

# Malayalam: Adopting CLDR collation tailoring by DUCET

**Cibu Johny cibu@google.com**

**2009-01-27**

[This request is also filed in CLDR bug tracking as bug #2005[[4](#)]]

Malayalam is the only language that uses Malayalam script as its primary script[[1](#), [2](#)]. So it is might be prudent to incorporate the CLDR collation tailoring rules for Malayalam directly in DUCET. That would remove the need for using tailoring by software components to get dictionary order for Malayalam words.

## Rules

The proposed rules and their rationale are described below. Please see the appendices for these rules in specific syntax.

Malayalam Visarga and Avagraha may not have any primary weight:

The first primary ignorable << U+0D03  
U+0D03 << U+0D3D

The Malayalam AU length marker may differ from AU-sign only in tertiary weight and may be placed in between OO-sign and Samvruthokaram sign (U-Sign and Virama):

U+0D4B < U+0D57  
U+0D57 <<< U+0D4C  
U+0D57 < U+0D41 U+0D4D

Chillu encoding prevailing before 5.1 and their corresponding visible virama forms may differ only by secondary weights. Also, 5.1 Chillus and pre-5.1 Chillus may defer only by tertiary weights for existing data compatibility. Following 6 rules sets correspond to Chillus:

### NNA

U+0D23 U+0D4D << U+0D23 U+0D4D U+200D  
U+0D23 U+0D4D U+200D <<< U+0D7A

### NA

U+0D28 U+0D4D <<< U+0D7B U+0D4D  
/nta/ may be treated as <NA + VIRAMA + RRA>  
U+0D28 U+0D4D << U+0D28 U+0D4D U+200D  
U+0D28 U+0D4D U+200D <<< U+0D7B

### RA

U+0D30 U+0D4D << U+0D30 U+0D4D U+200D  
U+0D30 U+0D4D U+200D <<< U+0D7C

## LA

U+0D32 U+0D4D << U+0D32 U+0D4D U+200D  
U+0D32 U+0D4D U+200D <<< U+0D7D

## LLA

U+0D33 U+0D4D << U+0D33 U+0D4D U+200D  
U+0D33 U+0D4D U+200D <<< U+0D7E

## KA

U+0D15 U+0D4D << U+0D15 U+0D4D U+200D  
U+0D15 U+0D4D U+200D <<< U+0D7F

Anuswara may be considered as a MA\_dead:

U+0D2E U+0D4D << U+0D02

The order of last few letters may be: HA, LLA, LLLA, RRA

U+0D39 < U+0D33  
U+0D33 < U+0D34  
U+0D34 < U+0D31

Note: This includes the fix for bug #2000[3] proposed.

## Reference

1. [Search in Ethnologue for Malayalam script.](#)
2. [http://en.wikipedia.org/wiki/Malayalam\\_script](http://en.wikipedia.org/wiki/Malayalam_script)
3. <http://www.unicode.org/cldr/bugs/locale-bugs/incoming?id=2000>
4. <http://www.unicode.org/cldr/bugs/locale-bugs/incoming?id=2005>
5. [Online version of this document](#)

## Appendix A: The CLDR version of the rules

Please refer to online version for characters[5]. PDF version might miss some conjuncts and newly added Unicode characters.

```
<rules>
<reset><first_primary_ignorable/></reset>
<s>ঃ</s>
<reset>ঃ</reset>
<s>?</s>

<reset>ং</reset>
<p>ং</p>
<reset>ঁ</reset>
<t>ঁ</t>
<reset>ঁ</reset>
<p>ঁ</p>

<reset>ঃ</reset>
```

```

<s>ሙ </s>
<reset>ሙ </reset>
<t>?</t>

<reset>ወ </reset>
<t>?ወ</t>
<reset>ወ </reset>
<s>ወ </s>
<reset>ወ </reset>
<t>? </t>

<reset>ወ </reset>
<s>ወ </s>
<reset>ወ </reset>
<t>? </t>

<reset>ይ </reset>
<s>ይ </s>
<reset>ይ </reset>
<t>? </t>

<reset>ቃ </reset>
<s>ቃ </s>
<reset>ቃ </reset>
<t>? </t>

<reset>ቃ </reset>
<s>ቃ </s>
<reset>ቃ </reset>
<t>? </t>

<reset>ወ </reset>
<s>ወ </s>

<reset>ወ </reset>
<p>ወ</p>

<reset>ሂ </reset>
<p>ሂ</p>
</rules>

```

## Appendix B: The ICU version of the rules

&[first primary ignorable] << ◌፡  
&◌፡ << ?

&ርወ < ◌ወ

&ଓঁ <<< এঁও  
&ଓঁ < ঔৰ

&ମୀ << ମୀ  
&ମୀ <<< ?

&ନ୍ <<< ?୍  
&ନ୍ << ନ୍  
&ନ୍ <<< ?

&ର୍ << ର୍  
&ର୍ <<< ?

&ଏଁ << ଏଁ  
&ଏଁ <<< ?

&ଇଁ << ଇଁ  
&ଇଁ <<< ?

&କ୍ << କ୍  
&କ୍ <<< ?

&ମ୍ << ମ୍

&ହୋ<ହୋ  
&ହୋ<ଫୋ &ଫୋ

Experimentation  
link:  
[http://demo.icu-project.org/icu-bin/locexp?\\_=ml&d\\_=en&x=col](http://demo.icu-project.org/icu-bin/locexp?_=ml&d_=en&x=col)