

TO: UTC

FROM: Swaran Lata, MeitY

SUBJECT: Recommendations to UTC #152 , July 31-August 4 2017 on Text segmentation in Indian languages

DATE: 16/06/2017

1. Background

- a) **Document:** [L2/16-161](#) : Proposal to add ABNF based Orthographic Indic syllable definition(appended at Annexure 'A') and guidelines for proper text segmentation in Indian languages and guiding principles of line breaking in Unicode text segmentation (UAX#29) . This document covers:
 - i. Indic Orthographic syllable definition
 - ii. Additional sample examples of tailored grapheme clusters
 - iii. Indic syllable boundary determination
 - iv. Examples of Initial letter styling and Vertical text segmentation
 - v. Guiding principles of line breaking in Indian languages
 - vi. The precise list of characters with their Unicode code points of all the categories of the characters of 9 scripts which constitute ABNF definition (Refer Annexure 'A')
- b) **Discussions on L2/16-161:** Presentation was made at UTC meeting #148 held on August 2-5, 2016 for highlighting the issues of text segmentation in Indian languages. It was suggested that the Unicode utility tool should be used to test data to identify the gaps in implementation of text segmentations for Indian languages and to report specific problems w.r.t UAX #29.

2. Action undertaken as a followup of UTC#148

- a) Using the utility tool of Unicode (<http://unicode.org/cldr/utility/breaks.jsp>), the sample data of Devanagari based languages[Hindi, Sanskrit (Excluding Vedic Extensions), Marathi, Konkani, Nepali, Bodo , Dogri , Maithili, Kashmiri, Sindhi] was used to test the existing UAX#29 rules viz - a - viz the proposed Orthographic Indic syllable definition (Refer Annexure 'A'). The test data mainly covers the combinations of conjunct characters, conjunct formations of consonants etc which sometimes don't get rendered properly.
- b) The output of the tool has been compared with correct segmentation expected as per Orthographic Indic syllable definition. It was observed that Rule 1 and Rule 3 are already implemented in the Unicode utility tool which is evident from the tested data presented at Annexure B. Following examples also illustrate this :

1. 
Vm C Cv C

2. 
V C CH

- c) The words covered in Annexure 'B' having CHC as a part are not being resolved properly by the Unicode Utility tool (Refer column 4) , which when segmented as per ABNF definition

gives the correct output (Refer column 5). This reveals improper implementation with reference to Rule 2 of ABNF.

3. Gaps in Utility tool implementation:

The following gaps in implementation of Rule 2 have been identified:

Improper break points of the following combinations:

- a. consonant + virama sequences + consonant (CHC)
- b. consonant + virama sequences + consonant + vowel sign (CHCv)
- c. consonant + virama sequences + consonant + modifier (CHCm)
- d. consonant + virama sequences + consonant + vowel sign + modifier (CHCvm)

The combinations at a, b are already covered in UAX#29 in the table 1a of section 3. The combinations at c & d need to examined and addressed.

4. Recommendations:

The data correlation was done for the above identified gaps. The errors due to non - implementation of Rule 2 in the Unicode Utility tool are marked in red colour. The corresponding output using ABNF Orthographic Indic syllable definition is shown in green colour. The corrective action w.r.t Rule 2 of ABNF implementation will require following additions in the UAX reports:

a) UAX#29 (UNICODE TEXT SEGMENTATION)

i. Proposed addition of Tailored Grapheme clusters :

The combinations at above defined 3(c) & 3(d) needs to be added as tailored grapheme cluster in the Table 1a.

| Examples | Characters | Comments |
|-------------------------------|---|-------------------------------|
| गौ (CHCm) Word : दुर्गौ | 0930 र DEVANAGARI LETTER RA 094D ॒ DEVANAGARI SIGN VIRAMA 0917 ग DEVANAGARI LETTER GA 0902 ० DEVANAGARI SIGN ANUSVARA | Combining character sequences |
| फ्रां(CHCvm) Word : फ्रांस | 092B फ DEVANAGARI LETTER PHA 094D ॒ DEVANAGARI SIGN VIRAMA 0930 र DEVANAGARI LETTER RA 093E ०॑ DEVANAGARI VOWEL SIGN AA 0902 ० DEVANAGARI SIGN ANUSVARA | Combining character sequences |

ii. Proposed addition of new section

The standard Indic syllable

- An V(Independent vowel) or a combination of V(Independent vowel) and m(modifier) is Indic syllable - V orVm
 - or
 - A sequences of consonant(C) + virama(H) followed by consonant(C) – CHC
 - or
 - A sequences of consonant(C) + virama(H) followed by consonant(C) with dependent vowel/modifier or both - CHCv, CHCm, CHCvm
 - or
 - Consonant with virama (applicable only at the end of the word) - CH

Indian languages syllable boundary determination:

Using the above Standard Indic syllable definition, the Indic syllable break never occurs between the pairs of characters as shown in the table below:

| Do not break between | | Examples |
|----------------------|---------|---|
| V | m | V X m |
| C | v, m | C X v, C X m, C X v X m |
| C | H | a) C X H X C, C X H X C X H X C, C X H X C X H X C X H X C x H x C b) C X H (applicable only at the end of the word in Sanskrit , Nepali, Dogri) |
| C | H, v, m | C X H X C X v, C X H X C X m, C X H X C X v X m, C X H X C X H X C X v, C X H X C X H X C X m , C X H X C X H X C X v v X m, C X H X C X H X C X H X C X v, C X H X C X H X C X v X m, C X H X C X H X C X m , C X H X C X H X C X v m, C X H X C X H X C X H X C X v , C X H X C X H X C X H X C X m , C X H X C X H X C X v X m |

Theoretically there is no practical limit of CHC, but the known limit is four in Sanskrit language (Excluding Vedic Extensions) and two in other languages. Example of maximum four CHC combinations is स्तन्यं is formed viz CHCHCHCHC.

b) UAX#14 (UNICODE LINE BREAKING ALGORITHM)

The definition of Standard Indic syllable (Refer UAX#29) may be used to break the line and a hyphen should be inserted at the breaking point so that word can be read intuitively however the language specific morpho-phonemic rules could be used for hyphenation as U+ 00AD (soft hyphen) is used in some Indian languages such as Tamil and Malayalam rather than inserting the hyphen character (U + 002D).

Annexure A

ABNF based Orthographic Indic syllable definition

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

This 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 etc)

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

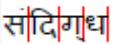
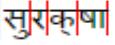
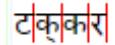
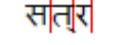
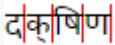
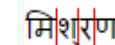
The Unicode code points of above categories of the characters is given below :

| S. No. | Scripts | Consonant(C) | Vowel Independent V | Vowel dependent v | Modifier m | Halant (Virama) H |
|--------|------------|---|---|---|--|--------------------------------|
| 1. | Devanagari | A) DEVANAGARI LETTER KA...DEVANAGARI LETTER HA (0915..0939) B) DEVANAGARI LETTER QA....DEVANAGARI LETTER YYA. (0958..095F) C) DEVANAGARI LETTER MARWARI DDA...DEVANAGARI LETTER BBA (0978..097F) | A) DEVANAGARI LETTER SHORT A...DEVANAGARI LETTER AU (0904..0914) B) DEVANAGARI LETTER VOCALIC RR...DEVANAGARI LETTER VOCALIC LL (0960..0961) C) DEVANAGARI LETTER CANDRA A...DEVANAGARI LETTER UUE (0972..0977) | A) DEVANAGARI VOWEL SIGN OE (093A) B) DEVANAGARI VOWEL SIGN OOE (093B) C) DEVANAGARI VOWEL SIGN AA...DEVANAGARI VOWEL SIGN II (093E...0940) D) DEVANAGARI VOWEL SIGN U...DEVANAGARI VOWEL SIGN AI (0941..0948) E) DEVANAGARI VOWEL SIGN CANDRA O...DEVANAGARI VOWEL SIGN AU (0949..094C) F) DEVANAGARI VOWEL SIGN PRISHTHAMATRA E...DEVANAGARI VOWEL SIGN AW (094E..094F) G) DEVANAGARI VOWEL SIGN CANDRA LONG E...DEVANAGARI VOWEL SIGN UUE (0955..0957) H) DEVANAGARI VOWEL SIGN VOCALIC L...DEVANAGARI VOWEL SIGN VOCALIC LL (0962..0963) | A) DEVANAGARI SIGN INVERTED CANDRABINDU. ..DEVANAGARI SIGN ANUSVARA (0900...0902) B) DEVANAGARI SIGN VISARGA (0903) C) DEVANAGARI SIGN AVAGRAHA (093D) D) DEVANAGARI SIGN NUKTA (093C) | DEVANA GARI SIGN VIRAMA (094D) |

Annexure B

Sample Test Results of 10 Indian languages(Devanagari script) based on Unicode Utility tool

1. Hindi

| Sl. No. | Rendering rules defined by Unicode | Words | Rendering as per Unicode segmentation utility tool (http://unicode.org/cldr/utility/breaks.jsp) | Correct Segmentation as per Indic syllable definition |
|----------------------------|---|--|---|--|
| Conjunct Formations | | | | |
| 2. | Example : $\text{ग} + \text{ध} \rightarrow \text{ग्ध}$ $\text{क} + \text{ष} \rightarrow \text{क्ष}$ $\text{क} + \text{क} \rightarrow \text{क्क}$ | संदिग्ध सुरक्षा टक्कर सत्र दक्षिण मिश्रणाने |  Cm Cv CH C  Cv C CH Cv  C CH C C  C CH C  C CH Cv C  Cv CH C Cv Cv | 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 |

| Rules for rendering | | | | |
|---------------------|--|--|--|---|
| 3. | 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 : $\text{र} + \text{क} \rightarrow \text{क} + \text{ं}$ | दुर्घटनाग्रस्त मुहूर्त शक्तरायुक्त क्रांस ट्रैक्टर सिट्रोनेलाचे | दुर्घटनाग्रस्त Cv CH C C Cv CH C CH C मुहूर्त Cv Cv CH C शक्तरायुक्त C CH C Cv Cv CH C फ्रांस CH Cvm C ट्रैक्टर CH Cv CH C C सिट्रोनेलाचे | दु र्घ ट ना ग्र स्त Cv CHC C Cv CHC CHC मु हू र्त Cv Cv CHC श क्त रा यु क्त C CHC Cv Cv CHC अं त र रा ष्ट्री य Vm C C Cv CH CH Cv C रा ष्ट्रपति Cv CH CH C C Cv फ्रां स CHCvm C ट्रै क्ट र CHCv CHC C सि ट्रो ने ला चे |
| 4. | 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 : $\text{फ} + \text{र} \rightarrow \text{फ} +$ | अंतरराष्ट्रीय राष्ट्रपति क्रांस ट्रैक्टर सिट्रोनेलाचे | अंतरराष्ट्रीय Vm C C Cv CH CH Cv C रा ष्ट्रपति Cv CH CH C C Cv फ्रांस CH Cvm C ट्रैक्टर CH Cv CH C C सिट्रोनेलाचे | अं त र रा ष्ट्री य Vm C C Cv CHCHCv C रा ष्ट्र प ति Cv CHCHC C Cv फ्रां स CHCvm C ट्रै क्ट र CHCv CHC C सि ट्रो ने ला चे |

| | | | | |
|--|--|-----------|--|--|
| | | टुक्रालाई | Cv CH Cv Cv Cv Cv टु क्रा ला ई Cv CH Cv Cv V | Cv CH Cv Cv Cv Cv टु क्रा ला ई Cv CH Cv Cv V |
|--|--|-----------|--|--|

5. Ligature rules

| | | | | |
|----|---|---------|--|--|
| 6. | 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 : $\text{ट} + \text{ठ} \rightarrow \text{ड्ठ}$ | इकट्ठा | इ कट्ठा C C CH Cv | इ क ट्ठा C C CH Cv |
| 7. | A conjunct ligature form can itself behave as a dead consonant and enter into further, more complex ligatures Example : $\text{स} + \text{त्र} \rightarrow \text{स्त्र}$ | शास्त्र | शा सृत्र Cv CH CH C सृत्री CH CH Cv | शा स्त्र Cv CHCHC स्त्री CHCHCv |

| | | | | |
|-----|--|--|--|---|
| | | | | |
| 8. | A conjunct ligature form can also produce a half-form Example : $\text{क्ष} + \text{य} \rightarrow \text{क्ष्य}$ | लक्ष्य तीक्ष्ण | ल क् ष्य C CH CH C ती क् ष्ण Cv CH CH C | ल क्ष्य C CHCHC ती क्ष्ण Cv CHCHC |
| 9. | $\text{त} + \text{र} + \text{ि} \rightarrow \text{त्र} + \text{ि} \rightarrow \text{त्रि}$ | त्रिवेदी कृत्रिम मात्रामा | त रि वेदी CH Cv Cv Cv कृ त्रि म Cv CH Cv C मा त्रा मा Cv CH Cv Cv | त्रि वे दी CH Cv Cv Cv कृ त्रि म Cv CH Cv C मा त्रा मा Cv CH Cv Cv |
| 10. | Example : $\text{द} + \text{्} + \text{ध} \rightarrow \text{द्ध}$ | सिद्धार्थन गर श्रद्धालुओं वृद्धिसँग | सि द् धा र् था न गर Cv CH Cv CH C C C श्रा द् धा लु ओ CH C CH Cv Cv Vm वृ द् धि सँ ग Cv CH Cv Cm C | सि द्धा र्थ न ग र Cv CH Cv CH C C C श्र द्धा लु ओ CH C CH Cv Cv Vm वृ द्धि सँ ग Cv CH Cv Cm C |

| | | | | |
|-----|--|------------|------------------------------|------------------------------|
| 11. | Modifying marks, apply to the orthographic syllable as a whole and should follow (in the memory representation) all other characters that constitute the syllable. | अंतःज्ञानी | अंतःज्ञानी Vm Cm CH Cv Cv | अं तः जा नी Vm Cm CHCv Cv |
| | क + ०९ + ० → क९ | | | |

2. Sanskrit(Excluding Vedic Extensions)

| Words | Rendering as per Unicode segmentation utility tool (http://unicode.org/cldr/utility/breaks.jsp) | Correct Segmentation as per Indic syllable definition |
|------------|--|---|
| गन्नदीधिम् | ग न्न दी धि म् | ग न्न दी धि म C CHC Cvm Cv C |
| प्रपंदये | प्र पं पद ये | प्र पं दये CHCHC Cm CHCvm |
| मनस्तापः | म न स तापः | म न स्ता पः C C CHCv Cm |
| हविष्करोमि | ह वि ष्क करो मि | ह वि ष्क रो मि C Cv CHC Cv Cv |
| अहर्पतिः | अ हर पति ः | अ ह र्प तिः V C CHC Cvm |
| गच्छति | ग च च्छ ति | ग च्छ ति C CHC Cv |
| अयम् | अ य म् | अ य म् V C CH |
| शिवश्चोदति | शि व शु चो दति | शि व श्चो द ति Cv C CHCv C Cv |
| मनष्टालयति | म न ष टाल यति | म न ष्टा ल य ति |

| | | |
|-------------|-----------------------|-------------------|
| | C C CH Cv C C Cv | C C CHCv C C Cv |
| अश्वष्ठकस्य | आ श्व ष्ठ क स्य | अ श्व ष्ठ क स्य |
| | V CH C CH C CH C CH C | V CHC CHC CHC CHC |
| दुष्पुत्रः | दु ष्पु त्रः | दु ष्पु त्रः |
| | Cv CH Cv CH Cm | Cv CHCv CHCm |
| द्विःपक्वम् | द ्वि ः पक्व म् | द्विः पक्व म् |
| | CH Cvm C CH C CH | CHCvm C CHC CH |
| द्विष्कामः | द ्वि ष्का मः | द्वि ष्का मः |
| | CH Cv CH Cv Cm | CHCv CHCv Cm |
| भर्तुर्भोगः | भर तु रभो गः | भ तु भो गः |
| | C CH Cv CH Cv Cm | C CHCv CHCv Cm |

3. Kashmiri

| Words | Rendering as per Unicode segmentation utility tool (http://unicode.org/cldr/utility/breaks.jsp) | Correct Segmentation as per Indic syllable definition |
|----------|---|---|
| शेत्युल | शे त्यु ल Cv CH Cv C | शे त्यु ल Cv CHCv C |
| महारैन्य | महा रै न्य C Cv Cv CH C | म हा रै न्य C Cv Cv CHC |
| सेक्युल | से क्यु ल Cv CH Cv C | से क्यु ल Cv CHCv C |
| ल्योदुर | ल्यो दु र CH Cv Cv C | ल्यो दु र CHCv Cv C |
| फ्योक | फ ्यो क CH Cv C | फ्यो क CHCv C |
| मोहन्युव | मो हन्यु व Cv C CH Cv C | मो ह न्यु व Cv C CHCv C |
| अन्यर | अ न्यर | अ न्य र |

| | V CH C C | V CHC C |
|----------|----------------------------|--------------------------|
| खवजि | ख व जि CH C Cv | ख्व जि CHC Cv |
| ख्युज | ख यु ज CH Cv C | ख्यु ज CHCv C |
| ख्योमुत | ख यो मु त CH Cv Cv C | ख्यो मु त CHCv Cv C |
| खैरिन्य | ख ैर न् य Cv Cv CH C | खै रि न्य Cv Cv CHC |
| र्येंड | र ्यें ड CH Cv C | र्यें ड CHCv C |
| छन्यर | छ न् या र C CH C C | छ न्य र C CHC C |
| छिर्युव | छि र यु व Cv CH Cv C | छि र्यु व Cv CHCv C |
| जॉल्य | जौ ल य Cv CH C | जॉ ल्य CH CHC |
| छवम्बुन | छ व म बु न CH C CH Cv C | छ्व म्बु न CHC CHCv C |
| डॉगिन्य | डौ गि न् य Cv Cv CH C | डॉ गि न्य Cv Cv CHC |
| टिक्युल | टि क यु ल Cv CH Cv C | टि क्यु ल Cv CHCv C |
| तिक्याजि | ति क या जि Cv CH Cv Cv | ति क्या जि Cv CHCv Cv |
| थ्यकुन | थ य कु न CH C Cv C | थ्य कु न CHC Cv C |
| त्यम्बर | त य म बु र Cv CH Cv C | त्य म्ब र Cv CHCv C |

| | | |
|--------|-------------|-----------|
| | CH C CH C C | CHC CHC C |
| थवसिलद | थ व सि ल द | थव सि ल द |

CH C Cv C C

CHC Cv C C

4. Marathi

| Words | Rendering as per Unicode segmentation utility tool (http://unicode.org/cldr/utility/breaks.jsp) | Correct Segmentation as per Indic syllable definition |
|----------------|---|---|
| उच्छ्वास | उ च्छ वा स | उ च्छ्वा स |
| | V CH CH Cv C | V CHCHCv C |
| व्यक्तिमत्व | व्य कृ ति मत्व | व्य कृति मत्व |
| | CH C CH Cv C CH C | CHC CHCv C CHC |
| आपल्या | आ प ल या | आ प ल्या |
| | V C CH Cv | V C CHCv |
| दातांच्यामध्ये | दा ता च्या मध्ये | दा तां च्या मध्ये |
| | Cv Cvm CH Cv C CH Cv | Cv Cvm CHCv C CHCv |
| टिपांच्या | टि पा च्या | टि पां च्या |
| | Cv Cvm CH Cv | Cv Cvm CHCv |
| ह्याच्यासाठी | ह्या च्या सा ठी | ह्या च्या सा ठी |
| | CH Cv CH Cv Cv Cv | CHCv CHCv Cv Cv |
| मिनटापेक्षाही | मि न टा पे क्षा ही | मि न टा पे क्षा ही |
| | Cv C Cv Cv CH Cv Cv | Cv C Cv Cv CHCv Cv |
| दुर्गंधी | दु र्ग ं धी | दु र्गं धी |
| | Cv CH Cm Cv | Cv CHCm Cv |
| शर्करायुक्त | श रकरा युक्त | श करा युक्त |
| | C CH C Cv Cv CH C | C CHC Cv Cv CHC |
| त्यामध्ये | त ्या मध्ये | त्या मध्ये |

| | CH Cv C CH Cv | CHCv C CHCv |
|---------------------|---|---|
| दुर्गंधीपासूनसुदृढा | दुर् गं धी पा सू न सु दृ ढा Cv CH Cm Cv Cv CV C Cv CH Cv | दुर्गं धी पा सून सुदृढा Cv CHCm Cv Cv CV C Cv CHCv |
| कोलेस्टेरॉलला | को ले सु टे रॉ ल ला Cv Cv CH Cv Cv C Cv | को ले स्टे रॉ ल ला Cv Cv CHCv Cv C Cv |
| नित्यकर्मामध्ये | नि त्य करु मा मध्ये Cv CH C C CH Cv C CH Cv | नि त्य करु मा मध्ये Cv CHC C CHCv C CHCv |
| मिश्रणाने | मि श्र णा ने Cv CH C Cv Cv | मि श्रणा ने Cv CHC Cv Cv |
| प्रत्येकवेळी | प्र त्ये कवे ळी CH C CH Cv C Cv Cv | प्रत्येकवेळी CHC CHCv C Cv Cv |
| ठेवण्यामध्ये | ठे वण ्या मध्ये Cv C CH Cv C CH Cv | ठेवण्या मध्ये Cv C CHCv C CHCv |
| मिश्रणाने | मि श्र णा ने Cv CH C Cv Cv | मि श्रणा ने Cv CHC Cv Cv |
| मधुमेहासारख्या | म धु मेहा सा रख्या C Cv Cv Cv Cv C CH Cv | मधुमेहा सारख्या C Cv Cv Cv Cv C CHCv |
| नित्यकर्मामध्ये | नि त्य करु मा मध्ये Cv CH C C CH Cv CH Cv | नि त्य करु मा मध्ये Cv CHC C CHCv CHCv |
| आजारांपासूनसुदृढा | आ जा रा पा सू न सु दृ ढा V Cv Cvm Cv Cv C Cv CH Cv | आजारांपासूनसुदृढा V Cv Cvm Cv Cv C Cv CHCv |
| समाविष्ट | स मा वि ष्ट C Cv Cv CH C | समाविष्ट C Cv Cv CHC |
| भाज्यांमध्ये | भा ज्या मध्ये Cv CH Cvm C CH Cv | भाज्यां मध्ये Cv CHCv C CHCv |
| स्वतःला | स वत ः ला CH C Cm Cv | स्वतःला CHC Cm Cv |
| उच्छ्वासाच्या | उ च्छ ्वा सा च्या | उच्छ्वासाच्या |

| | | |
|---------------|---|--|
| | C CH CH Cv Cv CH Cv | C CHCHCv Cv CHCv |
| मिनटापेक्षाही | मि न टा पे क्‌षा ही Cv C Cv Cv CH Cv Cv | मि न टा पे क्षा ही Cv C Cv Cv CHCv Cv |

5. Konkani

| Words | Rendering as per Unicode segmentation utility tool (http://unicode.org/cldr/utility/breaks.jsp) | Correct Segmentation as per Indic syllable definition |
|--------------|---|---|
| सुकिल्लीं | सु कि ल्‌लीं Cv Cv CH Cvm | सु कि ल्लीं Cv Cv CHCvm |
| कुडक्यांक | कु ड क्यां क Cv C CH Cvm C | कु ड क्यां क Cv C CHCvm C |
| रॅवपाच्या | रॅ व पा च्या Cvm C Cv CH Cv | रॅ व पा च्या Cvm C Cv CHCv |
| रॅप्याचो | रॅ प् या चो Cvm CH Cv Cv | रॅ प्या चो Cvm CHCv Cv |
| स्लिपां | सु लि पां CH CV Cvm | स्लि पां CHCv Cvm |
| सिट्रोनेलाचे | सि ट्रो ने ला चे Cv CH Cv Cv Cv Cv | सि ट्रो ने ला चे Cv CHCv Cv Cv Cv |
| मिसळिल्ल्यान | मि स ळि ल्‌ल्‌या न Cv C Cv CH CH Cv C | मि स ळि ल्ल्या न Cv C Cv CHCHCv C |
| रॅप्यांची | रॅ प् या ची Cvm CH Cvm Cv | रॅ प्यां ची Cvm CHCvm Cv |
| खोडाच्या | खो डा च्या Cv Cv CH Cv | खो डा च्या Cv Cv CHCv |
| रॅप्यांक | रॅ प् या क Cvm CH Cvm C | रॅ प्यां क Cvm CHCvm C |

| | | |
|--------------|---|-------------------------------------|
| वर्सातल्यान | वर् सा तल् या न C CH Cvm C CH Cv C | व र्सा त ल्या न C CHCvm C CHCv C |
| म्हण्टात | म् हृ ण् टा त CH C CH Cv C | म्ह ण्टा त CHC CHCv C |
| कृत्रीम | कृ त्री म Cv CH Cv C | कृ त्री म Cv CHCv C |
| सगल्यांत | स गल् या ं त C C CH Cvm C | स ग ल्यां त C C CHCvm C |
| रोप्याच्या | रो प्या च् या Cvm CH Cv CH Cv | रों प्या च्या Cvm CHCv CHCv |
| संस्थानान | सं स् था ना न Cm CH Cv Cv C | सं स्था ना न Cm CHCv Cv C |
| केल्ल्यान | के ल् ल् या न Cv CH CH Cv C | के ल्ल्या न Cv CHCHCv C |
| म्हळ्यार | म् हृ ळ्या र CH C CH Cv C | म्ह ळ्या र CHC CHCv C |
| सगल्यांत | स गल् या ं त C C CH Cvm C | स ग ल्यां त C C CHCvm C |
| नाशिल्ल्यान | ना शि ल् ल् या न Cv Cv CH CH Cv C | ना शि ल्ल्या न Cv Cv CHCHCv C |
| म्हटलां | म् हृ ट ला ं CH C C Cvm | म्ह ट लां CHC C Cvm |
| शेतकर्या नी | शे तकर् या चनी Cv C C CH Cvm Cv | शे त कर्या नी Cv C C CHCvm Cv |
| कुडक्यांनी | कु डक् या ं नी Cv C CH Cvm Cv | कु ड क्यां नी Cv C CHCvm Cv |
| जिल्ल्याच्या | जि ल् ल् या च् या Cv CH CH Cv CH Cv | जि ल्ल्या च्या Cv CHCHCv CHCv |

| | | |
|--------------------|--|--|
| ब्लॉकांत | ब्लॉका त CH Cv Cvm C | ब्लॉ कां त CHCv Cvm C |
| कुरुक्षेत्रांतल्या | कुरुक्षे त्रा तल्या Cv Cv CH Cv CH Cvm C CH Cv | कु रु क्षे त्रा त ल्या Cv Cv CHCv CHCvm C CHCv |
| भाज्ज्यांची | भा ज्ज्या ची Cv CH CH Cvm Cv | भा ज्ज्यां ची Cv CHCHCvm Cv |

6. Nepali

| Words | Rendering as per Unicode segmentation utility tool (http://unicode.org/cldr/utility/breaks.jsp) | Correct Segmentation as per Indic syllable definition |
|--------------|---|---|
| काटिन्छ | का टि न्छ Cv Cv CH C | का टि न्छ Cv Cv CHC |
| त्यसलाई | त्य स ला ई CH C C Cv V | त्य स ला ई CHC C Cv V |
| भनिन्छ | भ नि न्छ C Cv CH C | भ नि न्छ C Cv CHC |
| स्लिप्स | स लि प्स CH Cv CH C | स्लिप्स CHCv CHC |
| टुक्रामा | टु क्रा मा Cv CH Cv Cv | टु क्रा मा Cv CHCv Cv |
| बिरुवामध्ये | बि रुवा मध्ये Cv Cv Cv C CH Cv | बि रु वा मध्ये Cv Cv Cv C CHCv |
| टुक्रालाई | टु क्रा ला ई Cv CH Cv Cv V | टु क्रा ला ई Cv CHCv Cv V |
| छिद्रमा | छि द्र मा Cv CH C Cv | छि द्र मा Cv CHC Cv |
| सिट्रोनेलाका | सि ट्रो नेला का Cv CH Cv Cv Cv Cv | सि ट्रो ने ला का Cv CHCv Cv Cv Cv |
| त्यसो | त्य सो | त्य सो |

| | | |
|---------------|---|---|
| | CH C Cv मा तरा मा Cv CH Cv Cv | CHC Cv मा त्रा मा Cv CHCv Cv |
| गरिनुपछ | गरि नुप रछ C Cv Cv C CH C | ग रि नु प छे C Cv Cv C CHC |
| महिनाभन्दा | महि ना भन् दा C Cv Cv C CH Cv | म हि ना भ न्दा C Cv Cv C CHCv |
| एकपल्ट | एक पल ट V C C CH C | ए क प ल्ट V C C CHC |
| अतिरिक्त | अति रिक् त V Cv Cv CH C | अ ति रि क्ता V Cv Cv CHC |
| सुक्नु | सु क्नु Cv CH Cv | सु क्नु Cv CHCv |
| सामान्यतः | सा मा न्य तः Cv Cv CH C Cm | सा मा न्य तः Cv Cv CHC Cm |
| सिम्बोपोगान | सि म्बो पो गा न Cv CH Cv Cv Cv C | सि म्बो पो गा न Cv CHCv Cv Cv C |
| कृत्रिम | कृ त्रि म Cv CH Cv C | कृ त्रि म Cv CHCv C |
| प्रकारको | प्र कार को CH C Cv C Cv | प्र का र को CHC Cv C Cv |
| मध्यस्थकारीले | मध्य स्थ का री ले C CH C CH C Cv Cv Cv | म ध्य स्थ का री ले C CHC CHC Cv Cv Cv |
| क्षेत्रहरूमा | क्षे त्रहरू मा CH Cv CH C C Cv Cv | क्षे त्र ह रू मा CHCv CHC C Cv Cv |
| वृद्धिसँग | वृद्धि सँ ग Cv CH Cv Cm C | वृ द्धि सँ ग Cv CHCv Cm C |

7. Bodo

| Words | Rendering as per Unicode segmentation utility tool (http://unicode.org/cldr/utility/breaks.jsp) | Correct Segmentation as per Indic syllable definition |
|---------------|---|---|
| बाराद्राय | बा रा द्रा य Cv Cv CH Cv C | बा रा द्रा य Cv Cv CHCv C |
| खान्थियाव | खा न्थि या व Cv CH Cv Cv C | खा न्थि या व Cv CHCv Cv C |
| सरासनस्तायै | सरा सन स्ता यै C Cv C C CH Cv Cv | स रा स न स्ता यै C Cv C C CHCv Cv |
| रांखान्थियारि | रा ंखा न्थि या रि Cvm Cv CH Cv Cv Cv | रां खा न्थि या रि Cvm Cv CHCv Cv Cv |
| सिम्बपगान | सि म्ब पगा न Cv CH C C Cv C | सि म्ब प गा न Cv CHCv C Cv C |
| ग्रामनी | ग्रा रा मनी CH Cv C Cv | ग्रा म नी CHCv C Cv |
| बयनिखिउङ | ब य नि खु इङ् C C Cv CH Cv V | ब य नि खु ङ् C C Cv CHCv V |
| हाग्रा | हा ग्रा Cv CH Cv | हा ग्रा Cv CHCv |
| बोसोरब्रै | बो सो रब्रै Cv Cv C CH Cv | बो सो र ब्रै Cv Cv C CHCv |
| दिहुनग्राया | दि हु नग्रा या Cv Cv C CH Cv Cv | दि हु न ग्रा या Cv Cv C CHCv Cv |
| प्रभुआ | प्र भु आ CH C Cv V | प्र भु आ CHCv Cv V |
| बेनिफ्राय | बे नि फ्रा य Cv Cv CH Cv C | बे नि फ्रा य Cv Cv CHCv C |
| दैथाइहरगा | दै था इहरगा Cv Cv C C C CH Cv | दै था इ हरगा Cv Cv C C C CHCv |
| झज्जरनि | झ ज्जरनि Cv Cv C C C | झ ज्जरनि |

| | | |
|-----------|-------------------------------|------------------------------|
| | C CH C C Cv | C CHC C Cv |
| फार्माव | फा र्मा व Cv CH Cv C | फा मा व Cv CHCv C |
| थाग्रा | था ग्रा Cv CH Cv | था ग्रा Cv CHCv |
| हानिफ्राय | हा नि फ्रा य Cv Cv CH Cv C | हा नि फ्रा य Cv Cv CHCv C |
| रिपोर्टनि | रि पो रटा नि Cv Cv CH C Cv | रि पो ट नि Cv Cv CHC Cv |
| बादिब्ला | बा दि ब्ला Cv Cv CH Cv | बा दि ब्ला Cv Cv CHCv |
| जौखोन्दो | जौ खो न्दो Cv Cv CH Cv | जौ खो न्दो Cv Cv CHCv |
| जौखोन्दोल | जौ खो न्दो ल Cv Cv CH Cv C | जौ खो न्दो ल Cv Cv CHCv C |

8. Dogri

| Words | Rendering as per Unicode segmentation utility tool (http://unicode.org/cldr/utility/breaks.jsp) | Correct Segmentation as per Indic syllable definition |
|-----------|--|---|
| मनुक्खो | म नु क्खो C Cv CH Cv | म नु क्खो C Cv CHCv |
| मितर | मि त तर Cv CH C C | मि त र Cv CHC C |
| दिक्खी | दि क् खी Cv CH Cv | दि क्खी Cv CHCv |
| टोस्टर्ज़ | टो स टर ज़ Cv CH C CH C | टो स्ट झ़ Cv CHC CHC |

| | | |
|-----------|-------------------------------|----------------------------|
| आहन्नो | आ हन् नो V C CH Cv | आ ह न्नो V C CHCv |
| वक्तव्य | व कृत वय C CH C CH C | व क्त व्य C CHC CHC |
| पतर | पत तर C CH C C | प त र C CHC C |
| पुश्टी | पु श्टी Cv CH Cv | पु श्टी Cv CHCv |
| प्रक्रिया | पर क्रि या CH C CH Cv Cv | प्र क्रि या CHC CHCv Cv |
| प्रविश्ट | पर विश्ट CH C Cv CH C | प्र वि श्ट CHC Cv CHC |
| बक्खरा | बक् खरा C CH C Cv | ब क्ख रा C CHC Cv |
| सक्खना | सक् खना C CH C Cv | स क्ख ना C CHC Cv |
| खोहल्लना | खो हल् लना Cv C CH C Cv | खो ह ल्ल ना Cv C CHC Cv |

9. Maithili

| Words | Rendering as per Unicode segmentation utility tool (http://unicode.org/cldr/utility/breaks.jsp) | Correct Segmentation as per Indic syllable definition |
|-----------|--|---|
| निर्दिष्ट | नि र्दि ष्ट Cv CH Cv CH C | नि दि ष्ट Cv CHCv CHC |
| भ्रमणक | भ्र मणक CH C C C C | भ म ण क CHC C C C |
| कूटशब्दक | कू टशब्दक Cv C C CH C C | कू ट श ब्द क Cv C C CHC C |

| | | |
|------------|----------------------------------|----------------------------------|
| तालाबन्न | ता ला बा न्न Cv Cv C CH C | ता ला ब न्न Cv Cv C CHC |
| अस्वीकृत | अ स्वी कृ त V CH v Cv C | अ स्वी कृ त V CHv Cv C |
| संख्यासँ | सं ख्या सँ Cm CH Cv Cm | सं ख्या सँ Cm CHCv Cm |
| श्रेणीक | श्रे णी क CH Cv Cv C | श्रे णी क CHCv Cv C |
| प्रयोक्ता | प्र यो कृ ता CH C Cv CH Cv | प्र यो कता CHC Cv CHCv |
| निर्गम | नि रगम Cv CH C C | नि र्ग म Cv CHC C |
| फॅटवर्क | फॅ टवर्क Cm C C CH C | फॅ ट व र्क Cm C C CHC |
| अनुच्छेदसँ | अ नु च्छे दसँ V Cv CH Cv C Cm | अ नु च्छे द सँ V Cv CHCv C Cm |
| दहिन्ना | द हि न्ना C Cv CH Cv | द हि न्ना C Cv CHCv |
| स्वचालित | स्व चा लित CH C Cv Cv C | स्व चा लि त CHC Cv Cv C |
| सँपत्तिक | सँ पत्ति क Cm C CH Cv C | सँ प ति क Cm C CHCv C |
| निर्गत | नि रगत Cv CH C C | नि र्ग त Cv CHC C |
| प्रकार्यक | प्र कार यक CH C Cv CH C C | प्र का र्य क CHC Cv CHC C |
| संख्याक | सं ख्या क Cm CH Cv C | सं ख्या क Cm CHCv C |
| स्थानधारक | सं था नधा रक Cm CH Cv C | स्था न धा र क Cm CHCv C |

| | | |
|----------|-----------------------|----------------------|
| | CH Cv C Cv C C | CHCv C Cv C C |
| भविष्यके | भवि ष्यके | भ वि ष्य के |
| | C Cv CH C Cvm | C Cv CHC Cvm |

| | | |
|--|-------------------|------------------|
| | पैकितक | पै कित क |
| | Cm CH Cv C | Cm CHCv C |

10.Sindhi

| Words | Rendering as per Unicode segmentation utility tool (http://unicode.org/cldr/utility/breaks.jsp) | Correct Segmentation as per Indic syllable definition |
|-----------|---|---|
| ईश्वरु | ई श्वरु | ई श रु |
| | V CH C Cv | V CHC Cv |
| गर्मीपटु | गर मी पटु | ग मी प टु |
| | C CH Cv C Cv | C CHCv C Cv |
| ट्रेक्टरु | ट्र एक्टरु | ट्रे क्ट रु |
| | CH Cv CH C Cv | CHCv CHC Cv |
| कलेक्टरु | का ले क्टरु | क ले क्ट रु |
| | C Cv CH C Cv | C Cv CHC Cv |
| ग्र्यानु | ग्र या नु | ग्या नु |
| | CH Cv Cv | CHCv Cv |
| शास्त्रु | शा सत्रु | शा स्त्रु |
| | Cv CH CH Cv | Cv CHCHCv |
| शिक्षत | शि क्षत | शि क स्त |
| | Cv C CH C | Cv C CHC |
| इन्द्री | इ नद्री | इ न्द्री |
| | V CH CH Cv | V CHCHCv |
| श्रधालू | श्र धा लू | श्र धा लू |
| | CH C Cv Cv | CHC Cv Cv |
| सकार्थो | स कार थो | स का थो |
| | C Cv CH Cv | C Cv CHCv |

| | | |
|------------|---------------------------------|-------------------------------|
| मुछित्यारी | मु ख ति या री Cv CH Cv Cv Cv | मु छित या री Cv CHCv Cv Cv |
| बेवक्ति | बे व कृ ति Cv C CH Cv | बे व क्रित Cv C CHCv |
| आमदरफ्त | आ म द र रा फृ त V C C C CH C | आ म द र फ्त V C C C CHC |
| आस्तिकु | आ स ति कु V CH Cv Cv | आ स्ति कु V CHCv Cv |
| इन्किलाबी | इ न कि ला बी V CH Cv Cv Cv | इ न्कि ला बी V CHCv Cv Cv |
| इब्तिदा | इ बृ ति दा V CH Cv Cv | इ ब्ति दा V CHCv Cv |
| यजु | य ज जु C CH Cv | य जु C CHCv |
| सन्बंधु | स न बृ धु C CH Cm Cv | स न्बं धु C CHCm Cv |
| त्रिलोकु | तु रि लो कु CH Cv Cv Cv | त्रि लो कु CHCv Cv Cv |
| इस्पंगुरु | इ स पृ गु रु V CH Cm Cv Cv | इ स्पं गु रु V CHCm Cv Cv |