

## Background of Indic segmentation

(Input for UAX#29 technical report)

India has large linguistic diversity with 22 constitutionally recognized languages and 12 scripts. The mapping between languages and scripts is complex. Multiple languages may have common scripts, while a language can be written in multiple scripts. Each language and script is unique in nature and cannot be easily replicated, even if they share common characteristics. The orthographic changes may also occur in some languages and adoption of new orthography is a gradual process, thus posing additional challenges. However for segmentation purpose, all Indic Languages having phonetic correspondence at the alphabet level by and large follow a common definition of syllable across languages.

### **ABNF Indic Orthographic syllable definition:**

**V[m] | {CH}C [v][m] | CH**

The definition is a combination of 3 rules:

Rule 1 : V[m]

Rule 2 : {CH}C[v][m]

Rule 3 : CH (This rule is applicable only at the end of the word)

V(upper case) is independent vowel

m is modifier(Anusvara/Visarga/Chandrabindu)

C is a consonant which may or may not include a single nukta

v (lower case) is any dependent vowel or vowel sign (mātrā)

H is halant / virama

| is a rule separator

[ ] - The enclosed items is optional under this bracket

{ } - The enclosed item/items occurs zero or repeated multiple times

Testing of Devanagari data (Hindi, Marathi, Konkani, Nepali) using Unicode segmentation utility tool was carried out. The cases where utility tool does not give correct segmentation break points are tabulated and the correct desired out as per ABNF Orthographic Indic syllable definition is given below.

### Devanagari :

Sl. No.	Rendering rules defined by Unicode	Words	Rendering as per Unicode segmentation utility tool ( <a href="http://unicode.org/cldr/utility/breaks.jsp">http://unicode.org/cldr/utility/breaks.jsp</a> )	Segmentation as per Indic syllable definition
<b>Devanagari Script</b>				
<b>Conjunct Formations</b>				
1.	Example :  ग् + ध → ग्ध  क् + ष → क्ष  क् + क → क्क	संदिग्ध  सुरक्षा  टक्कर  सत्र  दक्षिण  मिश्रणाने  मिनटापेक्षाही  कुरुक्षेत्रांतल्या	संदिग्ध  सुरक्षा  टक्कर  सत्र  दक्षिण  मिश्रणाने  मिनटापेक्षाही  कुरुक्षेत्रांतल्या	सं दि ग्ध Cm Cv CHC  सु र क्षा Cv C CHCv  ट क्क र C CHC C  स त्र C CHC  द क्षि ण  C CHCv C  मि श्र णा ने  Cv CHC Cv Cv  मि न टा पे क्षा ही Cv C Cv Cv CHCv Cv  कु रु क्षे त्रां त ल्या Cv Cv CHCv CHCvm C CHCv

		क्षेत्रहरूमा	क्षेत्रहरूमा	क्षेत्रहरूमा CHCv CHC C Cv Cv
<b>Rules for rendering</b>				
2.	If the dead consonant precedes a consonant, then it is replaced by the superscript nonspacing mark (U+0930 devanagari letter ra ) positioned above or attached to the upper part of a base glyph form.  Example : र + क → क + ँ	दुर्घटनाग्रस्त  मुहूर्त  शर्करायुक्त  वर्सातल्यान	दुर्घटनाग्रस्त  मुहूर्त  शर्करायुक्त  वर्सातल्यान	दुर्घटनाग्रस्त Cv CHC C Cv CHC CHC मुहूर्त Cv Cv CHC शर्करायुक्त C CHC Cv Cv CHC वर्सातल्यान C CHCvm C CHCv C
3.	For certain consonants, the mark (U+0930 devanagari letter ra positioned below or attached to the lower part of a base glyph form) may graphically combine with the consonant to form a conjunct ligature form.  Example : फ + र → फ + ्र	अंतरराष्ट्रीय  राष्ट्रपति  फ्रांस  ट्रेक्टर	अंतरराष्ट्रीय  राष्ट्रपति  फ्रांस  ट्रेक्टर	अंतरराष्ट्रीय Vm C C Cv CHCHCv C राष्ट्रपति Cv CHCHC C Cv फ्रांस CHCvm C ट्रेक्टर

		सिट्रोनेलाचे	सि ट्रो ने ला चे	CHCv CHC C सि ट्रो ने ला चे Cv CHCv Cv Cv Cv
		टुकुरालाई	टु कुरा ला ई	टु कुरा ला ई Cv CHCv Cv V
<b>Ligature rules</b>				
4.	If a dead consonant immediately precedes another dead consonant or a live consonant, then the first dead consonant may join the subsequent element to form a two-part conjunct ligature form  Example :  ट् + ठ → ढ	इकढा	इ क ढा	इ क ढा C C CHCv
5.	A conjunct ligature form can itself behave as a dead consonant and enter into further, more complex ligatures Example :  स् + त्र → स्त्र	शास्त्र  स्त्री	शा स्त्र   स त्री	शा स्त्र Cv CHCHC  स्त्री  CHCHCv
6.	A conjunct ligature form can also produce a half-form  Example :  क्ष् + य → क्ष्य	लक्ष्य  तीक्ष्ण	ल क्ष्य   ती क्ष्ण	ल क्ष्य C CHCHC  ती क्ष्ण Cv CHCHC

7.	<p>त् + र + ि → त्र + ि → त्रि</p>	<p>त्रिवेदी</p> <p>कृत्रिम</p> <p>मात्रामा</p>	<p>त्रि वेदी</p> <p>कृ त्रि म</p> <p>मा त्रा मा</p>	<p>त्रि वे दी</p> <p>CHCv Cv Cv</p> <p>कृ त्रि म</p> <p>Cv CHCv C</p> <p>मा त्रा मा</p> <p>Cv Cv Cv</p>
8.	<p>Example :</p> <p>द + ळ + ध → द्ध</p>	<p>सिद्धार्थनगर</p> <p>श्रद्धालुओं</p> <p>वृद्धिसँग</p>	<p>सि द्धार थ न गर </p> <p>श्र द्द धा लु ओं</p> <p>वृ द्धि सँ ग </p>	<p>सि द्धा र्थ न ग</p> <p>र</p> <p>Cv CHCv CHC C C</p> <p>C</p> <p>श्र द्धा लु ओं</p> <p>CHC CHCv Cv</p> <p>Vm</p> <p>वृ द्धि सँ ग</p> <p>Cv CHCv Cm C</p>
9.	<p>Modifying marks, apply to the orthographic syllable as a whole and should follow (in the memory representation) all other characters that constitute the syllable.</p> <p>क + ा + ँ → काँ</p>	<p>अंतःजानी</p> <p>स्वतःला</p>	<p>अं तः जा नी</p> <p>स्व तः ला</p>	<p>अं तः जा नी</p> <p>Vm Cm CHCv Cv</p> <p>स्व तः ला</p> <p>CHC Cm Cv</p>

**Recommendation:**

The pitfalls observed through data analysis above reveal that the Indic ABNF definition 100 % correct and needs to be suitably integrated in UAX# 29.