

Correction to L2/07-090, L2/07-129 on how Malayalam Chillus used in boundary cases

Cibu C. Johny
cibu@google.com

2007-10-15
 version 1.0

The need for correction

The recommendations in documents L2/07-090 and L2/07-129 call for using Virama immediately following a Chillu. However, in the Indic model, Virama acts as the vowel remover for a consonant with default vowel /a/. The Chillus does not have an inherent vowel. So <chillu, virama> sequence could be violating the Indic model.

This document retracts the above usage and lists how various visuals can still be achieved while being consistent to Indic model. Following sections describes the changes in Unicode 5.1 and clarifications after making the correction.

Changes incorporated in Unicode 5.1

Effective of Unicode 5.1, representation of Malayalam Chillus would change and the encoding of /nta/ conjunct and dot style Reph would be clarified. The new representation could be summarized as below:

- All the Chillus are atomically encoded
- The /nta/ conjunct is represented by <NA, Virama, RRA>
- Dot repħ is represented by <RA, Virama, Consonant> (no change from Unicode 5.0)

The tables 1, 2 and 3 detail these changes with examples. Left side of the tables describes what is being changed in Unicode 5.0 and right side describes what the final encodings in Unicode 5.1 would be.

What happens to existing sequence for Chillu?

Rendering systems should render <Consonant, Virama, ZWJ> as chillus, where consonants is NNA(0D23) NA(0D28) RA(0D30) LA(0D32) or LLA(0D33).

In the above codepoint context, ZWJ should not be stripped off as far as possible.

However, <Consonant, Virama, ZWJ> is highly discouraged from being used for chillu.

Changes related to representation of Chillus:

Table 1	Visual	Codepoint sequence in 5.0	Codepoint sequence in 5.1.0 (atomic encoding)	Sample Words
C1	അം	അ, ഃ, ZWJ NNA, VIRAMA, ZWJ (0D23, 0D4D, 200D)	അ (0D7A)	ആം /aann/
C2	അൻ	അ, ഃ, ZWJ NA, VIRAMA, ZWJ (0D28, 0D4D, 200D)	അൻ (0D7B)	അവൻ /avan/
C3	അയർ	അ, ഃ, ZWJ RA, VIRAMA, ZWJ (0D30, 0D4D, 200D)	അയർ (0D7C)	ഞായർ /njaayarr/
C4	അപ്പത്തി	അ, ഃ, ZWJ LA, VIRAMA, ZWJ (0D32, 0D4D, 200D)	അപ്പത്തി (0D7D)	ഉൽപ്പത്തി /ulpatti/
C5	അവൾ	അ, ഃ, ZWJ LLA, VIRAMA, ZWJ (0D33, 0D4D, 200D)	അവൾ (0D7E)	അവൾ /avall/
C6	ടക്സാക്സി	Not defined	ടക്സാക്സി (0D7F)	ടക്സാക്സി /drksaakssi/

Clarification on Conjunct rendering fallbacks:

Table 2	Levels of rendering fallbacks for the code sequence <C1, VIRAMA, C2> In the examples, C1 = മ (0D28) and C2 = റ (0D2E)					
	Unicode 5.0				Unicode 5.1	
	Microsoft Kartika Interpretation		Debian Rachana Interpretation			
	Visual	Example	Visual	Example	Visual	Example
F1	Conjunct of C1 and C2	മറ	Conjunct of C1 and C2	മറ	Conjunct of C1 and C2	മറ
F2	C1, Virama, C2	മ്‌റ	Chillu-C1, C2	അംര	C1, Virama, C2	മ്‌റ
F3			C1, Virama, C2	മ്‌റ		

Changes specific to Chillu NA(ன) and RRA(ஓ) combinations:

Table 3	Visuals	Shortest Codepoint sequence in 5.0		Sequence in 5.1.0	When joiner is lost	Fallback	Sample Words
		Microsoft Kartika interpretation	Rachana, Varamozhi, etc.				
N1		ன, சூ, ஓ (0D28, 0D4D, 0D31)	ன, சூ, ZNWJ, ஓ (0D28, 0D4D, 200C, 0D31)	ன, சூ, ஓ (0D28, 0D4D, 0D31)			உப்புமாவின்ராவு /uppuumaavint-rava/
N2		ன, சூ, ZWJ, ZWNJ, ஓ (0D28, 0D4D, 200D, 200C, 0D31)	ன, சூ, ZWJ, ஓ (0D28, 0D4D, 200D, 0D31)	ன, ஓ (0D7B, 0D4D, 0D31)			ஹென்றி /henri/
N3		Not defined		ன, சூ, ZWJ, ஓ (0D28, 0D4D, 200D, 0D31)	N4	N4	எந்தெ /ente/
N4		ன, சூ, ZWJ, ஓ (0D28, 0D4D, 200D, 0D31)	ன, சூ, ஓ (0D28, 0D4D, 0D31)	ன, சூ, ஓ (0D28, 0D4D, 0D31)		N3	எந்தெ /ente/

Position for the left part of the vowel sign, if used.



Following table indicates the changes related to dot-reph representation:

Table 3	#	Visual	Codepoint sequence in 5.0	Codepoint sequence in 5.1.0	Fallback	Sample Words
o(RA - 0D30) with ω(YA - 0D2F)	D1	ଓ	Not defined	ଓ, ঁ, ZWJ, ও (0D30, 0D4D, 200D, 0D2F)	D4	No examples found yet.
	D2	ঁও	ଓ, ঁ, ZWJ, ও (0D30, 0D4D, 200D, 0D2F)	ঁ, ও (0D7C, 0D2F)		চাতুর্যুগ /chaturyugam/
	D3	ଓঁ	ଓ, ঁ, ও (0D30, 0D4D, 0D2F)	ଓ, ঁ, ও (0D30, 0D4D, 0D2F)	D4	ବାର୍ଯ୍ୟ /bha:rya/
	D4	ଓঁও	ଓ, ঁ, ZWNJ, ও (0D30, 0D4D, 200C, 0D2F)	ଓ, ঁ, ZWNJ, ও (0D30, 0D4D, 200C, 0D2F)		No examples found yet.
	D5	ঁও	ଓ, ঁ, ও, ঁ, ও (0D30, 0D4D, 0D2F, 0D4D, 0D2F)	ଓ, ঁ, ও, ঁ, ও (0D30, 0D4D, 0D2F, 0D4D, 0D2F)	ଓঁও (D4)	ବାର୍ଯ୍ୟ /bha:rya/
	D6	ঁওঁ	ଓ, ঁ, ZWJ, ও, ঁ, ও (0D30, 0D4D, 200D, 0D2F, 0D4D, 0D2F)	ঁ, ও, ঁ, ও (0D7C, 0D2F, 0D4D, 0D2F)		No examples found yet.
o(RA - 0D30) with ω(VA - 0D35)	D7	ঁৱ	ଓ, ঁ, ঔ (0D30, 0D4D, 0D35)	ଓ, ঁ, ZWJ, ঔ (0D30, 0D4D, 200D, 0D35)	D10	ନିବାନ୍ଦମ୍
	D8	ঁঔ	ଓ, ঁ, ZWJ, ঔ (0D30, 0D4D, 200D, 0D35)	ঁ, ঔ (0D7C, 0D35)		ପାର୍ବତି /bha:rya/
	D9	ଓঁৱ	Not defined	ଓ, ঁ, ঔ (0D30, 0D4D, 0D35)	D10	ବାର୍ବୋ /varvo/ (colloquial)
	D10	ଓঁঔ	ଓ, ঁ, ZWNJ, ঔ (0D30, 0D4D, 200C, 0D35)	ଓ, ঁ, ZWNJ, ঔ (0D30, 0D4D, 200C, 0D35)		No examples found yet.
	D11	ঁৱ	ଓ, ঁ, ঔ, ঁ, ঔ (0D30, 0D4D, 0D35, 0D4D, 0D35)	ଓ, ঁ, ঔ, ঁ, ঔ (0D30, 0D4D, 0D35, 0D4D, 0D35)	ଓঁৱ (D10)	ପାର୍ବତି /bha:rya/
	D12	ঁঔ	ଓ, ঁ, ZWJ, ঔ, ঁ, ঔ (0D30, 0D4D, 200D, 0D35, 0D4D, 0D35)	ঁ, ঔ, ঁ, ঔ (0D7C, 0D35, 0D4D, 0D35)		ସର୍ବଔ /sarvvam/
RA with other consonants	D13	ঁঁ	ଓ, ঁ, ঔ	ଓ, ঁ, ঔ (0D30, 0D4D, 0D37)	ঁ	ବାର୍ଷମ୍ /varsham/

Rendering order in some of the edge case conjuncts

Following tables clarify the storage representation of specific Malayalam visuals for which the representation is not apparent from the Indic model. These visuals are based on the assumption that, the font has all the possible conjuncts. If the font lacks a specific conjunct, the fallback visual is also specified.

U-sign(0D41) variants and Samvruthokaram

Table 4	Visuals related to Samvrutho-karam	Storage Representation	Fallback	Sample Words	Notes
U1	മു OR മം	മ, മു (0D28, 0D41)	U5	മു/ /manu/	Rendering engine decides the visual
U2	മു	മ, ZWNJ, മു, മു (0D28, 200C, 0D41, 0D4D)	U2	അവന്മു /avani/	Requesting one specific visual
U3	മു	മ, ZWJ, മു, മു (0D28, 200C, 0D41, 0D4D)	U1	അവന്മു /avani/	Same reading as U2; Requesting one specific visual
U4	മു	മ, മു (0D28, 0D4D)	None	അവന്മു /avani/	
U5	മു	മ, ZWNJ, മു (0D28, 200C, 0D41)	U6	മു/ /manu/	
U6	മം	മ, ZWJ, മു (0D28, 200D, 0D41)	U5	മം/ /manu/	Same reading as U3

VA(0D35) and RA(0D30)-sign forms

Table 5	Visuals related conjoining α	Storage Representation	Fallback	Sample Words	Notes
V1	ତାନ୍ତି	ତ, ଅ, ନ୍ତି (0D28, 0D4D, 0D35)	ତ	ତନ୍ତି /tanvi/	
V2	ତାର୍ଵାମ	ତା, ର୍ଵା, ମ (0D34, 0D4D, 0D35)	V3	ତାର୍ଵାମ /ta:r'va:ram/	Same reading as V3. Fallback cannot be ର୍ଵା.
V3	ତାର୍ଵାମ	ତା, ର୍ଵା, ZWNJ, ମ (0D34, 0D4D, 200C, 0D35)	None	ତାର୍ଵାମ /ta:r'va:ram/	
V3	ତାର୍ଵୋ	ତା, ZWJ, ର୍ଵା, ମ (0D34, 200D, 0D4D, 0D35)	V3	ତାର୍ଵୋ	
V4	ଉଷ୍ମା	ଉ, ଷ୍ମା	ଉଷ୍ମା	ଉଷ୍ମା /uwwa/	
V5	ଉଷ୍ମା	ଉ, ZWJ, ଷ୍ମା	ଉଷ୍ମା	ଉଷ୍ମା /uw'wo:/	

Conjuncts, as in V2, is formed only in following regular expression context:

/[ଫଳ]େ(ଓ|ଓ)/. That is, in perl like syntax:

/ [\x{0D34} \x{0D2F}] \x{0D4D} (\x{0D35} | \x{0D30}) /.

Fallbacks for these forms will not be with VA or RA subjoining form. Instead, it will be with explicit virama as in V3.

Common conjunct that are often misrepresented in currently available fonts

Table 6	Common Conjuncts	Storage Representation	Sample Words
L1	ପଙ୍କ	ପ, ଙ୍କ, ଙ୍କ (0D19, 0D4D, 0D15)	ପଙ୍କ /panka/
L2	ଇଞ୍ଚ	ଇ, ଞ୍ଚ, ଞ୍ଚ (0D1E, 0D4D, 0D1A)	ଇଞ୍ଚ /incha/
L3	ମନ୍ଦ	ମ, ନ୍ଦ, ନ୍ଦ (0D23, 0D4D, 0D1F)	ମନ୍ଦ /manda/
L4	ମନ୍ତ	ମ, ନ୍ତ, ନ୍ତ (0D28, 0D4D, 0D24)	ମନ୍ତ /monta/
L5	ପଞ୍ପ	ପ, ଙ୍ପ, ଙ୍ପ (0D2E, 0D4D, 0D2A)	ପଞ୍ପ /pampa/