) and of Nepal (Nepali). In addition, the Devanagari script is used to write the following languages: Awadhi, Bagheli, Bhatneri, Bhili, Bihari, Braj Bhasha, Chhattisgarhi, Garhwali, Gondi (Betul, Chhindwara, and Mandla dialects), Harauti, Ho, Jaipuri, Kachchhi, Kanauji, Konkani, Kului, Kurnaoni, Kurku, Kurukh, Marwari, Mundari, Newari, Palpa, and Santali.

¶117 Devanagari is written horizontally from left to right, and words are separated by space.

¶118 *Standards.* The Devanagari block of the Unicode standard is based on ISCII-1988. See section <u>4.1, *The* ISCII *Standard*</u>, for more details on the relation to ISCII.

¶119 *Encoding.* The encoding of the Devanagari script uses the consonant linking model. The consonant linker coded character is U+094D \bigcirc DEVANAGARI SIGN VIRAMA.

6.1.1 Letters

¶120 Devanagari has the full complement of basic letters presented in section 2.1, Letters.

¶121 The vowels:

	sho	ort	lor	ng
let	ter	sign	letter	sign
अ a	0905		आā 0906	○T -ā 093E
इ i	0907	○ -i 093F	ई 1 0908	ी -ī 0940
ਤ u	0909	ु -u 0941	ऊ ū 090A	ू -ū 0942
ऎ e	090E	े-e 0946	एē 090F	े -ē 0947
ओ ०	0912	ो -0 094A	ओ ō 0913	ो -ō ०९४८

122 Other vowels:

lette	r		sign	ı	
ऐai	0910	<u></u> `	ai	0948	
औ au (0914	ौ	-au	094C	
Ϋĕ (090D	ँ -	ĕ	0945	todo usage?
ऑŏ	0911	ॉ	-ŏ	0949	todo usage?
ॲ ॲ	0972				todo usage?
ऄ ः	0904				todo usage?
		Š	ॅ	0955	for Avestan

¶123 Sanskrit:

sho	ort	long	,
letter	sign	letter	sign
ऋ r 090B ृ	-ŗ 0943 😿 ŗ 0960) ्रु -ऱ्र ०९६२	
ਲ ਹੈ 090C ੂੰ	-] 0944 ऌ] 0961	्र - रे 0963	

¶124 The plosive consonants:

	Voiceless				Voiced		Nasals			
		unaspira	ited	aspira	ited	unas	pirated		aspirated	
velar	क ka	0915 ख k	kha 0916	ग ga (0917 घ gha	0918	ৰু na	0919		
palatal	च ca	091A छ C	ha 091B	ज ja ।	091C झ jha	091D	স ña	091E		
retroflex	ਟ ța	091F ਰ țl	ha 0920	ड da	0921 ढ ḍha	0922	ण ņa	0923		
dental	त ta	0924 थ t	ha 0925	द da ।	0926 ध dha	0927	न na	0928		
labial	ч ра	092A	ha 092B	ब ba	092C भ bha	092D	म ma	092E		

$\P125$ The sonorant and fricative consonants:

Sonorants Fricatives

palatal य ya 092F श śa 0936 retroflex र ra 0930 ष sa 0937 dental ल la 0932 र sa 0938 labial व va 0935

¶126 Other consonants:

ষ	ষ	097A	HEAVY YA	todo
ॹ	ज़	0979	ZHA	for Avestan
ग	ग	097B	GGA	for Sindhi
ত	ড	097C	JJA	for Sindhi
<u>ड</u>	<u>ड</u>	097E	DDDA	for Sindhi
ब	ब	097F	BBA	for Sindhi
?	?	097D	GLOTTAL STOP	for Limbu

¶127 Signs:

vowel nasalization \ddot{o} \ddot{m} 0901 candrabindu consonant nasalization \dot{o} \dot{m} 0902 anusvara unvoiced aspiration : h 0903 visarga

¶128 Forms with nukta:

୍	093C	nukta
क़ qa	0958	= 0915, 093C
ख़ <u>k</u> ha	0959	= 0916, 093C
ग ġa	095A	= 0917, 093C
ज़ za	095B	= 091C, 093C
ड़ ṛa	095C	= 0921, 093C
ढ़ ṛha	095D	= 0922, 093C
फ़ fa	095E	= 092B, 093C
य़ yंa	095F	= 092F, 093C

¶129 Nepali and Marathi use an additional letter \exists called the *eyelash ra*, to represent an /r/ sound. This letter is used contrastively with the letter \exists ra: for example Marathi \exists af daryā "ocean" vs. \exists \exists II daryā "valleys". For legacy reasons, Unicode does not encode a separate character for the eyelash ra; instead, it is represented by the sequence $\langle U+0931 \exists DEVANAGARI LETTER RRA, U+094D \rangle$ DEVANAGARI SIGN VIRAMA>. For compatibility with Unicode 2.0, the sequence $\langle U+0930 \exists DEVANAGARI LETTER RA, U+094D \rangle$ U+094D \Diamond DEVANAGARI SIGN VIRAMA, U+200D ZERO WIDTH JOINER> can also be supported. Although the eyelash ra is represented as if it were a half form of either \exists or \exists , it is actually a separate letter.

¶130 The Sindhi language uses implosive sounds which are written using the following four characters: U+097B \exists DEVANAGARI LETTER GGA, U+097C \exists DEVANAGARI LETTER JJA, U+097E \underline{s} DEVANAGARI LETTER DDDA and U+097F \underline{s} DEVANAGARI LETTER BBA. Implementations may recognize the previously recommended sequences using U+093C \bigcirc DEVANAGARI SIGN NUKTA or U+0952 \bigcirc DEVANAGARI STRESS SIGN ANUDATTA, but authors are strongly encouraged to use the specific characters listed above.

¶131 In older texts, alternate rendering for the vowel signs e, ai, o and au are sometime used. Those are represented using U+094E⊺ DEVANAGARI VOWEL SIGN PRISHTHAMATRA E in combination with other vowel signs:

¶132 Note that this is an exception to the general rule that vowel signs are represented atomically. Because of the combining class of the characters involved, to avoid ambiguity in representation U+094E should be first.

¶133 Marathi Allographs. In Marathi and some South Indian orthographies, variant glyphs are preferred for U+0932 ল DEVANAGARI LETTER LA and U+0936 श DEVANAGARI LETTER SHA:

¶134 todo figure 9.9

6.1.2 Digits and Numbers

¶135 Devanagari has a full set of decimal digits, encoded at U+0966 ° DEVANAGARI DIGIT ZERO ... U+096F DEVANAGARI DIGIT NINE.

¶136 There are no additional number characters.

6.1.3 Punctuation and Signs

¶137 The danda | is known as *purna viram* in Hindi, and the double danda || is known as *deergh viram* in Hindi. Devanagari uses the common characters encoded as U+0964 | DEVANAGARI DANDA and U+0965 || DEVANAGARI DOUBLE DANDA.

¶138 The avagraha is encoded as U+093D S DEVANAGARI SIGN AVAGRAHA.

¶139 The *om sign* is encoded as U+0950 35 DEVANAGARI OM.

¶140 ° appears after letters or combinations of letters and marks the sequence as an abbreviation. It is encoded as U+0970 ° DEVANAGARI ABBREVIATION SIGN.

¶141 o is used for todo and is encoded as U+0951 o DEVANAGARI STRESS SIGN UDATTA.

¶142 \bigcirc is used for todo and is encoded as U+0952 \bigcirc DEVANAGARI STRESS SIGN ANUDATTA.

¶143 U+0953 OEVANAGARI GRAVE ACCENT and U+0954 OEVANAGARI ACUTE ACCENT were originally encoded for todo. Instead, U+xx and U+yy are recommended for those uses.

¶144 todo U+0971 [·] DEVANAGARI SIGN HIGH SPACING DOT

¶145 *Bodo*. The orthography of Bodo uses ' as a tone mark. This sign is called गोजौ कमा (gojau kamaa) and is represented using U+02BC ' MODIFIER LETTER APOSTROPHE. For example: खर' (head), दख'ना (dress commonly worn by Bodo women).

¶146 **Dogri.** The orthography of Dogri uses ' as a tone mark. This sign is called सुर-चि'न्न (sur chinha) and is represented using U+02BC ' MODIFIER LETTER APOSTROPHE. It occurs after short vowels (may be inherent) and indicates a high-falling tone: ख'न (down), कु'न (who), ति'लकना (to slip). After long vowels, a high-falling tone is written using U+0939 ह DEVANAGARI LETTER HA.

¶147 The avagraha 5 is used to indicate extra-long vowels; for example, तला (sole) versus तला5 (pond).

¶148 *Maithili*. The orthography of Maithili uses ' to denote the prolongation of a short अ a and to denote the truncation of words. This sign is called बिकारी कामा (bikari kaamaa) and is represented using U+02BC ' MODIFIER LETTER APOSTROPHE. For example:

कहाँ जहे हो - कत' जा रहल छें (prolongation); कत' पड़ा' गेल', abbreviation for कतए पड़ाए गेलह (where did you go away?).

¶149 The avagraha 5, called बिकारी (bikari), is used to indicate extra-long vowels.

¶150 *Atomicity of Vowel Letters.* Vowel letters are encoded atomically in Unicode, even if they can be analyzed visually as consisting of multiple parts. The table below shows the letters that can be analyzed, the single code point that should be used to represent them in text, and the sequence of code points resulting from analysis that should not be used.

To represent	Use	Do not use
ऄ	0904	<0905, 0946>
आ	0906	<0905, 093E>
স	090A	<0909, 0941>
ऍ	090D	<090F, 0945>
ऎ	090E	<090F, 0946>
ऐ	0910	<090F, 0947>
ऑ	0911	<0905, 0949>
ऒ	0912	<0905, 094A>
ओ	0913	<0905, 094B>
औ	0914	<0905, 094C>
ॲ	0972	<0905, 0945>

6.1.4 Typographic clusters

¶151 *Typographic Cluster Cores.* Here are some of the common typographic clusters made of two consonants which exhibit a fully-joined form:

市 ka • 市 ka

市 ka • 市 ta

中 ka • 市 ta

中 ka • 雨 ta

中 ka

<

द da • ब ba	→द्घ dba
द da • भ bha	\rightarrow द्ध dbha
द da • म ma	→ द्म dma
द da • य ya	→ द्य dya
द da • व va	→ द्व dva
ਟ ța • ਟ ța	→ इ țța
ਟ ța • ਰ țha	→ इ țțha
ਰ țha • ਰ țha	→ ਭ țhțha
ड ḍa • ग ga	→ ड्ग dga
ड d़a • ड d़a	→ ड d़da
ड ḍa • ढ ḍha	→ ड्रु ḍḍha
त ta • त ta	→त्त tta
त ta • र ra	→त्र tra
न na • न na	→ ঈ nna
फ pha • र ra	→ फ्र phra
श śa • र ra	→ श्र śra
ह ha • म ma	→ ह्य hma
ह ha • य ya	→ ह्य hya
ह ha • ल la	→ ह्न hla
ह ha • व va	→ ह्र hva

¶152 In the absence of fully-joined form, the core of a typographic cluster is displayed in a conjoined form, where all but the last consonant are in *half form* and the last consonant is its full form.

¶153 The half form of the consonants $\overline{\Phi}$ ka and $\overline{\Psi}$ pha is obtained by dropping the fragment on the right side of the stem: $\overline{\Phi}$ k and $\overline{\Psi}$ ph.

क ka • ख kha \rightarrow क्ख kkha क ka • ळ ḷa \rightarrow क्ळ kḷa

¶154 The half form of the consonants with a vertical stem on the right (खगघचजझञणतथधननपबभमयलवशषस) is obtained by dropping their vertical stem: खग्टचज्झ ऊण्त्थधनन्य हभ्मटल्द १ ४ २.

ग ga • ट țа \rightarrow गट gța л ga • ळ ļа \rightarrow गळ gļa л ņa • ट țа \rightarrow णट ņța л na • ट țа \rightarrow णट пțа

¶155 The remaining consonants (ङ छट ठ ड ढ द ळ ह) do not have half form. When they are present in non-final position in a cluster core,

the non-joined form of the cluster is displayed, with a halant sign below all but the last consonant.

ट ta • स sa \rightarrow ट्स tsa

¶156 *Satellite Consonants.* \exists ra in initial position is a satellite consonant. It is then displayed as a small hook, called a *repha*, on the top right side of the cluster, including any satellite vowels:

र ra • त ta $\rightarrow \hat{d}$ rta र ra • क ka • ख kha $\rightarrow \hat{d}$ rta र ra • क ka • ख kha $\rightarrow \hat{d}$ rds rkkha र ra • म ma • आ ā $\rightarrow \hat{d}$ rmā र ra • थ tha • ई ī $\rightarrow \hat{d}$ rthī र ra • म ma • ओ ō • ं ṁ $\rightarrow \hat{d}$ rmōṁ

¶157 However, when a cluster is displayed in a non-conjoined form (either because such form is the most natural form or because it was requested by the use of a ZWNJ), the *repha* is displayed on the first consonant: \forall ra • $\overline{\sigma}$ ka • $\overline{\sigma}$ ka • $\overline{\sigma}$ ka. todo wrong rendering

¶158 \triangleleft ra in final position is a satellite consonant. It is displayed as small stroke attached to the vertical stem, or as a circumflex-like sign, called a *vattu*, below the consonants which do not have a stem.

¶159 Satellite Vowels. Some typographic clusters with vowel signs have full-joined forms:

र ra • उ u → रु ru र ra • ऊ ū → रू rū ह ha • ऋ ŗ → ह hŗ

¶160 Otherwise, the vowel sign attaches around the cluster core.

on क ka on क rka आ ā ा -ā का kā कf rkā इ i ि -i कि ki कि rki ई ī ी -ī की kī की rkī उ u ु -u कु ku कु rku ऊ ū ू -ū कू kū कू rkū ऍ ĕ ॅ -ĕ के kĕ के rke ए ē े -ē के kē के rkē

ऐ ai	े -ai	कै kai	कै rkai
ऑ ŏ	ॉ -ŏ	कॉ kŏ	र्को rkŏ
ऒ ०	ो -०	कॊ ko	र्को rko
ओ ō	ो -ō	को kō	को rkō
औ au	ौ -au	कौ kau	कौँ rkau
ऋ rॢ	်-င်	कृ kr	र्कृ rkr
ऋ ŗ	ၙ -ļ	कृ kļ	र्कृ rkļ
ਲ ਹੈ	ॢ -ऱ	क्रू kऱ्	र्कू rkऱ्
ਲੂ ਹੈ	ू -	कॢ kļ	क्रॄं rkļ

¶161 There is no vowel sign for \exists and \exists a.

¶162 The vowel sign for ξ i is generally displayed on the left side of its typographic cluster: $\overline{\sigma}$ ka • ξ i \rightarrow $\overline{\Phi}$ ki; $\overline{\sigma}$ ka • \overline{c} ta • ξ i \rightarrow $\overline{\overline{qc}}$ kti. When a cluster is displayed in non-joined form (either because such form is the most natural form or because it was requested by the use of a ZWNJ), this vowel sign is displayed just before the last consonant of the cluster: $\overline{\sigma}$ ka • $\overline{\sigma}$ ka • ξ i \rightarrow $\overline{\phi}\overline{\Phi}$ kki.

¶163 Special Clusters. In a cluster formed of \exists ra and a vocalic vowel, the \exists ra is displayed as a *repha* on the vowel letter:

 $\begin{array}{l} \overline{\mathsf{r}} \ \mathbf{ra} \bullet \overline{\mathfrak{R}} \ \mathbf{r} \to \overline{\mathfrak{R}} \ \mathbf{rr} < 0930, \, 0943 > \\ \overline{\mathsf{r}} \ \mathbf{ra} \bullet \overline{\mathfrak{R}} \ \mathbf{\bar{r}} \to \overline{\mathfrak{R}} \ \mathbf{rl} < 0930, \, 0944 > \\ \overline{\mathsf{r}} \ \mathbf{ra} \bullet \overline{\mathfrak{R}} \ \mathbf{\bar{r}} \to \overline{\mathfrak{R}} \ \mathbf{rl} < 0930, \, 0964 > \\ \overline{\mathsf{r}} \ \mathbf{ra} \bullet \overline{\mathfrak{R}} \ \mathbf{\bar{l}} \ \to \overline{\mathfrak{R}} \ \mathbf{rr} \ \mathbf{rl} < 0930, \, 0962 > \\ \overline{\mathsf{r}} \ \mathbf{ra} \bullet \overline{\mathfrak{R}} \ \mathbf{\bar{l}} \ \to \overline{\mathfrak{R}} \ \mathbf{rl} \ \mathbf{rl} < 0930, \, 0963 > \end{array}$

¶164 *Conjoined and Non-Joined Forms of Clusters.* The display of a cluster can be encouraged to be in conjoined form (with fallback to a non-joined form if necessary) by inserting a U+200D ZERO WIDTH JOINER after the consonant linker coded character.

¶165 The display of a cluster can be encouraged to be in non-joined form by inserting a U+200C ZERO WIDTH NON-JOINER after the consonant linker coded character.

 $x \cdot y$ x, 094D, yx, 094D, 200D, yx, 094D, 200C, yक ka • क ka \rightarrow क kkaक kaक ka • क ka \rightarrow क kkaक kaक ka • त ta \rightarrow क ktaक त kta

¶166 The display of vowel signs can be encouraged to be in full-joined form by inserting a U+200D ZERO WIDTH JOINER before them, and encouraged to be in conjoined form by inserting a U+200C ZERO WIDTH NON-JOINER before them.

 $x \cdot y$ x, y_{vs} $x, 200D, y_{vs}x, 200C, y_{vs}$ \forall ra • \exists u \rightarrow \forall rau or \exists rau \forall rau \forall ra • \exists u \rightarrow \forall rau or \exists rau \forall rau \forall ra • \exists u \rightarrow \forall rau or \exists rau \forall rau \forall ra • \exists u \rightarrow \forall rau \forall rau \forall ra • \exists u \rightarrow \forall rau \forall rau \forall ra • \exists u \forall rau \forall rau \forall ra • \exists bar or \exists rau \forall rau</

167 *Typographic Fragments.* The following typographic fragments can be represented as follows:

Full form Half form Re	epresented by
------------------------	---------------

क	क	<0915, 094D, 200D>
ख	ख	<0916, 094D, 200D>
ग	ī	<0917, 094D, 200D>
ਬ	ਦ	<0918, 094D, 200D>
च	ੁ	<091A, 094D, 200D>
ज	ন	<091C, 094D, 200D>
झ	ङ्	<091D, 094D, 200D>
স	5	<091E, 094D, 200D>
ण	σ	<0923, 094D, 200D>
त	7	<0924, 094D, 200D>
थ	શ	<0925, 094D, 200D>
ध	દ	<0927, 094D, 200D>
न	Ŧ	<0928, 094D, 200D>
ऩ	Ę	<0929, 094D, 200D>
प	τ	<092A, 094D, 200D>
দ্দ	ጥ	<092B, 094D, 200D>
ब	3	<092C, 094D, 200D>
भ	t	<092D, 094D, 200D>
म	Ŧ	<092E, 094D, 200D>
य	Σ	<092F, 094D, 200D>
ल	ल	<0932, 094D, 200D>
व	5	<0935, 094D, 200D>
গ	१	<0936, 094D, 200D>
ষ	ס	<0937, 094D, 200D>
स	रू	<0938, 094D, 200D>
क्ष	ද	<0915, 094D, 0937, 094D, 200D>
হা	ঝ	<091c, 094D, 091e, 094D, 200D>
त्त	त्त्	<0924, 094D, 0924, 094D, 200D>
त्र	त्	<0924, 094D, 0930, 094D, 200D>
প্স	श्र्	<0936, 094D, 0930, 094D, 200D>

¶168 do we want the 'half YA in post form'? subjoined tta, ttha, ga, dha, bha?

¶169 the last five entries are those in Table 9.4; however, that table is not meant to be exhaustive; should we have more?

 $\P170$ the last two entries actually show some ambiguity: we say that <930, 94D, 200D> is to be rendered by an eyelash ra (rule R5a, for compatibility with Unicode 2.0); so which takes precedence? The eyelash ra or the half form of the conjunct?

6.2 Bengla

¶171 The Bengla script is a North Indian script closely related to Devanagari. It is used to write the Bengla language primarily in the West Bengal state and in the nation of Bangladesh. It is also used to write Assamese in Assam and a number of other minority languages, such as Bishnupriya Manipuri, Daphla, Garo, Hallam, Khasi, Mizo, Munda, Naga, Rian, and Santali, in northeastern India.

¶172 Bengla is written horizontally from left to right, and words are separated by space.

¶173 In English, the term "Bengali" is commonly used instead of "Bengla", but the later term is preferred. This older term explains the formal name of the script and the character names in Unicode.

¶174 *Standards.* The Bengla block of the Unicode standard is based on ISCII-1988. See section <u>4.1, *The* ISCII *Standard*</u>, for more details on the relation to ISCII.

¶175 *Encoding.* The encoding of the Bengla script follows the consonant linking model. The consonant linker coded character is U+09CD \bigcirc BENGALI SIGN VIRAMA.

¶176 Letters. Bengla has all the basic letters presented in section <u>2.1, Letters</u>, except the consonants \overline{a} va and $\overline{\infty}$ la.

sho	rt	long		
letter	sign	lette	er	sign
অ a 0985	আ ā 098	₀ :t -ā	09BE	
₹i 0987 ि-i	09BF ঈ ī 098	≋ी -ī	09C0	
ঊ u 0989 ू -u	09C1 & ū 098	A ू -ū	09C2	
ଏē 098F (େ -ē	छ ०९८७ ঐ ai ०९९	ってつ-ai	09C8	
<u> ७ </u>	ō 09CB ঔ au 099	⁴ िो -a।	U 09CC	

¶177

sho	ort	long		
letter	sign	letter	sign	
∜ r 098B Ç	-ŗ 09C3 禊 ŗ 09E	80 Ç -] 09C4		
▶Ì 098C ୢ	-గ్గా 09E2 న్ని 🗍 09E	£1 ू −į 09E3		

¶178

	Voiceless					Voiced		Nasals				
		unas	pirated		aspir	ated		unas	pirated		aspirated	
velar	ক ka	0995	খ kha	0996	গ ga	0997	ঘ gha	0998	& 'na	0999		
palatal	চ ca	099A	ছ cha	099B	জ ja	099C	∛ jha	099D	ා ña	099E		
retroflex	ট ța	099F	ठे țha	09A0	ড ḍa	09A1	⊽ ḍha	09A2	ণ ņa	09A3		
dental	ত ta	09A4	থ tha	09A5	দ da	09A6	ধ dha	09A7	ন na	09A8		
labial	প pa	09AA	ফ pha	09AB	ব ba	09AC	ভ bha	09AD	ম ma	09AE		

Sonorants Fricatives

palatal 직 ya 09AF 차 śa 09B6 retroflex র ra 09B0 직 ṣa 09B7 dental ল la 09B2 지 sa 09B8

¶180

Others

হ ha 09B9

¶181

Signs

ঁ mঁ 0981 চন্দ্রবিন্দু candrabindu ং mं 0982 অনুস্বার anusbāra

ঃ ḥ 0983 বিসর্গ bisarga

¶182 The vowel nasalisation sign is known as the *candrabindu* (Bengla চন্দ্রবিন্দু candrabindu), and is encoded as U+0981 is BENGALI SIGN CANDRABINDU.

¶183 The consonant nasalization sign ং is known as the anusvara (Bengla অনুস্থার anusbāra), and is encoded as U+0982 ং BENGALI SIGN ANUSVARA.

¶184 The unvoiced aspiration sign ঃ is known as the *visarga* (Bengla বিসর্গ bisarga), and is encoded as U+0983 ঃ BENGALI SIGN VISARGA.

¶185 $\$ t, called *khanda ta*, was originally a reduced form of $\overline{\top}$ ta. This letterform has been pressed in service as a letter on its own. For example, it can be seen in the Bengla word $\[mathbb{N}\]$ me", and is in contrast with the letter $\[mathbb{o}\]$ ta: the Bengla word $\[mathbb{N}\]$ mata means "opinion". To accomodate that case, Unicode encodes U+09CE $\[mathbb{C}\]$ BENGALI LETTER KHANDA TA. In all known uses, $\[mathbb{C}\]$ tappears as the only letter of a typographic cluster.

¶186 ব ra and star ra are used to write Assamese, and are encoded as U+09F0 Text{TER RA WITH MIDDLE} DIAGONAL and U+09F1 Text{TER RA WITH LOWER DIAGONAL respectively.}

¶187 *Digits and Numbers.* Bengla has a full set of decimal digits, encoded at U+09E6 ° BENGALI DIGIT ZERO ... U+09EF & BENGALI DIGIT NINE.

¶188 todo something about the currency fractions 09F4..09F9

Punctuation and Signs. Bengla uses the danda | and double danda || encoded as U+0964 | DEVANAGARI DANDA and U+0965 || DEVANAGARI DOUBLE DANDA.

¶190 The avagraha is encoded as U+09BD ? BENGALI SIGN AVAGRAHA.

¶191 The rupee mark $_{\sim}$ is used todo and is encoded as U+09F2 $_{\sim}$ BENGALI RUPEE MARK. It is in addition to the rupee sign b

which is encoded as U+09F3 b BENGALI RUPEE SIGN.

¶192 \checkmark is used todo and is encoded as U+09FA \checkmark BENGALI ISSHAR.

¶193 \Box is use todo and is encoded as U+09FB \Box BENGALI GANDA MARK.

¶194 The coded character U+09D7 \bigcirc BENGALI AU LENGTH MARK does not correspond to an abstract character, i.e., it is not used for the representation of text. It is encoded to represent the right side of the vowel sign for \heartsuit au in isolation, and to provide a canonical decomposition for U+09CC \bigcirc BENGALI VOWEL SIGN AU.

¶195 *Atomicity of Vowel Letters.* Vowel letters are encoded atomically in Unicode, even if they can be analyzed visually as consisting of multiple parts. The table below shows the letters that can be analyzed, the single code point that should be used to represent them in text, and the sequence of code points resulting from analysis that should not be used.

To represent Use Do not use আ 0986 <0985, 09BE>

196 *Typographic Cluster Cores.* Here are some of the common typographic clusters made of two consonants which exhibit a fully-joined form:

ক ka • ক ka	→क kka
ক ka • ত ta	→ জ kta
ক ka • র ra	→ <u>ক</u> kra
ক ka • ষ ṣa	→ 꽈 kṣa
গ ga• ধ dha	→ র্দ্ব gdha
& 'na •	→ক 'nka
ঙ ṁa • গ ga	→ ॐ ṅga
জ ja • ঞ ña	→ ॼ jña
ঞ ña • চ ca	→ 🅸 ñca
ঞ ña • ছ cha	→ ञ्थ ñcha
ঞ ña • জ ja	→ अ ñja
ণ ṇa • ড ḍa	→ ♡ ņḍa
⊡ ta • ⊡ ta	→ उ tta
ত ta • থ tha	→খt'tha
ত ta • র ra	→ ত্র tra
দ da • দ da	$\rightarrow \overline{\mu} dda$
দ da • ধ dha	→ দ্ব d'dha
ন na • ত ta	→ ত nta
ব ba • দ da	→ ৹দ bda
ব ba•ধ dha	→ এধ bdha
*† śa • চ ca	→*চ śca

¶197 In the absence of a fully-joined form, the core of a typographic cluster is displayed in a conjoined form, where all the but the last consonant are in *half form* and the last consonant is in its full form.

¶198 The reduced form of the consonants খ গ ণ ন প শ স is obtained by dropping the vertical stem and slightly reducing their size: খ্ গ্ ণ্ ন্ প্ শ্ স্.

¶199 The reduced form of the other consonants is a smaller version of them.

ক ka•খ kha → ক্খ kkha ক ka•গ ga → ক্গ kga গ ga•ট ṭa → গ্ট gṭa গ ga•ড ḍa → গ্ড gḍa

 \mathbb{Q}_{200} Because every consonant has a half form, it is rare to encounter clusters in non-joined form, although they are in principle possible. This is displayed with a *hasant* sign below all but the last consonant:

ਹੋ t਼a • স sa → ऍস tsa

¶201 Not dealt with correctly yet: "clusters" of the form cons + hasant. May be this is the right place?

¶202 *Satellite Consonants. ¬* ra in initial position is a satellite consonant. It is then displayed as a slanted stoke on the top right side of the cluster core:

র ra•ত ta → র্ত rta র ra•ক ka•খ kha → ক্র্থ rkkha র ra•ম ma•আ ā → র্মা rmā র ra•থ tha•ঈī → র্থী rthī

¶203 A number of consonants in final position are satellite consonants:

য ya ক্য kya খ্য khya গ্য gya ঘ্য ghya চ্য cya

র ra	ক্র kra	গ্র gra	ਬ ghra	জ jra	ত্র tra
ল la	ৰু kla	ส gla	ੜਾ mla	絔 pla	ल्ल lla
ব ba	ৰু kba	জ্ব jba	ত্ব tba	ন্ব dba	ন্থ nba
ন na	গ্ন gna	ন nna	ত্ম tna	ন্শ sna	হ্ন hna
ম ma	ক্ম kma	অ মা	দ্ম dma	ন্ম nma	স্ম sma

 $\P204$ In a cluster formed of \exists ra and one of \exists ya, \exists ra, \exists la, \exists ba, \exists na, \exists ma, \exists ma, \exists ra adopts its satellite form and the second consonant adopts its full form. To encourage the initial \exists ra to have its full form and the second consonant to have its satellite form, one inserts a U+200D ZERO WIDTH JOINER between the \exists ra and the consonant linker.

त्र ra, 09CD, y त्र ra, 200D, 09CD, y

ৰ্য rya	র্য raya
र्त्र rra	ব্ৰ rara
र्ल rla	র্ল rala
র্ব rba	র⊲raba
र्न rna	র্ন rana
พ์ rma	র্ম rama

¶205 Satellite Vowels. Some typographic clusters with vowel signs have full-joined forms:

গ ga • উ u	→ ੴ gu
র ra • উ u	→ 👎 ru
র ra • ঊ ū	→ র rū
শ śa • উ u	→ 🖱 śu
হha∙উu	→ হ hu
হha•ঋrৢ	→ হা hr̥
ৰ ra • উ u	→ rau
ৰ ra • ঊ ū	→ ब raū
ৰ ra • উ u	→ 🗟 rau
ৰ ra • ঊ ū	→ क raū
ত ta • র ra • উ u	→ जू trau
ন na • ত ta • উ u	$u \rightarrow \overline{\mathfrak{V}}$ ntau
স sa • ত ta • উ u	→ ङू stau

¶206 Otherwise, the vowel sign attaches around the cluster core.

on ক ka on ক rka আ ā া -ā কা kā কা rkā ই i ি -i কি ki কি rki

ঈī	ी -ī	কী kī	কী rkī
উu	ू -u	কু ku	∳ rku
ঊ ū	ू -ū	কূ kū	∳ rkū
এ ē	<i>ि</i> -ē	কে kē	र्क rkē
ঐ ai	ৈ -ai	কৈ kai	ৰ্কৈ rkai
Зō	ো -ō	কো kō	ৰ্কো rkō
ઝ au	। (ॊ -au	। কৌ kau	। কৌ rkau
ঋŗ	् -rॢ	কृ kr	∳ rkŗ
ৠŢ	्र -ļ	ক <mark>ৃ</mark> kļ	∳ rkļ
ය l	ç -ŗ	কু k দ ়	কুঁ rk দ ়্
శ్య ļ		কু kļ	∲ rkļ

 $\P 207$ There is no vowel sign for অ a.

¶208 The vowel sign for \exists i is generally displayed on the left side of its typographic cluster: $\overline{\Phi}$ ka • \exists i \rightarrow $\overline{\Phi}$ ki; $\overline{\Phi}$ na • $\overline{2}$ ga • \exists i \rightarrow $\overline{2}$ n'ga'i. todo where when there is an hasant?

¶209 The vowel signs for $\[mathscrew=0.5]{\[mathscrew=$

¶210 The vowel sign for \Im \bar{o} is made of two parts, one that is generally displayed on the left side of the typographic cluster and one that is display on the right side: $\bar{\Phi}$ ka • \Im $\bar{o} \rightarrow \bar{\phi}$ ko; $\bar{\eta}$ na • $\hat{\eta}$ ga • \Im $\bar{o} \rightarrow \bar{\phi}$ $\hat{\eta}$ n'go. todo where when there is an hasant?

¶211 The vowel sign for \Im au is similarly displayed in two parts.

¶212 *Ya-phalaa.* The form that \exists ya adopts when last in a conjunct, \exists , is known as *ya-phalaa*. This form is also used, followed by the vowel sign $\dagger \cdot \bar{a}$, as a diacritic sign on the vowel letters \exists a and $\mathfrak{Q} = \bar{e}$; those combinations are used to transcribe accurately the vowel sounds [æ] and todo in loan words. Those combinations are represented by:

To repres	ent Use
অ্যা	<0985, 09CD, 09AF, 09BE>
এ্যা	<098F, 09CD, 09AF, 09BE>

¶213 Interaction of Repha and Ya-phalaa.. todo

¶214 *Conjoined and Non-Joined Forms of Clusters.* The display of a cluster can be encouraged to be in conjoined form (with fallback to a non-joined form if necessary) by inserting a U+200D ZERO WIDTH JOINER after the consonant linker coded character.

¶215 The display of a cluster can be encouraged to be in non-joined form by inserting a U+200C ZERO WIDTH NON-JOINER after the consonant linker coded character.

$x \bullet y$	x, 09CD, y	, x, 09CD, 200D, y	<i>x</i> , 09CD, 200C, <i>y</i>
ক ka • ক ka	→क kka	ক্ক kka	ক্ক kka
ক ka • ত ta	→ জ kta	ক্ত kta	ক্ত kta
	\rightarrow		

ক ka • খ kha ক্ kkha ক্ খ kkha ক্ খ kkha

¶216 The display of vowel signs can be encouraged to be in full-joined form by inserting a U+200D ZERO WIDTH JOINER before them, and encouraged to be in conjoined form by inserting a U+200C ZERO WIDTH NON-JOINER before them.

$x \bullet y$	x, y_{vs}	x , 200D, y_{vs}	$x, 200C, y_{vs}$
গ ga • উ u	→ છ gau or ગૂં gau	ී gau	গু gau
র ra • উ u	→ क rau or त्रू rau	ৰু rau	त्रू rau
র ra • ঊ ū	→ क raū or রূ raū	র raū	রূ raū
শ śa • উ u	→ ৺ śau or শু śau	🖱 śau	™ śau
হha•উu	\rightarrow $\overline{2}$ hau or $\overline{2}$ hau	হ hau	रू hau
হha•ঋrৢ	\rightarrow रा har or र har	হ্য har	र्ड् har
ৰ ra • উ u	→ क rau or बू rau	ৰু rau	बू rau
ৰ ra • ঊ ū	→ ক raū or ৰূ raū	ৰা raū	ৰূ raū
র ra • উ u	→ বৃৃ rau or বু rau	ৰু rau	जू rau
ৰ ra • ঊ ū	→ ज raū or बूraū	কা raū	बू raū
ত ta • র ra • উ u	$\rightarrow \overline{\square}$ trau or $\overline{\square}$ trau	ক্র trau	जू trau
ন na • ত ta • উ u	$\rightarrow \overline{\mathfrak{G}}$ ntau or $\overline{\mathfrak{F}}$ ntau	। छ ntau	ন্তু nt <i>a</i> u
স sa • ত ta • উ u	→ ख stau or छू stau	छ stau	र्डू stau

¶217 Typographic Fragments. The following typographic fragments be represented as follows:

Full form Half form Represented by 각 각 <0996, 09CD, 200D> more...

¶218 The left side of the vowel signs $1 - \bar{0}$ and $1 - \bar{0}$ are represented in isolation by $<0020, 09C7 > \rightarrow 0$

¶219 The right side of the vowel sign $(\uparrow -\bar{o} \text{ is represented in isolation by } <0020, 09BE> \rightarrow \uparrow$.

¶220 The right side of the vowel sign (\bigcirc] -au is represented in isolation by <0020, 09D7> →]. The coded character U+09D7 \bigcirc] BENGALI AU LENGTH MARK does not correspond to an abstract character and is encoded solely to provide a canonical decomposition for U+09CC (\bigcirc] BENGALI VOWEL SIGN AU, and to represent its right side in isolation.

6.3 Gurmukhi

¶221 The Gurmukhi script is a North Indian script used to write the Punjabi (or Panjabi) language of the Punjab state of India. Gurmukhi, which literally means "proceeding from the mouth of the Guru," is attributed to Angad, the second Sikh Guru (1504 - 1552 CE). It is derived from an older script called Landa, and is closely related to Devanagari structurally. The script is closely associated with Sikh and Sikhism, but it is used on an everyday basis in East Punjab. (West Punjab, now in Pakistan, uses the Arabic script.)

¶222 Gurmukhi is written horizontally from left to right, and words are separated by space.

¶223 *Standards.* The Gurmukhi block of the Unicode standard is based on ISCII-1988. See section <u>4.1, *The* ISCII *Standard*</u>, for more details on the relation to ISCII.

¶224 *Encoding.* The encoding of the Gurmukhi script follows the consonant linking model. The consonant linker coded character of Gurmukhi is U+0A4D \bigcirc GURMUKHI SIGN VIRAMA.

¶225 *Letters.* Gurmukhi has all the basic letters presented in section <u>2.1, Letters</u>, except the consonant $\overline{\P}$ sa, and the vocalic vowels $\overline{\Re}$ r, $\overline{\Re}$ r, $\overline{\Re}$ l, and $\overline{\varrho}$ l.

	short		long				
letter	sign	letter	sign				
ਅ a 0A	05	ਆ ā 0A0	6 • ⊺ -ā 0A	3E			
ਇi0A	07 f -i 0A3F	ਈī 0A0	8 ी -ī 0A	40			
ਉ u 0A	09 ౖ -u 0A41	ਊū 0A0.	A ૄ -ū 0A	42			
ਏ ē 0A	0F े -ē 0A3E	ਐ ai 0A1	0 ै -ai 0A	48			
ਓ ō 0A13 ੋ -ō 0A4B ਔ au 0A14 ੌ -au 0A4C							
	Voice	less	Ve	biced		Na	sals
	unaspirated	aspirated	unaspirated	aspi	rated		
velar	ਕ ka 0A15 ਖ	kha 0A16	ਗ ga 0A17	ਘ gha	0A18	ਙ na	0A19
palatal	ਚ ca 0A1A ਛ	cha 0A1B	ਜja 0A1E	ਝ jha	0A1D	ਞ ña	0A1E
retrofley	ĸさța 0A1F る	țha 0A20	ਡ ḍa 0A21	ਢ ḍha	0A22	ਣ ņa	0A23

dental ਤ ta 0A24 ਥ tha 0A25 ਦ da 0A26 ਧ dha 0A27 ਨ na 0A28

labial ਪ pa 0A2A ਫ pha 0A2B ਬ ba 0A2C ਭ bha 0A2D ਮ ma 0A2E

¶226 쿸 ra is encoded as U+0A5C 쿸 GURMUKHI LETTER RRA.

¶227 The vowel nasalisation sign ざ is known as the adak bindi, and is encoded as U+0A01 ざ GURMUKHI SIGN ADAK BINDI

¶228 Two different marks are used for the consonant nasalization sign. $\dot{\odot}$ is known as the *bindi* and is encoded as U+0A02 $\dot{\odot}$ GURMUKHI SIGN BINDI. $\dot{\odot}$ is known as the *tippi* and is encoded as U+0A70 $\ddot{\odot}$ GURMUKHI TIPPI. Present practice is to use bindi only after the vowel signs $\neg -\bar{a}$, $\uparrow -\bar{c}$, $\ddot{\bigcirc} -ai$, $\ddot{\bigcirc} -\bar{o}$, $\ddot{\bigcirc} -au$, and after the vowel letters $\frac{1}{2}$ u and $\frac{1}{2}$ \bar{u} , and to use tippi in the other cases. However, older texts may follow different conventions.

¶229 The unvoiced aspiration sign of is known as the visarga, and is encoded as U+0A03 of GURMUKHI SIGN VISARGA.

¶230 The additional consonants with nukta (called *pairin bindi*, literally, "with a dot in the foot," in Punjabi) are primarily used to differentiate Urdu or Persian loan words. They include π sa and $\aleph \underline{l}a$, but do not include π ra which is genuinely Punjabi.

¶231 The *yakash* sign ु probably originated from a subjoined form of ज ya. However, because its usage is relatively rare and is not entirely predictable, it is encoded as a separate character, U+0A75 ु GURMUKHI SIGN YAKASH. This character should occur after the consonant to which it attaches and before any vowel sign.

¶232 *Digits and Numbers.* Gurmukhi has a full set of decimal digits, encoded at U+0A66 ° GURMUKHI DIGIT ZERO ... U+0A6F せGURMUKHI DIGIT NINE.

1233 *Punctuation and Signs.* Gurmukhi uses the danda | and double danda || encoded as U+0964 | DEVANAGARI DANDA and U+0965 || DEVANAGARI DOUBLE DANDA.

¶234 The sign indicates that the following consonant is geminate. It is encoded as U+0A71 GURMUKHI ADDAK. For example, the word ਪੱਰਾ paga, "turban", can be represented with the sequence <U+0A2A, U+0A71, U+0A17>. todo

¶235 In dictionnaries and similar materials, words starting with \mathfrak{P} u or \mathfrak{P} \overline{u} are sorted together, under the heading \mathfrak{P} . For the purpose of representing this heading, Unicode encodes U+0A73 \mathfrak{P} GURMUKHI URA. This character should not be used to represent \mathfrak{P} u and \mathfrak{P} \overline{u} . Similarly, words starting with \mathfrak{P} i or \mathfrak{P} \overline{i} are sorted together, under the heading \mathfrak{P} . Again, only for the purpose of representing that heading, Unicode encodes U+0A72 \mathfrak{P} GURMUKHI IRI.

¶236 The religious symbol *khanda* sometimes used in Gurmukhi texts is encoded at U+262C \oplus ADI SHAKTI in the Miscellaneous Symbols block. U+0A74 \oplus GURMUKHI EK ONKAR, which is also a religious symbol, can have different presentation forms, which do not change its meaning. The font used in the code charts shows a highly stylized form; simpler forms look like the digit one, followed by a sign based on ura, along with a long upper tail.

1237 In older texts, the *udaat* sign indicates a high tone. It is encoded as U+0A51 GURMUKHI SIGN UDAAT. This character should occur after the consonant to which is attaches and before any vowel sign.

¶238 *Tones.* Tones. The Punjabi language is tonal, but the Gurmukhi script does not contain any specific signs to indicate tones. Instead, the voiced aspirates (\mathfrak{U} gha, \mathfrak{F} jha, \mathfrak{E} dha, \mathfrak{T} dha, \mathfrak{T} bha) and the letter \mathfrak{T} ha combine consonantal and tonal functions.

1239 *Atomicity of Vowel Letters.* Vowel letters are encoded atomically in Unicode, even if they can be analyzed visually as consisting of multiple parts. The table below shows the letters that can be analyzed, the single code point that should be used to represent them in text, and the sequence of code points resulting from analysis that should not be used.

To represent	Use	Do not use
ਆ	0A06	<0A05, 0A3E>
ਇ	0A07	<0A72, 0A3F>
ਈ	0A08	<0A72, 0A40>
ਉ	0A09	<0A73, 0A42>
ਊ	0A0A	<0A73, 0A42>
ਏ	0A0F	<0A72, 0A47>
ਐ	0A10	<0A05, 0A48>
ଟ	0A13	<0A73, 0A4B>
ਔ	0A14	<0A05, 0A4C>

¶240 *Typographic Cluster Cores.* Gurmukhi is unusual in that there are no fully joined forms. In addition, in modern usage, the core of a typographic cluster is always a single consonant.

¶241 *Satellite Consonants.* In modern Punjabi, the consonants \exists ra, \exists ha and \exists va are satellite consonants in final position, and take a subjoined (*pairin*) form.

ਪ pa • ਰ ra $\rightarrow \chi$ pra ਮ ma • ਹ ha $\rightarrow H$ mha ਦ da • ਵ va $\rightarrow \xi$ dva

¶242 In modern Punjabi, the consonant ज ya is a satellite consonant in final position, and take a postjoined (addha) form.

ਦ da • ਯ ya → ਦ੍ਯ dya

¶243 In older texts, one can find more consonants with a pairin form:

¶244 In older texts, one can find more consonants with an addha form:

¶245 Older texts also exhibit another feature which is not found in modern Punjabi: instead of showing the second consonant in satellite form, one can see the first consonant in a satellite form. For \exists ra, this satellite form is similar to the Devanagari *reph*. For the other consonants, this is a half form, generally formed by dropping the vertical stem. This conjoined form can be encouraged by inserting a U+200D ZERO WIDTH JOINER after the consonant linker. Inserting a U+200C ZERO WIDTH NON-JOINER encourages a non-joined form:

 $x \cdot y$ x, 0A4D, yx, 0A4D, ZWJ, yx, 0A4D, ZWNJ, y $\exists ra \cdot \exists va \rightarrow \overline{q} rva$ $\overline{q} rva$ $\overline{q} \exists rva$ $\exists ra \cdot \exists va \rightarrow \overline{q} rva$ $\overline{q} rva$ $\overline{q} \exists rva$ $\exists ra \cdot \exists va \rightarrow \overline{q} sva$ $\overline{q} sva$ $\overline{q} \exists sva$

¶246 Satellite Vowels. The vowel signs attach around the cluster core.

아 ਕ ka ਆ ā ਾ -ā ਕਾ kā ਇ i ি -i ਕਿ ki ਈ ī ੀ -ī ਕੀ kī ਉ u ੁ -u ੁ ku ਊ ū ੂ -ū ਤੂ kū ਏ ē ੇ -ē ਕੇ kē ਐ ai ੈ -ai ਕੇ kai ਓ ō 존 -ō ਕੋ kō ਔ au ~ -au ~ kau

¶247 There is no vowel sign for \mathcal{W} a.

¶248 In older texts, such as the *Sri Guru Granth Sahib* (the Sikh holy book), one can find typographic clusters with a vowel sign attached to a vowel letter, or with two vowel signs attached to a consonant. The most common cases are \bigcirc -u attached to \boxdot \overline{o} , as in ত্রিਮার্য or burnahā, and both the vowel signs \bigcirc - \overline{o} and \bigcirc -u attached to a consonant, as in ব্রাधिন goubida; this is used to indicate the metrical shortening of/ \overline{o} / or the lengthening of/u/ depending on the context. Other combinations are attested as well, such as \boxed{d} _biana.

¶249 Because of the combining classes of the characters U+0A4B ⊂ GURMUKHI VOWEL SIGN OO and U+0A41 ⊂ GURMUKHI VOWEL SIGN U, the sequences *<consonant*, U+0A4B, U+0A41> and *<consonant*, U+0A41, U+0A4B> are not canonically equivalent. To avoid ambiguity in representation, the first sequence, with U+0A4B before U+0A41, should be used in such cases. More generally, when a consonant or independent vowel is modified by multiple vowel signs, the sequence ov vowel signs in the underlying representation of the text should be: left, top, bottom, right.

1250 *Conjoined and Non-Joined Forms of Clusters.* The display of a cluster can be encouraged to be in non-joined form by inserting a U+200C ZERO WIDTH NON-JOINER after the consonant linker coded character.

 $x \cdot y$ x, 0A4D, yx, 0A4D, 200C, y $a \text{ ka} \cdot \overline{\sigma} \text{ ra } \rightarrow \overline{q} \text{ kra}$ ब्र $\overline{\sigma} \text{ kra}$ $a \text{ ka} \cdot \overline{w} \text{ ya} \rightarrow \overline{a} \overline{w} \text{ kya}$ ब् \overline{w} kya

6.4 Gujarati

¶251 The Gujarati script is a North Indian script closely related to Devanagari. It is most obviously distinguished from Devanagari by not having a horizontal bar for its letterforms, a characteristic of the older Kaithi script to which Gujarati is related. The Gujarati script is used to write the Gujarati language of the Gujarat state in India.

¶252 Gujarati is written horizontally from left to right, and words are separated by space. The consonant letters often have a vertical stem on their right side, which is graphically and historically related to the sign denoting the inherent /a/ vowel.

¶253 *Standards.* The Gujarati block of the Unicode standard is based on ISCII-1988. See section <u>4.1, *The* ISCII *Standard*</u>, for more details on the relation to ISCII.

¶254 *Encoding*. The encoding of the Gujarati script follows the consonant linking model. The consonant linker coded character of Gujarati is U+0ACD Q GUJARATI SIGN VIRAMA.

^{¶255} Letters. Gujarati has all the basic letters presented in section <u>2.1, Letters</u>.

¶256 The vowel nasalisation sign is encoded as U+0A81 GUJARATI SIGN CANDRABINDU.

 257 The consonant nasalization sign \dot{o} is known as the *ansuvara* in Gujarati, and is--> encoded as U+0A82 \dot{o} GUJARATI SIGN ANUSVARA.

¶258 The unvoiced aspiration sign : is encoded as U+0A83 : GUJARATI SIGN VISARGA.

¶260 해 ŏ is used todo and is encoded as U+0A91 해 GUJARATI VOWEL CANDRA O and U+0AC9 ``I GUJARATI VOWEL SIGN CANDRA O.

¶261 *Digits and Numbers.* Gujarati has a full set of digits, used to write decimal numbers. Those are encoded at U+0AE6 ∘ GUJARATI DIGIT ZERO .. U+0AEF ∈ GUJARATI DIGIT NINE.

1262 Punctuation and Signs. Gujarati uses the danda | and double danda || encoded as U+0964 | DEVANAGARI DANDA and U+0965 || DEVANAGARI DOUBLE DANDA.

¶263 The *avagraha* is encoded as U+0ABD 5 GUJARATI SIGN AVAGRAHA.

¶264 The *rupee sign* is encoded as U+0AF1 3₀ GUJARATI RUPEE SIGN.

¶265 *Atomicity of Vowel Letters.* Vowel letters are encoded atomically in Unicode, even if they can be analyzed visually as consisting of multiple parts. The table below shows the letters that can be analyzed, the single code point that should be used to represent them in text, and the sequence of code points resulting from analysis that should not be used.

To represent	Use	Do not use
આ	0A86	<0A85, 0ABE>
ઍ	0A8D	<0A85, 0AC5>
એ	0A8F	<0A85, 0AC7>
એ	0A90	<0A85, 0AC8>
ઑ	0A91	<0A85, 0AC9>
ઓ	0A93	<0A85, 0ACB>
ઓ	0A94	<0A85, 0ACC>

¶266 *Typographic Cluster Cores.* Here are some of the common typographic clusters made of two consonants which exhibit a fully-joined form:

8 ka • 8 ka	→ SS kka
s ka•∀ șa	→ क्ष kṣa
5 na • 8 ka	→ş; 'nka
ડ na • ગ ga	→ ફ nga
ਤ na • ध gha	→ ਝ ngha
S na • ਮ ma	→ ਝੂ nma
Ƴ ja • ઞ ña	→ হা jña
ઞ ña • ચ ca	→ ਝ ñca
ઞ ña • જ ja	→ ਖ਼ ñja
sțs•sța	→੬ țța
εța•δțha	→ζδ țțha
5 ḍa • 5 ḍa	→ ૬ વ઼ંવેa
5 ḍa • ઢ ḍha	→ਙ d਼dha
a dha • a dha	→& ḍhḍha
d ta•d ta	→ त tta
દ da • ગ ga	→દ્ગ dga
६ da • ध gha	→ ट्ध dgha
ε da • ε da	\rightarrow ६ dda
દ da • ધ dha	→ & d'dha

 \mathcal{E} da • 여 na \rightarrow ፉ dna \mathcal{E} da • 여 ba $\rightarrow \mathcal{E}$ 여 dba \mathcal{E} da • 여 bha $\rightarrow \mathcal{E}$ 여 dbha \mathcal{E} da • 권 ma \rightarrow 놓 dma 여 na • 여 na \rightarrow 치 nna

¶267 In the absence of a fully-joined form, the core of a typographic cluster is displayed in a conjoined form, where all the but last consoant are in half form and the last consonant is in full form.

 \P^{268} The half form of \$ ka is \$; the half form of \$ ja is \$; the half form of \$ pha is \$. The half form of the consonants with a vertical stem on the right is obtained by dropping this vertical stem.

 $s ka \cdot \delta tha \rightarrow s\delta ktha$ $v ja \cdot s ka \rightarrow vs jka$ $s pha \cdot s ka \rightarrow ss phka$ $v kha \cdot s ka \rightarrow vs khka$ $v kha \cdot \varepsilon da \rightarrow v\varepsilon khda$ $cta \cdot s ka \rightarrow cs tka$ $cta \cdot H ma \rightarrow cH tma$

¶269 The remaining consonants do not have a half form. When they are present in non-final position in a cluster core, they non-joined form of the cluster is displayed, with a *halant* sign below all but the last consonant.

ଓ cha • \$ ka \rightarrow ୍ଞ୍ s chka

¶270 *Satellite Consonants.* ? ra in initial position is a satellite consonant. It is then displayed as a small hook, called a *repha*, on the top right side of the cluster, including any satellite vowels:

૨ ra • ઙ ka → ઙ rka ૨ ra • ઙ ka • ખ kha → ડર્ખ rkkha

¶271 Where should the reph be placed if there is an a sign? clean up the examples above.

 \forall kha • २ ra \rightarrow \forall khra \forall kha • २ ra \rightarrow \forall khra \Diamond ña • २ ra \rightarrow \Diamond ñra \Im jha • २ ra \rightarrow \Im jhra δ țha • २ ra \rightarrow \Diamond țhra

¶273 Satellite Vowels. The vowel signs attach around the cluster core.

on S ka on S rka

ਆ ā ਈ i	ा -ā ि -i	ङ। kā s ki	ร์โrkā Is rki
ឋī	ी -ī	डी kī	ร์โ rkī
βu	ु-u	şku	§rku
। १	ू -ū	ş kū	§ rkū
એ ĕ	े-ĕ	`š kĕ	\ Šrkĕ
એ ē	े <i>-</i> ē	€ 8kē	\$rkē
એ ai	े-ai	Škai	Frkai
ઑ ŏ	ॉ -ŏ	કૉ kŏ	ร์โ rkŏ
ઓ ō	ो -ō	કો kō	ร์โ rkō
ઓ au	ाेो -au	ા કો kau	કોં rkau
жŗ	် -ပႆ	ş kŗ	§ rkŗ
>te ŗ	<i>ॄ</i> -ļ	ş kļ	ş rkļ
ო ļ	ૣ ૣ	ş ka	§ rką
ဖ်၂	្ត ្ត	s ką	ş rka

¶274 There is no vowel sign for \Im a.

¶275 The vowel sign for \mathcal{E} i is generally displayed on the left side of its typographic cluster: \mathfrak{S} ka • \mathcal{E} i \rightarrow \mathfrak{S} ki; \mathfrak{I} na • \mathfrak{I} ga • \mathcal{E} i \rightarrow $\mathfrak{G}\mathfrak{I}$ n'gi.

¶276 todo where when there is hasant?.

1277 *Conjoined and Non-Joined Forms of Clusters.* The display of a cluster can be encouraged to be in conjoined form (with fallback to a non-joined form if necessary) by inserting a U+200D ZERO WIDTH JOINER after the consonant linker coded character.

¶278 The display of a cluster can be encouraged to be in non-joined form by inserting a U+200C ZERO WIDTH NON-JOINER after the consonant linker coded character.

$x \bullet y$	<i>x</i> , 0ACD, <i>y</i>	, x, 0ACD, 200D, y	<i>x</i> , 0ACD, 200C, <i>y</i>
ઞ ña • ચ ca →	→ ਝ ñca	ગ્સ ñca	ઞ્ચ ñca
ɛ da • н ma →	→ \$I dma	દ્મ dma	દ્મ dma

1279 *Typographic Fragments.* The following typographic fragments can be represented as follows:

|--|

ક	ક	<0A95, 0ACD, 200D>
Ч	Ն	<0A96, 0ACD, 200D>
၁၂	С	<0A97, 0ACD, 200D>
ย	٤	<0A98, 0ACD, 200D>
ચ	ર	<0A9A, 0ACD, 200D>
୪	જ	<0A9C, 0ACD, 200D>
거	Ъ	<0A9E, 0ACD, 200D>
ତା	g	<0AA3, 0ACD, 200D>

d	С	<0AA4, 0ACD, 200D>
થ	૭	<0AA5, 0ACD, 200D>
ધ	٤	<0AA7, 0ACD, 200D>
ব	σ	<0AA8, 0ACD, 200D>
ч	ι	<0AAA, 0ACD, 200D>
ફ	ફ	<0AAB, 0ACD, 200D>
어	ዮ	<0AAC, 0ACD, 200D>
ભ	୯	<0AAD, 0ACD, 200D>
ы	9	<0AAE, 0ACD, 200D>
ચ	5	<0AAF, 0ACD, 200D>
А	C-	<0AB2, 0ACD, 200D>
ഗ	U	<0AB3, 0ACD, 200D>
q	с	<0AB5, 0ACD, 200D>
থ	ঀ	<0AB6, 0ACD, 200D>
ы	ъ	<0AB7, 0ACD, 200D>
਼	5	<0AB8, 0ACD, 200D>
ଝା	Ş	<0A95, 0ACD, 0AB7, 0ACD, 200D>
হা	ş	<0A9C, 0ACD, 0A9e, 0ACD, 200D>
Я	>	<0AA4, 0ACD, 0AB0, 0ACD, 200D>
烙	٩	<0AB6, 0ACD, 0AB0, 0ACD, 200D>

6.5 Oriya

¶280 The Oriya script is a North Indian script that is structurally similar to Devanagari, but with semicircular lines at the top of most letters instead of the straight horizontal bars of Devanagari. The actual shapes of the letters, particularly for vowel signs, show similarities to Tamil. The Oriya script is used to write the Oriya language of the Orissa state in India, as well as minority languages such as Khondi and Santali.

§281 *Standards.* The Oriya block of the Unicode standard is based on ISCII-1988. See section <u>4.1, *The* ISCII *Standard*</u>, for more details on the relation to ISCII.

[282 *Encoding.* The encoding of the Oriya script follows the consonant linking model. The consonant linker coded character of Oriya is U+0B4D \bigcirc ORIYA SIGN VIRAMA.

¶283 Letters. Oriya has the full complement of basic letters presented in section 2.1, Letters.

¶284 The vowel nasalisation sign is known as the *candrabindu* in Oriya, and is encoded as U+0B01 ORIYA SIGN CANDRABINDU.

¶285 The consonant nasalization sign ⁽⁾° is known as the *ansuvara* in Oriya, and is encoded as U+0B02 ⁽⁾° ORIYA SIGN ANUSVARA.

¶286 The unvoiced aspiration sign 38 is known as the *visarga* in Oriya, and is encoded as U+0B03 38 ORIYA SIGN VISARGA.

¶287 When Sanskrit words written in Devanagari using the letter \overline{q} va are borrowed in the Oriya language, they are most often written with the Oriya letter \overline{q} ba, to match the corresponding change in pronunciation: Sanskrit $\overline{q}\overline{q}$ vana becomes Oriya $\overline{q}\overline{q}$ bana. For the purpose of accurate transliteration, the letter \overline{q} va emerged, and is used to transliterate Devanagari \overline{q} va, while the letter \overline{q} ba is used only for the transliteration of Devanagari \overline{q} ba. This letter is not in common use today, but can be found in academic and technical texts. It is encoded as U+0B35 \overline{q} ORIYA LETTER VA. In older sources, this character is rendered by a wide range glyphs: \overline{q} , \overline{q} , \overline{q} , \overline{q} , \overline{q} , \overline{q} (the latter being a modern rendering of the previous two).

¶288 the representative glyph uses a dot inside the circle and the text gives a variant with a ring above the circle. The font utkal.ttf, which is relatively recent uses a dot above, and does not even have a glyph matching the representative glyph. Is there is a disconnect somewhere?

¶289 To write Perso-Arabic or English loan words with a [w], the letter \mathfrak{G} is sometimes used. It appears to have originally been devised as a ligature of \mathfrak{G} \overline{o} and \mathfrak{G} ba, but because ligatures of independent vowels and consonants are not normally used in Oriya, this letter has been encoded as U+0B71 \mathfrak{G} ORIYA LETTER WA, without a canonical decomposition. It is used initially in words or orthographic syllables to represent the foreign consonant; as a native semivowel, virama + ba is used because that is historically accurate. This character has a wide range of glyphs: \mathfrak{G} , \mathfrak{G} , \mathfrak{G} .

 \mathbb{P}^{290} I can't make complete sense of the "as a native semivowel" part. Does it mean to that in native words, the *sound* [w] does not occur initially, and is written by BA? Also, this business of not doing vowel letter + virama + consonant is no longer accurate (bengla vowel + virama + ya), right? if so, how do we position U+B71?

¶291 & ya is used for [...] and is encoded as U+0B5F & ORIYA LETTER YYA.

¶292 *Digits and Numbers.* Oriya has a full set of decimal digits, encoded at U+0B66 ○ ORIYA DIGIT ZERO ... U+0B6F ♥ ORIYA DIGIT NINE.

1293 Punctuation and Signs. Oriya uses the danda | and double danda || encoded as U+0964 | DEVANAGARI DANDA and U+0965 || DEVANAGARI DOUBLE DANDA.

¶294 The *avagraha* is encoded as U+0B3D € ORIYA SIGN AVAGRAHA.

¶295 ♀ is placed before names of persons who are deceased and is encoded as U+0B70 ♀ ORIYA ISSHAR.

¶296 *Atomicity of Vowel Letters.* Vowel letters are encoded atomically in Unicode, even if they can be analyzed visually as consisting of multiple parts. The table below shows the letters that can be analyzed, the single code point that should be used to represent them in text, and the sequence of code points resulting from analysis that should not be used.

To represent	Use	Do not use
ଆ	0B06	<0B05, 0B3E>
ঝ	0B10	<0B0F, 0B57>
ଔ	0B14	<0B13, 0B57>

¶297 *Typographic Cluster Cores.* Here are some of the common typographic clusters made of two consonants which exhibi a fully joined form:

 $\Re ka \bullet \otimes sa \longrightarrow \otimes ksa$

ଙ ṅa • କ ka → ଙ୍କ ṅka ଙ ṅa • ଖ kha → ଙ୍ଗ ṅkha \mathfrak{C} na • ଗ ga $\rightarrow \mathfrak{G}$ nga ଚ ca • ଚ ca → ଚ୍ଚ cca \Im ca • \Im cha \rightarrow \Im ccha କ ja • କ ja → ଜ jja କ ja • ଞ ña → ଜଞ୍ jña ଞ ña • ଚ ca → ଞଁ ñca \mathfrak{B} ña • \mathfrak{A} cha $\rightarrow \mathfrak{A}$ ñcha ଞ ña • କ ja → 🗞 ñja \mathfrak{B} ña • \mathfrak{C} jha $\rightarrow \mathfrak{B}$ ñjha ଟ ta • ଟ ta → ଟ tta ଶ na • $\mathfrak{S} da \rightarrow \mathfrak{S} nda$ ଶ na • $\[mathbb{Q}\]$ dha \rightarrow ଷ ndha ତ ta • କ ka \rightarrow କୁ tka ତ ta • ପ pa → Q tpa ତ ta • ସ sa →ସ tsa ଦ da • ଦ da → ଦ୍ଦ dda Ω da • ∂ dha $\rightarrow \Omega$ d'dha ଧ dha • ୟ ýa → ଧ୍ୟ dhýa \Re na • 2 tha \rightarrow \Re ntha \Re na • \Im da $\rightarrow \Re$ nda ନ na • ଧ dha \rightarrow \bigotimes ndha ବ ba • ଦ da \rightarrow ବ୍ଦ bda ମ ma • ପ pa → প mpa 위 ma • 안 pha \rightarrow 위 mpha ଶ śa • ଛ cha \rightarrow ଗ୍ଟ ścha \mathfrak{A} sa • \mathfrak{O} pa $\rightarrow \mathfrak{A}$ spa \mathfrak{A} sa • \mathfrak{A} pha $\rightarrow \mathfrak{A}$ spha

¶298 In the absence of a fully-joined form, the core of a typographic cluster is displayed in conjoined form, where all but the first consonant are in half form and the last consonant is in its full form.

¶299 The reduced forms of @ dha, \odot ta, & tha, ? ma and @ va are dha, ta, ta, ma and va respectively.

¶300 The reduced forms of \Im ka, \Im ja, \Im na, \Re la and \Re la is obtained by dropping their loop: ${}_{ \Im}$ ka, ${}_{ \Im}$ ja, ${}_{ \Im}$ na, ${}_{ \Re}$ la and ${}_{ \Re}$ la.

¶301 The reduced form of remaining consonants is a smaller version of them.

¶302 [Examples]

¶303 [visible halant] is displayed with a *halant* sign under it:

ଖ kha • କ ka → ଖ୍କ khka

¶304 *Satellite Consonants.* \Im ra is a satellite consonant in initial position. It is then displayed as a small mark, called a *repha* on the top right side of the cluster core.

ର ra • କ ka \rightarrow କି rka

¶305 බ ra is a satellite consonant in final position todo is this correct? examples..

¶306 Where should the reph be placed if there is an a sign? where should the reph be placed if there is a postjoined ya? clean up the examples above.

¶ $307 \, \Theta$ ta is a satellite consonant in initial position.

¶308 \bigcirc ra is a satellite consonant in final position.

¶ $309 \, \Im$ ya is a satellite consonant in final position.

¶310 In a cluster formed of \Im ra and one of the consonants above, the initial \Im ra adopts its satellite form and the second consonant adopts its full form. To encourage the initial \Im ra to have its full form and the second consonant to have its satellite form, one inserts a U+200D ZERO WIDTH JOINER between the \Im ra and the consonant linker.

ର ra, 0B4D, y ର ra, 200D, 0B4D, y

ế rņa	ର୍ଣ୍ଣ raṇa
ର୍ତ rta	ର୍ଭ rata
ର୍ନ rna	ର୍କ rana
ର୍ବ rba	ର୍ବ raba
ଭ୍ରି rbha	ର୍ଣ୍ଣ rabha
ର୍ମ rma	ର୍କ rama
ณิ์ rya	ରଯ୍ raya
ର୍ର rra	ର୍ର rara
ณ์ rla	ରୁ rala

ଳ rļa କ୍ଲ raļa

¶311 Satellite Vowels. The vowel sign attaches arond the cluster core:

		on କ ka	on କି rka
ଆ ā	∶l-ā	କା kā	ର୍କା rkā
ଇ i	ि -i	କି ki	ର୍କି rki
ଈ ī	ी-ī	କୀ kī	ର୍କୀrkī
ର u	ू -u	କୁ ku	ର୍କୁ rku
ଊū	्रू -ū	କୂ kū	ର୍କୁ rkū
√ē	6≎-ē	କେ kē	ର୍କେ rkē
থী ai	6 ai	କୈ kai	ର୍କି rkai
Зō	60 -ō	କୋ kō	ର୍କୋ rkō
ଔ au	ේ -au	। କୌ kau	କୌ rkau
ର r	©; -ŗ	କୃ kr	ର୍କ <mark>ୁ</mark> rkr

¶312 There is no vowel sign for \mathfrak{A} a, \mathfrak{R} , \mathfrak{r} , \mathfrak{E}], \mathfrak{L} , \mathfrak{k} , \mathfrak{k} .

¶313 *Typographic Fragments.* The following half forms can be represented as follows:

Full form Half form Represented by

କ	୩	<0020, 200D, 0B4D, 0B15>
ଜ	ហ	<0020, 200D, 0B4D, 0B1C>
ନ	a	<0020, 200D, 0B4D, 0B28>
ବ	ч	<0020, 200D, 0B4D, 0B2C>
ଲ	m	<0020, 200D, 0B4D, 0B32>
ଳ	m	<0020, 200D, 0B4D, 0B33>
ବ	પ	<0020, 200D, 0B4D, 0B35>

 $\P314$ The list above comes from the TDIL document, p30, bottom of left column "KA, .. LLA are presented in their half form". Peter's doc, L2/04-279 says "most consonants" have a (subjoined) half form (table 6 on page 5.

¶315 However, look at JA. In the TDIL doc, Table 1 on page 54 gives only three oher consonants after which JA subjoins. If that is the case, it is really worth saying that JA has a half form, or is more appropriate to say that there are three conjuncts with JA in second position for which one can identify the fragment that corresponds to JA? Personally, I think the later.

¶316 Furthermore, the three consonants listed by TDIL are KA, LA and LLA. First, we can't tell what they possibly have in mind, because they have KA instead of JA in their examples. Second, I checked utkal.ttf (which is apparently well accepted in the Linux community) and I don't see any subjoined JA.

¶317 Starting from utkal, I see that NNA, TA, NA, BHA, MA, RA, LA, LLA and VA have subjoined half forms. RA and TA are not in TDIL's list only because they are treated separated. That leaves KA and JA in TDIL's list but not in utkal; and NNA, BHA, MA in utkal but not in TDIL. And all of that is short of Peter's 'most''. Given the disparity of sources, the likely that we can reliably render any set of half forms seems pretty small, unless Unicode set a list.

¶318 The left side of the vowel signs 6[°] -ai, 6[°] - \bar{o} and 6[°] -au can be represented in isolation by <0020, 0B47> \rightarrow 6.

¶319 The top side of the vowel sign 6[°] -ai can be represented in isolation by <0020, 0B56> \rightarrow . The coded character U+0B56[°] ORIYA AI LENGTH MARK does not correspond to an abstract character and is encoded solely to provide a canonical decompositon for U+0B48 6[°] ORIYA VOWEL SIGN AI, and to represent its top side in isolation.

¶320 The right side of the vowel sign 6 \bigcirc | -ō can be represented in isolation by <0020, 0B3E> → 1

¶321 The right side of the vowel sign 6 au can be represented in isolation by $<0020, 0B57 > \rightarrow \exists$ The coded character U+0B57 \exists ORIYA AU LENGTH MARK does not correspond to an abstract character and is encoded solely to provide a canonical decomposition for U+0B4C 6 \exists ORIYA VOWEL SIGN AU, and to represent its right side in isolation.

6.6 Tamil

¶322 The Tamil script is a South Indian script. South Indian scripts are structurally related to the North Indian scripts, but they are used to write the Dravidian languages of southern India and of Sri Lanka, which are genetically unrelated to the North Indian languages such as Hindi, Bengla, and Gujarati. The shapes of letters in the South Indian scripts are generally quite distinct from the shapes of letters in Devanagari and its related scripts.

¶323 The Tamil script is used to write the Tamil language of the Tamil Nadu state in India, as well as minority languages such as the Dravidian language Badaga and the Indo-European language Saurashtra. Tamil is also used in Sri Lanka, Singapore, and parts of Malaysia. The Tamil script also lacks conjunct consonant forms.

¶324 Tamil is written horizontally from left to right, and words are separated by space.

¶325 *Standards.* The Tamil block of the Unicode standard is based on ISCII-1988. See section <u>4.1, *The* ISCII *Standard*</u>, for more details on the relation to ISCII.

¶326 *Encoding.* The encoding of the Tamil script in Unicode follows the consonant linking model. The consonant linker coded character of Tamil is U+0BCD $\dot{}$ TAMIL SIGN VIRAMA.

¶327 *Letters.* The Tamil script has fewer consonants than the other Indic scripts, no vocalic vowels, but has all the other vowels presented in section <u>2.1, Letters</u>. In addition, Tamil as the full complement of additional letters presented in section .

¶328 The consonant nasalization sign \circ is not used in Tamil, but it encoded as U+0B82 \circ TAMIL SIGN ANUSVARA.

¶329 The sign \circ is known as the *aytham* in Tamil. It is historically related to the visarga in other Indic scripts, but has become an ordinary spacing letter in Tamil. It is used to modify the sound of other consonants and, in particular, to represent the spelling of words borrowed into Tamil from English or other languages. It is encoded as U+0B83 \circ TAMIL SIGN VISARGA; despite its position in the code charts and its name, this sign is not similar to the unvoiced aspiration sign of the other scripts.

¶330 When representing the "missing" consonants in transcriptions of languages such as Sanskrit or Saurashtra, superscript European digits are often used, so $\square^2 =$ pha, $\square^3 =$ ba, and \square = bha. The characters U+00B2 ² SUPERSCRIPT TWO, U+00B3 ³ SUPERSCRIPT THREE, and U+2074 SUPERSCRIPT FOUR can be used to preserve this distinction in plain text.

¶331 *Digits and Numbers.* Tamil has a full set of digits, encoded at U+0BE6 0 TAMIL DIGIT ZERO ... U+0BEF க TAMIL DIGIT NINE.

¶332 Traditionally, the digits あ 1 through あ 9 are used in a positional system, along with D 10, m 100 and \pm 1000, which are encoded as U+0BF0 D TAMIL NUMBER TEN, U+0BF1 m TAMIL NUMBER ONE HUNDRED and U+0BF2 \pm TAMIL NUMBER ONE THOUSAND.

¶333 *Punctuation and Signs.* Tamil uses the danda || and double danda || encoded as U+0964 | DEVANAGARI DANDA and U+0965 || DEVANAGARI DOUBLE DANDA.

[334 [Amd3] The om sign is encoded as U+0BD0 a TAMIL OM.

¶335 െ... ഐ

9336 rf2 is used to [...] and is encoded as U+0BFA rf2 TAMIL NUMBER SIGN.

¶337 *Typographic Clusters.* The Tamil script is notable for the quasi absence of typographic interactions between the consonant letters which form a typographic cluster. About the only case is that of π ka + \Im sa which results in π ksa but even that form is disappearing in modern usage. As a result, clusters are essentially always rendered by strategy 3.

¶338 On the other hand, there are rich interactions between the consonants of a cluster and the diacritic marks for the vowels:

∟ța∙ी-i	→ ւգ ți
∟ța•°°-ī	\rightarrow le $\dot{t}\bar{1}$
ல la ∙ ി -i	→ കി li
ல la • °° -ī	→ ລໍ lī
க ka • ு -u	→ கு ku
க ka • ூ -ū	→கூ kū
ங 'na • ு -u	→ щ 'nu
ங 'nа • ூ -ū	→ நூ 'nū
ச ca • ு -u	→ சு cu
ச ca • ூ -ū	→ சூ cū
ஜ ja • ு -u	→ஜு ju
ஜ ja • ூ -ū	→ஜூjū
ஞ ña • ு -u	→ ஞ ñu
ஞ ña • ூ -ū	→ ஞூ ñū
∟ța•ு-u	→ (ի țu
∟ța•ூ-ū	→ (β țū
ഞ്ഞ ṇa • ு -u	→ ഞ_ ņu
ഞ ṇa • ூ -ū	→ ഞ_ī ņū
த ta • ு -u	→ து tu
த ta • ூ -ū	→ தூ tū
ந na • ு -u	→ ҧ nu
ந na • ூ -ū	→ ҧӷ nū
ன <u>n</u> a • ு -u	→னு <u>n</u> u
ன <u>n</u> a • ෟ -ū	→ னூ <u>n</u> ū
⊔ pa • ு -u	→цри
ப pa • ூ -ū	→பூ pū
ம ma • ு -u	→ (p mu
ம ma • ூ -ū	→ േµp mū

ш уа• ு -и	→щ уи
ய ya • ூ -ū	→щ yū
ர ra • ு -u	→ (ҧ ru
ர ra • ூ -ū	→ ҧ rū
ற <u>r</u> a • ு -u	→ ற] <u>r</u> u
ற <u>r</u> a • ூ -ū	→ற <u>ா</u> ū
හ la •	→ லு lu
ல la • ூ -ū	→ லூ lū
ள ļa • ு -u	→ளு ļu
ள ļa • ூ -ū	→ ளூ ļū
ழ <u>l</u> a • ு -u	→ ழு <u>l</u> u
ழ <u>l</u> a • ூ -ū	→ழ <u>l</u> ū
ഖ va • ு -u	→ഖุ vu
ഖ va • ூ -ū	→ வூ vū
ர ra • ு -u	→ (ҧ ru
ர ra • ூ -ū	→ ҧ rū
ഐ ṣa • ு -u	→ണം șu
ஷ șa • ு -ū	→ഐ [©] șū
സ sa • ு -u	→ണം su
ஸ sa • ூ -ū	→ൺ sū
ഈ ha • ு -u	→ഈ hu
ഈ ha • ூ -ū	→ഈ9 hū
க்ஷ kṣa • ு -u	→ குஷு kṣu
க்ஷ kṣa • ூ -ū	→ கூஷூ kṣū

¶339 The consonant π ra is simplified when this does not lead to confusion; for example:

$$\begin{split} & \boldsymbol{\eta} \ \mathbf{ra} \boldsymbol{\cdot} \dot{\boldsymbol{\circ}} & \rightarrow \dot{\boldsymbol{\Pi}} \ \mathbf{r} \\ & \boldsymbol{\eta} \ \mathbf{ra} \boldsymbol{\cdot} \boldsymbol{\circ} \boldsymbol{\uparrow} \mathbf{i} \rightarrow \boldsymbol{\mathfrak{fl}} \ \mathbf{ri} \\ & \boldsymbol{\eta} \ \mathbf{ra} \boldsymbol{\cdot} \boldsymbol{\circ}^{\mathsf{e}} \mathbf{\cdot} \mathbf{\bar{i}} \rightarrow \boldsymbol{\mathfrak{fl}} \ \mathbf{r\bar{i}} \end{split}$$

¶340 However, various governmental bodies mandate that the basic shape of π should be used for these ligatures as well, especially in school textbooks. Media and literary publications in Malaysia and Singapore mostly use the unchanged form of π .

¶341 Some changes are no longer in common use:

ன n_a •െ -ai →னை n_ai ல la •െ -ai → லை lai ள la •െ -ai → ണെ lai

 $\P342$ The coded character U+0BD7 on TAMIL AU LENGTH MARK does not correspond to an abstract character, i.e., it is not used for the representation of text. It is encoded to represent the right side of the vowel sign and to provide a canonical decomposition for U+0BCC and TAMIL VOWEL SIGN AU.

¶343 *Named sequences.* Tamil is less complex than some of the other Indic scripts, and both conceptually and in processing can be treated as an atomic set of elements: consonants, stand-alone vowels, and syllables. The table below shows these atomic elements, with the corresponding Unicode characters or sequences. These elements have been approved as Tamil named character sequences: see NamedSequences.txt in the Unicode Character Database.

¶344 In implementations such as natural language processing, where it may be useful to treat such Tamil text elements as single code points for ease of processing, they can be mapped to a segment of the Private Use Area.

¶345 In the table below, the first row shows the transliterated representation of the Tamil vowels in abbreviated form, while the first column shows the transliterated representation of the Tamil consonants. These short strings can then be concatenated to form unique names for Tamil pure consonants and syllables: K, KA, KAA, KI, KE, KU, and so on. Details on the complete names for each element can be found in NamedSequences.txt.

		А	AA	Ι	Π	U	UU	E	EE	AI	0	00	AU
	°° 0B83	의 0B85	ஆ 0B86	(A) 0B87	FT 0B88	உ 0B89	<u>୭ଗ</u> ୀ 0B8A	റ 0B8E	ஏ 0B8F	₿ 0B90	ള 0B92	ତ୍ତୁ 0B93	ତୃଣ <u>ୀ</u> 0B94
K	க் 0B95 0BCD	க 0B95	あIT 0B95 0BBE	கி 0B95 0BBF	සී 0B95 0BC0	(5 0B95 0BC1	ቻጫ 0B95 0BC2	கெ 0B95 0BC6	கே 0B95 0BC7	னைக் 0B95 0BC8	கொ 0B95 0BCA	Съп 0B95 0BCB	கௌ 0B95 0BCC
NG	њі 0В99 0BCD	БЫ 0B99	ГЫП 0B99 0BBE	пЫ 0B99 0BBF	ഥ്പ് 0B99 0BC0	떠 0B99 0BC1	நு 0B99 0BC2	Gпы 0В99 0ВС6	Съ 0в99 0вс7	ைங 0B99 0BC8	Сып ов99 0ВСА	Спып ов99 овсв	நௌ 0B99 0BCC
С	ச் 0B9A 0BCD	ሆ 0B9A	₽П 0В9А 0BBE	சி oB9A OBBF	& 0В9А 0ВС0	풍 0B9A 0BC1	碼 0B9A 0BC2	ිළ 0B9A 0BC6	රිජ 0B9A 0BC7	ണട് 0B9A 0BC8	Овса овуа	Сғп ов9а овсв	சௌ 0B9A 0BCC
NY	ஞ் obje obcd	ஞ 0B9E	ஞா 0B9E 0BBE	ஞி obje obje	ஞீ 0B9E 0BC0	ஞு 0B9E 0BC1	ஞூ 0B9E 0BC2	බල 0B9E 0BC6	රීල් 0B9E 0BC7	ஞை 0B9E 0BC8	ௌா 0B9E 0BCA	ஞோ 0B9E 0BCB	ஞௌ obje obcc
TT	亡 0B9F 0BCD	∟ 0B9F	L_∏ 0B9F 0BBE	لم 0B9F 0BBF	لد 0B9F 0BC0	(ј) 0B9F 0BC1	டு 0B9F 0BC2	G∟0B9F 0BC6	C∟ 0B9F 0BC7	ണ്ഥ 0B9F 0BC8	ССА ОВОЯ	С∟п ов9ғ овсв	டௌ oB9F0BCC
NN	ळा ०ваз ०вср	OBA3	OBA3 OBBE	ഞ്ഞി oba3 obbF	ഞ് 0BA3 0BC0	∭ 0BA3 0BC1	∭ 0BA3 0BC2	ணெ 0BA3 0BC6	Coor OBA3 OBC7	ഞ്ഞ് 0BA3 0BC8	ணொ 0BA3 0BCA	Comn OBA3 OBCB	ത്തെണ oba3 obcc
Т	த் oba4 obcd	த 0BA4	தா oba4 obbe	தி 0BA4 0BBF	தீ 0BA4 0BC0	து 0BA4 0BC1	தூ 0BA4 0BC2	தெ 0BA4 0BC6	தே 0BA4 0BC7	தை 0BA4 0BC8	தொ ова4 овса	தோ ова4 овсв	தௌ OBA4 OBCC
N	ҧ́ овая овср	Б 0ВА8	БП ОВА8 ОВВЕ	ҧ҄ӏ овая оввғ	ҧ овая овсо	ђј 0ВА8 0ВС1	Ҧ∏ 0ВА8 0ВС2	Србовая ОВС6	Сђ оваз 0вс7	நை OBA8 0BC8	СБП ОВА8 ОВСА	СБП ОВА8 ОВСВ	நௌ 0BA8 0BCC
Р	ப் оваа obcd	∐ 0BAA	ШП OBAA OBBE	OBAA 0BBF	ഥ് obaa obco	Ц 0ВАА 0BC1	Ц 0ВАА 0BC2	Gы оваа овсе	Сы оваа овст	பை 0BAA 0BC8	С⊔П 0ВАА 0ВСА	Сып оваа овсв	பௌ obaa obcc
М	D OBAE OBCD	LD 0BAE	LDT 0BAE 0BBE	மி obae obbf	மீ obae 0BC0	UD OBAE OBCI	UD OBAE OBC2	GLD OBAE 0BC6	CLD OBAE 0BC7	ണ്ഥ obae 0BC8	GLDП 0BAE 0BCA	CLDIT 0BAE 0BCB	மௌ obae obcc
Y	ய் obaf obcd	ULI 0BAF	ULIT OBAF OBBE	UG OBAF OBBF	ഥ് obaf 0BC0	Щ OBAF OBC1	UL OBAF OBC2	Gш овағ овс6	Сш овағ овст	பைய OBAF 0BC8	Gшп 0BAF0BCA	Сшп овағ овсв	யௌ obafobcc
R	П овво овср	Г 0ВВ0	ГГП 0ВВ0 0ВВЕ	повво оввғ	п овво овсо	(T5 0BB0 0BC1	(Ҧ 0ВВ0 0ВС2	ОГЛ ОВВО ОВС6	Сџовво 0вс7	のり 0BB0 0BC8	ОГЛП ОВВО ОВСА	Слп овво овсв	ரௌ 0BB0 0BCC
L	ல் oBB2 oBCD	ର 0BB2	லா 0BB2 0BBE	ති 0BB2 0BBF	ര് 0BB2 0BC0	ല്ല 0BB2 0BC1	ല്ല∏ 0BB2 0BC2	ରେ 0BB2 0BC6	Cබ 0BB2 0BC7	തെல 0BB2 0BC8	லொ 0BB2 0BCA	Сог овв2 0всв	ରେଶ୍ 0BB2 0BCC
v	ഖ് obbs obcd	ഖ 0BB5	വ∏ 0BB5 0BBE	ബി obbs obbF	ഖ് 0BB5 0BC0	പ്പ 0BB5 0BC1	പ്ര 0BB5 0BC2	പെ 0BB5 0BC6	Cඛ 0BB5 0BC7	ഞഖ 0BB5 0BC8	வொ 0BB5 0BCA	வோ 0BB5 0BCB	ബെണ 0BB5 0BCC
LLL	ழ் 0BB4 0BCD	ழ 0BB4	ழп овв₄ овве	ပြာ 0BB4 0BBF	ழீ 0BB4 0BC0	ழு 0BB4 0BC1	ഥ്ര 0BB4 0BC2	Gழ 0BB4 0BC6	Сழ 0BB4 0BC7	ழை 0BB4 0BC8	Срп овва овса	Сழп овв4 овсв	ழௌ 0BB4 0BCC
LL	ണ് 0BB3 0BCD	GT 0BB3	OTTIT OBB3 OBBE	ണി oBB3 oBBF	ണ് 0BB3 0BC0	ளு 0BB3 0BC1	ണ്ട്ര 0BB3 0BC2	Сат оввз овс6	Сат оввз овс7	ഞണ 0BB3 0BC8	Gatt 0BB3 0BCA	Сатп оввз овсв	ണെണ 0BB3 0BCC

RR	ற் லக	ற obbi	றா ம	றி 0BB1	றீ லக	Щ овві	∭ 0BB1	றெ 0BB1	றே 0BB1	றை லக	றொலை	றோ லக	றௌ 0BB1
	0BCD		0BBE	0BBF	0BC0	0BC1	0BC2	0BC6	0BC7	0BC8	0BCA	0BCB	0BCC
NNN	ன் 0BA9	ങ	னா 0BA9	ണി 0BA9	ങ് 0BA9	ത്വ 0BA9	ഞ ∏ 0BA9	തെ obag	කෙ OBA9	തെ	னொ	ளோ	னௌ
	0BCD	0BA9	0BBE	0BBF	0BC0	0BC1	0BC2	0BC6	0BC7	0BA9 0BC8	0BA9 0BCA	0BA90BCB	0BA9 0BCC
J	ல் 0B9C 0BCD	છ 0B9C	OBBE	නි 0B9C 0BBF	ജ് 0B9C 0BC0	の 0BC1	භූව _{0B9C} 0BC2	බනූ 0B9C 0BC6	රිනූ 0B9C 0BC7	ജെ 0B9C 0BC8	ஜொ 0B9C 0BCA	ஜோ 0B9C 0BCB	ஜௌ 0B9C0BCC
SH	UU 0BB6	ហា	ИОП 0ВВ6	UTI 0BB6	0BB6	0BB6 °U	UUD 0BB6		GUTI 0BB6	സെ 0886	மொ	ஶோ	மொ
511	0BCD	0BB6	0BBE	0BBF	0BC0	0BC1	0BC2	0BC6	0BC7	0BC8	0BB6 0BCA	0BB6 0BCB	0BB6 0BCC
SS	പ് 0BB7	ବ୍ୟ	ஷா 0BB7	ല്പെ 0BB7	പ്പെ 0BB7	ച്ചെ 0BB7	ല്പെ∋ 0BB7	ഷെ 0887	ക്ഷേ 0BB7	ഞ്ഞ	ஷொ	ஷோ	ஷௌ
	0BCD	0BB7	0BBE	0BBF	0BC0	0BC1	0BC2	0BC6	0BC7	0BB7 0BC8	0BB7 0BCA	0BB7 0BCB	0BB7 0BCC
S	ஸ் 0BB8	സ	லா obbs	സി 0888	സ് രങ്ങ	സ്ന 0BB8	സ് ⁹ 0BB8	സെ 0BB8	Cen 0BB8	സെ	ஸொ	ஸோ	ஸௌ
	0BCD	0BB8	0BBE	0BBF	0BC0	0BC1	0BC2	0BC6	0BC7	0BB8 0BC8	0BB8 0BCA	0BB8 0BCB	0BB8 0BCC
Н	ണ് 0BB9	ണ വള്ളം	ണ∏ 0BB9	ണ്ണി 0BB9	ണ്ട് 0BB9	ണ്ടം 0BB9	ണ്ട്ര വള്ള വളന്ന	ബ്രെ വങ്ങം വങ്ങൾ	ബ്രോ 0BB9.0BC7	ഞ്ഞ നള്ള നൂറ്റം	ஹொ obbgobca	ஹோ	ஹௌ முஜல முகுகு
	UBCD		UDDE	UDDF	UBCU	ODCI	000/0002	ODD/ ODCO	ODD/ ODC/	ODD) ODCO	ODD/ ODCA	ODD/ ODCD	OBD) OBCC
	க்ஷ் 0B95	க்ஷ	கூடிா	கூடி 0895	கூடீ 0B95	ምሞ	கூடி	க்ஷெ	க்ஷே	ക്ഷെ	க்ஷொ	க்ஷோ	க்ஷௌ
KSS	0BCD 0BB7 0BCD	0B95 0BCD 0BB7	0B95 0BCD 0BB7 0BBE	0BCD 0BB7 0BBF	0BCD 0BB7 0BC0	0B95 0BCD 0BB7 0BC1	0B95 0BCD 0BB7 0BC2	0B95 0BCD 0BB7 0BC6	0B95 0BCD 0BB7 0BC7	0B95 0BCD 0BB7 0BC8	0B95 0BCD 0BB7 0BCA	0B95 0BCD 0BB7 0BCB	0B95 0BCD 0BB7 0BCC

6.7 Telugu

¶346 The Telugu script is a South Indian script used to write the Telugu language of the Andhra Pradesh state in India, as well as minority languages such as Gondi (Adilabad and Koi dialects) and Lambadi. The script is also used in Maharashtra, Orissa, Madhya Pradesh, and West Bengal. The Telugu script became distinct by the thirteenth century CE and shares ancestors with the Kannada script.

¶347 Many Telugu letters have a v-shaped headstroke, which is a structural mark corresponding to the horizontal bar in Devanagari and the arch in Oriya script.

¶348 *Standards.* The Telugu block of the Unicode standard is based on ISCII-1988. See section <u>4.1, *The* ISCII *Standard*</u>, for more details on the relation to ISCII.

¶349 *Encoding.* The encoding of the Telugu script in Unicode follows the consonant linking model. The consonant linker coded character of Telugu is U+0C4D 5 TELUGU SIGN VIRAMA.

¶350 *Letters.* Telugu has the full complement of basic letters presented in section <u>2.1, Letters</u>. In addition, Telugu as the full complement of additional letters presented in section, except for ϖ na and \wp la.

¶351 The vowel nasalisation sign ్c is known as అరసున్న arasunna and is encoded as U+0C01 ్c TELUGU SIGN CANDRABINDU.

¶352 The consonant nasalization sign ం is known as సున్న sunna and is encoded as U+0C02 ం TELUGU SIGN ANUSVARA.

¶353 The unvoiced aspiration sign ုး is known as ລະວັດ visarga and is encoded as U+0C03 ူး TELUGU SIGN VISARGA.

¶354 Telugu does not have a nukta.

¶355 The letters 式 and 惑 have been used historically for the voiceless alveolar affricate [ts] and the voiced alveolar affricate [dz] respectively. They are encoded as U+0C58 式 TELUGU LETTER TSA and U+0C59 惑 TELUGU LETTER DZA. These characters are commonly found in old grammar books and dictionaries. Historically, various glyphs were used to render these two characters but the notation used currently was proposed by Charles Philip Brown in the 19th century. These letters are not in current use in contemporary Telugu.

1356 *Digits and Numbers.* Telugu has a full set of decimal digits, encoded at U+0C66 O TELUGU DIGIT ZERO ... U+0C6F E TELUGU DIGIT NINE.

¶357 A set of digits for fractions is encoded at U+0C78 ♀ TELUGU FRACTION DIGIT ZERO FOR ODD POWERS OF FOUR .. U+0C7E ≥ TELUGU FRACTION DIGIT THREE FOR EVEN POWERS OF FOUR.

¶358 *Punctuation and Signs.* Telugu uses the danda | and double danda || encoded as U+0964 | DEVANAGARI DANDA and U+0965 || DEVANAGARI DOUBLE DANDA.

¶359 The *avagraha* is encoded as U+0C3D ≥ TELUGU SIGN AVAGRAHA.

¶360 The sign 𝙂 is used to denote a unit of volume or weight and is encoded as U+0C7F 𝙂 TELUGU SIGN TUUMU.

¶361 *Atomicity of Vowel Letters.* Vowel letters are encoded atomically in Unicode, even if they can be analyzed visually as consisting of multiple parts. The table below shows the letters that can be analyzed, the single code point that should be used to represent them in text, and the sequence of code points resulting from analysis that should not be used.

To represent Use Do not use & 0C13 <0C12, 0C55> 配 0C14 <0C12, 0C4C>

¶362 Similarly, some vowel signs are atomically encoded:

To represent	Use	Do not use
ూ	0C42	<0C41, 0C3E>
్రా	0C44	<0C43, 0C3E>
ें	0C47	<0C46, 0C55>
್	0C4B	<0C4A, 0C55>

¶363 *Typographic Clusters.* Telugu consonant clusters are most commonly represented by a subscripted, and often transformed, consonant glyph for the second element of the cluster:

K ga • K ga	→ ň gga
క ka • య ya	→ ≤ړ kya
క ka•ష ṣa	→ ǯ kṣa
≤ ka • ≤ ka	→Šb
మ ma • మ ma	→ మ్మ
v la • v la	$\rightarrow \overset{\mathfrak{O}}{\mathfrak{m}}$
ĕ ta • ĕ ta	\rightarrow \tilde{s}
య ya • య ya	→య్య
వ va • వ va	\rightarrow \sim
ŏ ra • ŏ ra	→ğ
న na • న na	→న్న
స sa • త ta • ర ra	کي stra

¶364 Satellite Consonants. Historically, a & ra in initial position could be rendered is a satellite consonant. It was then displayed on the right side of cluster core by the sign &, known as 3030Ros valapalagilaka.

¶365 *Satellite vowels.* Vowel signs attach either above, above and below or the right; when they attach above, and the consonant has a v-shaped headstroke, that headstroke disappears. Here are example on \leq ka:

		š ka
ਭ ā	ে∵ -ā	হ <u>ু</u> kā
ສ i	;;°-i	8 ki
ఈ ī	^{\$} -ī	ŝ kī
ఉ u	ာ -u	కు ku
ఊ ū	ూ -ū	కూ kū
ຍນນ rຸ	ു -r	ຮ _{ິງ} kŗ
ౠ rై	ာ -ļ	۲ۍ kļ
ఎ e	ි -е	3 ke
် ē	<i>ੈ -</i> ē	ਤੇ kē
ສ ai	ू -ai	<u>₹</u> kai
۵ ۵	ో -0	s" ko
٤ō	ో -ō	s ^e kō
ਛਾ au	ౌ -au	ਤਾ kau

¶366 A few combinations show a more complex interaction of the vowel signs:

	ా -ā	i 🔅 -	i ് -	ī (ຸນ -ເ	1 ् र ग -i	ī℃ -0	्र - ि	్ౌ -	au
ఙ na				జు	ఙూ				
ప pa				పు	పూ				
ఫ pha				ఫు	పూ				
మ ma	l	మి	మీ			మొ	మో		
య ya	l	ဿ	ဿာ			ಯು	ಯ್		
వ va				వు	పా				
హ ha	హా			హు	హూ	హొ	హో	హౌ	

¶367 *Clusters without vowels.* For Sanskrit or foreign words ending in a consonant sound, the absence of a final vowel sound is indicated by the presence of the sign 5, called 3×20 pollu, above the consonant.

¶368 Non-joined Forms of Clusters.

¶369 *Typographic fragments.* The coded character U+0C55 ° TELUGU LENGTH MARK does not correspond to an abstract character and is encoded solely to represent the second element of U+0C47 7 TELUGU VOWEL SIGN EE and of U+0C4B 7⁶

TELUGU VOWEL SIGN OO in isolation.

 $\P370$ The coded character U+0C56 \bigcirc TELUGU AI LENGTH MARK does not correspond to an abstract character is encoded solely to provide a canonical decomposition for U+0C48 \bigcirc TELUGU VOWEL SIGN AI, and to represent its bottom side in isolation.

6.8 Kannada

¶371 The Kannada script is a South Indian script. It is used to write the Kannada (or Kanarese) language of the Karnataka state in India and to write minority languages such as Tulu. The Kannada language is also used in many parts of Tamil Nadu, Kerala, Andhra Pradesh, and Maharashtra. This script is very closely related to the Telugu script both in the shapes of the letters and in the behavior of conjunct consonants.

¶372 Many Kannada letters have a headstroke formed of a horizontal line and a hook corresponding to the horizontal bar in Devanagari and the arch in the Oriya script.

¶373 *Standards.* The Kannada block of the Unicode standard is based on ISCII-1988. See section <u>4.1, *The* ISCII *Standard*</u>, for more details on the relation to ISCII.

¶374 *Encoding.* The encoding of the Kannada script in Unicode follows the consonant linking model. The consonant linker coded character of Kannada is U+0CCD ^{••} KANNADA SIGN VIRAMA.

¶375 *Letters.* Kannada has the full complement of basic letters presented in section <u>2.1, Letters</u>. In addition, Kannada as the full complement of additional letters presented in section except for ϖ na and \mathfrak{g} ta.

¶376 The consonant nasalization sign \circ and is encoded as U+0C82 \circ KANNADA SIGN ANUSVARA.

¶377 The unvoiced aspiration sign ℃ is encoded as U+0C83 ℃ KANNADA SIGN VISARGA.

¶378 The tongue-root sibilant \Box is encoded as U+0CF1 \Box KANNADA SIGN JIHVAMULIYA.

¶379 The labial sibilant \Box is encoded as U+0CF2 \Box KANNADA SIGN UPADHMANIYA

¶380 & la is actually an obsolete Kannada letter that is transliterated in Dravidian scholarship as z, l or r. It is encoded as U+0CDE es

KANNADA LETTER FA. This coded character should have been named "LLLA" rather than "FA", so the name in this standard is simply a mistake. This letter has not been actively used in Kannada since the end of the tenth century. Collations should treat U+0CDE e9

KANNADA LETTER FA as following U+0CB3 & KANNADA LETTER LLA.

¶381 U+0CD5 © KANNADA LENGTH MARK and U+0CD6 ; KANNADA AI LENGTH MARK are provided for the encoding of fragments of vowel signs. However, those two length marks have no independent existence in the Kannada writing system and do not play any part as independent codes in the traditional collation order.

 ¶382 Digits and Numbers.
 Kannada has a full set of decimal digits, encoded at U+0CE6 O KANNADA DIGIT ZERO .. U+0CEF

 © KANNADA DIGIT NINE.

¶383 *Punctuation and Signs.* Kannada uses the danda | and double danda || encoded as U+0964 | DEVANAGARI DANDA and U+0965 || DEVANAGARI DOUBLE DANDA.

¶384 The avagraha is encoded as U+0CBD 5 KANNADA SIGN AVAGRAHA.

¶385 Atomicity of Vowel Letters. Vowel letters are encoded atomically in Unicode, even if they can be analyzed visually as consisting of

multiple parts. The table below shows the letters that can be analyzed, the single code point that should be used to represent them in text, and the sequence of code points resulting from analysis that should not be used.

 To represent Use
 Do not use

 ಔ
 0C94 <0C92, 0CCC>

¶386 Similarly, some vowel signs are atomically encoded:

To represent	Use	Do not use
ూ	0C42	<0C41, 0C3E>
ూ	0C44	<0C43, 0C3E>
ें	0C47	<0C46, 0C55>

¶387 *Typographic Clusters.* Kannada consonant clusters are most commonly represented by a subscripted, and often transformed, consonant glyph for the second element of the cluster:

ಗ ga • ಗ ga → ಗ್ಗ gga ಕ ka • ಕ ka → ಕ್ಕ kka ಕ ka • ಯ ya → ಕ್ಕ kya ಕ ka • ಷ şa → ಕ್ಷ kşa

¶388 Satellite Consonants. In Kannada, of ra is a satellite consonant in initial position. It is then displayed on the right side of the cluster core: ಕ್ಷF rkna, ಕ್ವಾF rknā, ಕೀfF rknī. Alternatively, an initial of ra can keep it full form, and this rendering can be imposed by the use of U+200D ZERO WIDTH JOINER: of rka

¶389 *Satellite vowels.* Vowel signs attach either above, above and below or the right; when they attach above, and the consonant has a v-shaped headstroke, that headstroke disappears. Here are example on $\vec{\sigma}$ ka:

ಕ ka ಆ ā ಾ -ā ಕಾ kā ಇ i ಿ -i ಕಿ ki ಈ ī ೀ -ī ಕೀ kī ಉ u ು -u ಕು ku ಊ ū ೂ -ū ಕು ku ಯ ಫ ್ -r ಕೃ kr ಖು r ೃ -r ಕೃ kr ಖಾ ಸ್ ೃ -! ಕೃ kļ ಲು! ್ಲ ್ಲ ಕೃ kậ ಲಾ ಸ್ ಜ್ kậ ಎ ೇ -e ಕೇ kē ஐai ீ-ai தீkai 遼au ゔ-au ಕಾkau

¶390 A few combinations show a more complex interaction of the vowel signs:

ೀ -	∙īാ -	u ை -ū
ಮ ma ಮೀ		
ಪ pa	ಪು	ಪೂ
ಫ pha	ಫ್ರ	ಫೂ
ವ va	ವು	ವೂ

¶391 Non-joined forms of clusters..

1392 Clusters without vowels. For Sanskrit or foreign words ending in a consoant sound, the absence of a final vowel sound is indicated by the presence of the sign of above the consonant.

1393 Typographic fragments. The coded character U+0CD5 °C KANNADA LENGTH MARK does not correspond to an abstract character and is encoded solely to represent the second element of U+0CC7 ° e KANNADA VOWEL SIGN EE in isolation.

1394 The coded character U+0CD6 , KANNADA AI LENGTH MARK does not correspond to an abstract character and is encoded solely to represent the second element of U+0CC8 , KANNADA VOWEL SIGN AI in isolation.

6.9 Malayalam

¶395 The Malayalam script is a South Indian script used to write the Malayalam language of the Kerala state. Malayalam is a Dravidian language like Kannada, Tamil, and Telugu. Throughout its history, it has absorbed words from Tamil, Sanskrit, Arabic, and English.

¶396 Standards. The Malayalam block of the Unicode standard is based on ISCII-1988. See section <u>4.1, The ISCII Standard</u>, for more details on the relation to ISCII.

1397 *Encoding*. The encoding of the Malayalam script follows the consonant linking model. The consonant linker coded character of Malayalam is U+0D4D ("MALAYALAM SIGN VIRAMA.

1398 Letters. Malayalam has the full complement of basic letters presented in section 2.1, Letters, except for the vowel sign for 60 60 . In addition, Malayalam as the full complement of additional letters presented in section, except for on na.

¶399 The consonant nasalization sign ○o and is encoded as U+0D02 ○o MALAYALAM SIGN ANUSVARA.

400 The unvoiced aspiration sign \Im : is encoded as U+0D03 \Im : MALAYALAM SIGN VISARGA.

1401 Chillu Characters. The letters ൺ ൺ, ൻ ൻ, ർ ർ, ൽ ൽ, ൾ ൾ, and ൿ ൿ are know as chillu or cillaksaram characters.

1402 In Malayalam-language text, chillu letters never start a word. The chillu letters month ind, no no, d d, od month, and d d are quite

common; ൿ ൿ is not very common.

¶403 Prior to Unicode 5.1, the representation of text with chillus was problematic, and not clearly described in the text of the standard. Because older data will use different representation for chillus, implementations must be prepared to handle both kinds of data. The following table shows the relation between the representation in Unicode Version 5.0 and earlier and the new representation in Version 5.1, for the chillu letters considered in isolation.

Character	Representation in 5.0 and prior	Preferred 5.1 and later representation
ൺ ൺ	NNA, VIRAMA, ZWJ (0D23, 0D4D, 200D)	U+0D7A ൺ MALAYALAM LETTER CHILLU NN
ൻ ൻ	NA, VIRAMA, ZWJ (0D28, 0D4D, 200D)	U+0D7B ro MALAYALAM LETTER CHILLU N
ል ል	RA, VIRAMA, ZWJ (0D30, 0D4D, 200D)	U+0D7C & MALAYALAM LETTER CHILLU RR
ൽ ൽ	LA, VIRAMA, ZWJ (0D32, 0D4D, 200D)	U+0D7D ൽ MALAYALAM LETTER CHILLU L
რ რ	LLA, VIRAMA, ZWJ (0D33, 0D4D, 200D)	U+0D7E ൾ MALAYALAM LETTER CHILLU LL
ക് ക്	undefined	U+0D7F ൿ MALAYALAM LETTER CHILLU K

¶404 The letter \cap \cap is normally read r. Repetition of that sound is written by two occurrences of the letter: $\cap \cap$. Each occurrence can bear a vowel sign.

1405 Repetition of the letter, written either 00 or β , is also used for the sound tt. In this case, the two 0 fundamentally behave as a

digraph. The digraph can bear a vowel sign in which case the digraph as a whole acts graphically as an atom a left vowel part goes to the left of the digraph and a right vowel part goes to the right of the digraph. Historically, the side-by-side form was used until around 1960 when the stacked form began appearing and supplanted the side-by-side form. The same situation is common in many other orthographies. For example, *th* in English can be a digraph (*cathode*) or two separate letters (*cathouse*); *gn* in French can be a digraph (*oignon*) or two separate letters (*gnome*).

¶406 The sequence <0D31, 0D31> represents $\cap \cap$, regardless of the reading of that text. The sequence <0D31, 0D4D, 0D31> represents Ω . In both cases, vowels signs can be used as:

പാററ	0D2A 0D3E 0D31 0D31	paatta	cockroach	
പാറ്റ	0D2A 0D3E 0D31 0D4D 0D31			
മാറെറാലി	0D2E 0D3E 0D31 0D46 0D31 0D3E 0D32 0D3F	maattoli	echo	
മാറ്റൊലി	0D2E 0D3E 0D31 0D4D 0D31 0D46 0D3E 0D32 0D3F			
ബാറററി	0D2C 0D3E 0D31 0D31 0D31 0D3F	baattari/	battery	
ബാറ്ററി	0D2C 0D3E 0D31 0D4D 0D31 0D31 0D3F			
സൂറററ്	0D38 0D42 0D31 0D31 0D31 0D4D	suu <u>r</u> att	(name of a place)	
സൂററ്റ്	0D38 0D42 0D31 0D31 0D4D 0D31 0D4D			
ടെംപററി	0D1F 0D46 0D02 0D2A 0D31 0D31 0D3F	tempa <u>r</u> ari	temporary (English loan word)	
ലെക്ചററോട് 0D32 0D46 0D15 0D4D 0D1A 0D31 0D31 0D4B 0D1F 0D4D /lekca <u>r</u> aroot/ to the lecturer				

¶407 A very similar situation exists for the combination of rd rd and rd rd. When used side by side, rd rd can be read either /nr/ or /nt/, while rd rd is always read /nt/.

¶408 The sequence <0D7B, 0D31> represents (700, regardless of the reading of that text. The sequence <math><0D7B, 0D4D, 0D31>

represents non. In both cases, vowels signs can be used as appropriate:

ആൻറാ 0D06 0D7B 0D47 0D31 0D3E aantoo (proper name) ആൻറോ 0D06 0D7B 0D4D 0D31 0D47 0D3E എൻറോൺ0D0E 0D7B 0D31 0D47 0D3E 0D7A /enrool/ enroll (English loan word)

¶409 *Digits and Numbers*. Malayalam has a full set of decimal digits, encoded at U+0D66 o MALAYALAM DIGIT ZERO .. U+0D6F ෆ MALAYALAM DIGIT NINE.

¶410 Additional numbers and fraction characters are encoded at U+0D70 ഡ MALAYALAM NUMBER TEN .. U+0D75 ൺ MALAYALAM FRACTION THREE QUARTERS.

[411 *Punctuation and Signs.* Malayalam uses the danda | and double danda || encoded as U+0964 | DEVANAGARI DANDA and U+0965 || DEVANAGARI DOUBLE DANDA.

1412 The avagraha is known as the praslesham in Malayalam and is encoded as U+0D3D [] MALAYALAM SIGN AVAGRAHA.

¶413 The date mark re- is encoded at U+0D79 re- MALAYALAM DATE MARK.

¶414 *Atomicity of Vowel Letters.* Vowel letters are encoded atomically in Unicode, even if they can be analyzed visually as consisting of multiple parts. The table below shows the letters that can be analyzed, the single code point that should be used to represent them in text, and the sequence of code points resulting from analysis that should not be used.

To represent Use Do not use ஹ 0D08 <0D07, 0D57> බ 0D0A <0D09, 0D57> බ 0D10 <0D0E, 0D46> බ 0D13 <0D12, 0D3E> බ 0D14 <0D12, 0D57>

[415 *Typographic Cluster Cores.*

[416 Satellite consonants.

¶417 *Satellite vowels.* In modern times, the dominant practice is to write the dependent form of the 60 60 vowel using only \mathfrak{V} which is placed on the right side of the consonant it modifies; such texts are represented in Unicode using U+0D57 $\Im \mathfrak{V}$ MALAYALAM AU LENGTH MARK. In the past, this dependent form was written using both on the left side and \mathfrak{V} on the right side; U+0D4C \mathfrak{S} \mathfrak{V} MALAYALAM VOWEL SIGN AU can be used for documents following this earlier tradition. This historical simplification started much earlier than the orthographic reforms mentioned above.

[418 Conjoined and Non-Joined Forms of Clusters.

[419 *Typographic Fragments.* todo

6.10 Sinhala

¶420 The Sinhala script, also known as Sinhalese, is used to write the Sinhala language, the majority language of Sri Lanka. It is also used

to write the Pali ans Sanskrit languages. The script is a descendant of Brahmi and resembles the scripts of South India in form and structure.

¶421 *Encoding.* The encoding of the Sinhala script follows the consonant linking model. The consonant linker coded character of Sinhala is U+0DCA \Box SINHALA SIGN AL-LAKUNA.

¶422 Letters. Sinhala has the full complement of basic letters presented in section <u>2.1, Letters</u>.

[423 The consonant nasalization sign \Box is known as the *ansuvaraya* in Sinhala, and is encoded as U+0D82 \Box SINHALA SIGN ANUSVARAYA.

[424 The unvoiced aspiration sign \Box is known as the *visargaya* in Sinhala, and is encoded as U+0D83 \Box SINHALA SIGN VISARGAYA.

¶425 Sinhala differs from other languages of the region in that it has a series of prenasalized stops that are distinguished from the combination of a nasal followed by a stop. In other words, both forms occur and are written differently – for example, AB <U+0D85, U+0DAC> a8}a [a:;a] "sound" versus ACDE <U+0D85, U+0DAB, U+0DCA, U+0DA9> aV}a [a9;a] "egg." In addition, Sinhala has separate distinct signs for both a short and a long low front vowel sounding similar to the initial vowel of the English word "apple," usually represented in IPA as U+00E6 æ latin small letter ae (ash). The independent forms of these vowels are encoded at U+0D87 and U+0D88; the corresponding dependent forms are U+0DD0 and U+0DD1.

¶426 Because of these extra letters, the encoding for Sinhala does not precisely follow the pattern established for the other Indic scripts (for example, Devanagari). It does use the same general structure, making use of phonetic order, matra reordering, and use of the virama (U+0DCA sinhala sign al-lakuna) to indicate conjunct consonant clusters. Sinhala does not use half-forms in the Devanagari manner, but does use many ligatures.

¶427 Digits and Numbers. While the Sinhala script has a full set of decimal digits, those are not encoded at the present time..

428 Punctuation and Signs. — was formerly used as a full stop or period and is encoded as U+0DF4 — SINHALA PUNCTUATION KUNDDALIYA.

1429 *Atomicity of Vowel Letters.* Vowel letters are encoded atomically in Unicode, even if they can be analyzed visually as consisting of multiple parts. The table below shows the letters that can be analyzed, the single code point that should be used to represent them in text, and the sequence of code points resulting from analysis that should not be used.

To represent	Use	Do not use
ආ	0D86	<0D85, 0DCF>
ඇ	0D87	<0D85, 0DD0>
¢۲	0D88	<0D85, 0DD1>
Ĉ ෟ	0D8C	<0D8B, 0DDF>
දිවුම	0D8E	<0D8D, 0DD8>
දිඩා	0D90	<0D8F, 0DDF>
లి	0D92	<0D91, 0DCA>
ෙළු	0D93	<0D91, 0DD9>
ඖ	0D96	<0D91, 0DDF>

¶430 Typographic Cluster Cores. A few consonant clusters are

[431 Satellite Consonants. $\circ \circ$ in initial position is a satellite consonant. It is then displayed as a small hook, called *repaya*, on the top right side of the cluster.

¶432 todo examples

- ¶433 & & in final position is a satellite consonant. It is displayed as a curve below the cluster core, called a *rakaaraansaya*.
- ¶434 todo examples
- ¶436 todo examples

Appendix A. Acknowledgments

¶437 Thanks to the Editorial Committee, Asmus Freytag.

Appendix B. Things to do

¶438 This section is of course not part of the text; it's just a place to remember the things to do to this text.

- ¶439 in the description of the models, introduce VIRAMA (small caps) as the generic name for the consonant linker coded character?
- **[**440 set baseline-to-baseline to some fixed value

Appendix C. References

References

[Snell]	Beginner's Hindi Script London: Hodder & Stoughton, 2003. ISBN 0-07-141984-5		
[Arden]	A Progressive Grammar of the Telugu Language Madras: The Society For Promoting Christian Knowledge, 1905		
[Brown]	A Grammar of the Telugu Language Madras: The Christian Knowledge Society's Press, 1857		
[Feedback]	http://www.unicode.org/reporting.html For reporting errors and requesting information online.		
[Reports]	Unicode Technical Reports <u>http://www.unicode.org/reports/</u> <i>For information on the status and development process for technical reports, and for a list of technical reports.</i>		
[Unicode]	<i>The Unicode Standard, Version 4.0</i> (Boston, MA, Addison-Wesley, 2003. ISBN 0-321-18578-1), as amended by <i>Unicode 4.0.1</i> (http://www.unicode.org/versions/Unicode4.0.1) and by <i>Unicode 4.1.0</i> (http://www.unicode.org/versions/Unicode4.1.0).		
[Versions]	Versions of the Unicode Standard <u>http://www.unicode.org/versions/</u> <i>For details on the precise contents of each version of the Unicode Standard, and how to cite them.</i>		

Appendix D. Modifications

¶447 This section indicates the changes introduced by each revision.

448 *Revision 1.*

• **¶449** First version

¶450 Copyright © 2005-2012 Unicode, Inc. All Rights Reserved. The Unicode Consortium makes no expressed or implied warranty of

any kind, and assumes no liability for errors or omissions. No liability is assumed for incidental and consequential damages in connection with or arising out of the use of the information or programs contained or accompanying this technical report. The Unicode Terms of Use apply.

[451 Unicode and the Unicode logo are trademarks of Unicode, Inc., and are registered in some jurisdictions.