

# INDIC TEXT SEGMENTATION

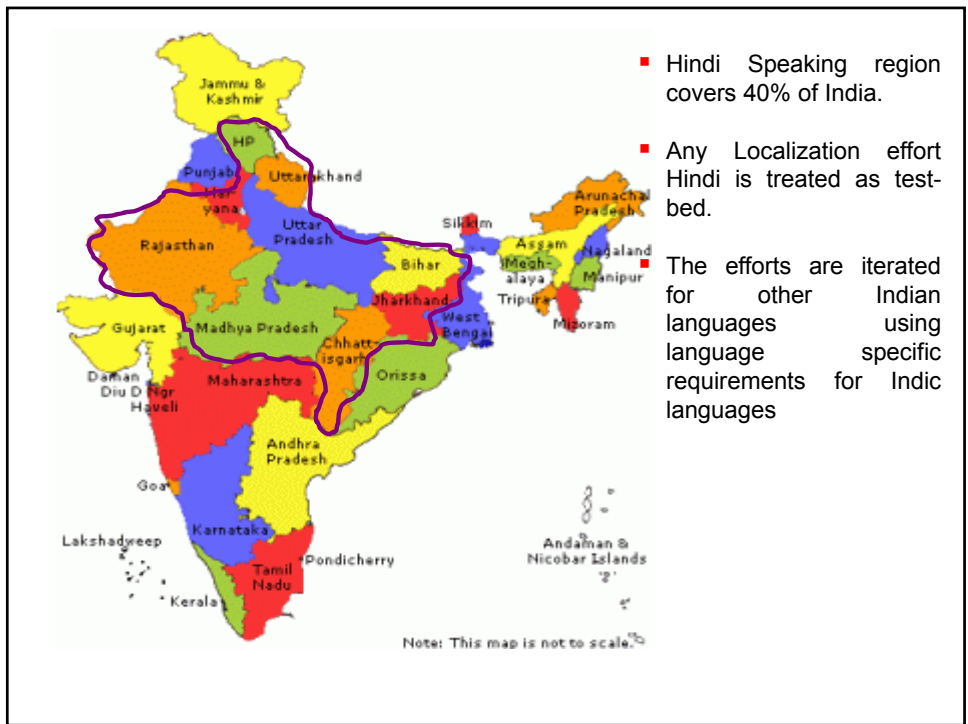
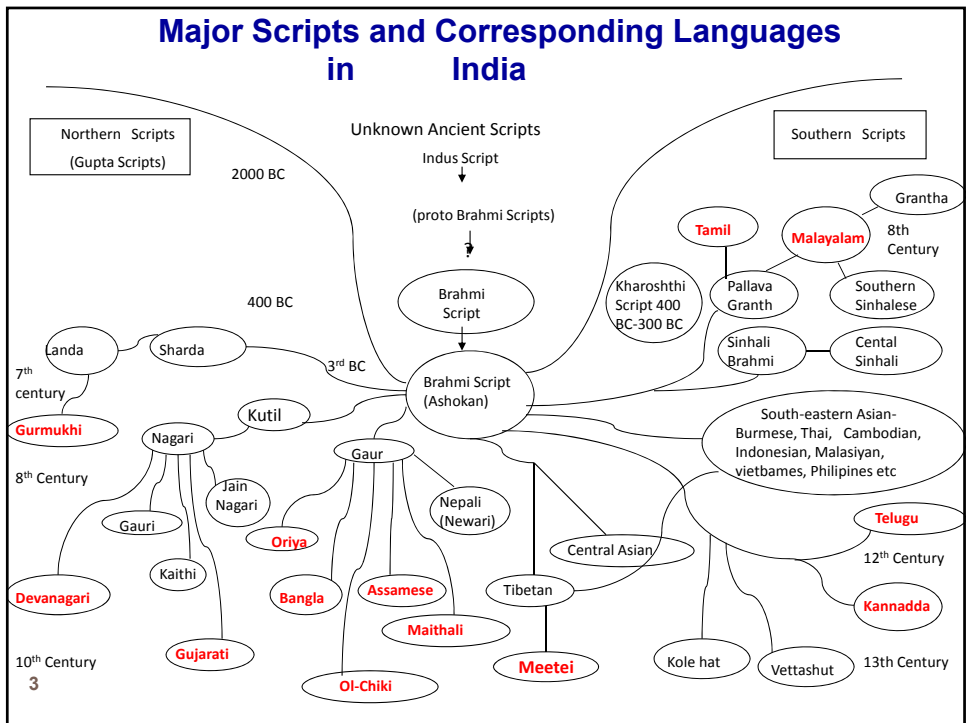
Presented by :  
Swaran Lata

Senior Director & HoD (TDIL Programme)  
Department of Electronics and Information Technology (DeitY)

E-mail: [slata@deity.gov.in](mailto:slata@deity.gov.in)

## Diverse Multilinguality in India

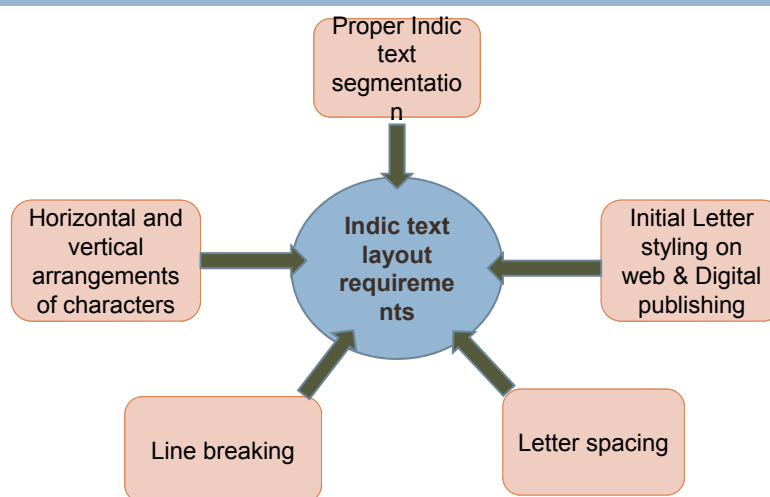
My India is great.	English
<u>मेरा भारत महान।</u>	Hindi
আমার ভারত মহান।	bangla
ಮೇರಾ ಭಾರತ್ ಮಹನ್।	Kannada
మేరా భారత్ మహాన్।	Telugu
<u>എന്റെ ഭാരതം വലിയത് ആണ്</u>	Malayalam
<u>माझा भारत मोठा आहे</u>	Marathi
<u>ମୋର ଭାରତ ବୃହତ୍ ଅଟେ</u>	Oriya
<u>ਮੇਰਾ ਭਾਰਤ ਮਹਾਨ ਹੈ</u>	Punjabi
<u>என்னுடைய இந்தியா மிகச்சிறப்பான இருக்கிறது</u>	Tamil



## Indian language complexities

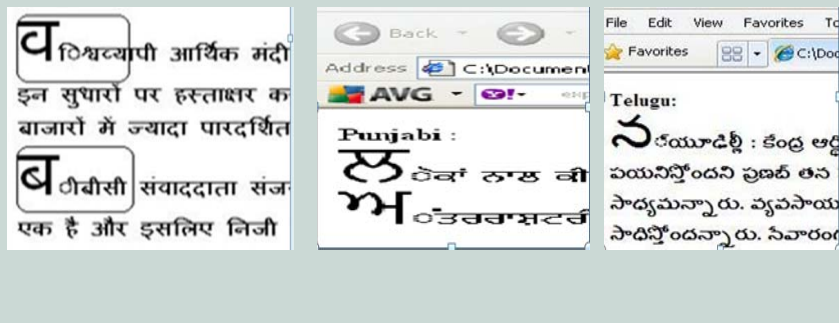
- India has large linguistic diversity with 22 constitutionally recognized languages and 12 scripts
- The mapping between languages and scripts is complex as multiple languages may have common scripts, and a language can be written in multiple scripts
- Each language and script is unique in nature and cannot be easily replicated , even if they share common characteristics

## Indic Text layout requirements



## Challenges in Indian languages

**Use case Scenarios:** Initial letter styling on Web publishing



## Challenges in Indian languages

**Use case Scenarios:** Text input in a word processor

Correct representation

द्वारा विद्यालय

द्वारा विद्यालय

द्वारका विद्वान

द्वारका विद्वान

## Challenges in Indian languages

**Use case Scenarios:** Formatting and spacing on

□ Spacing



बॉलीवुड

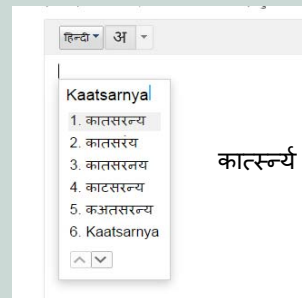
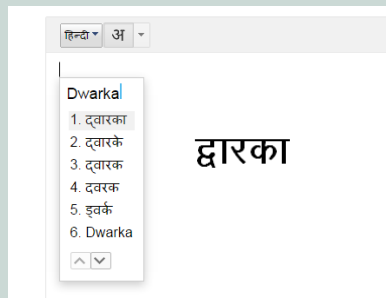
□ Change shape

स्कूल

परिस्थिति

## Challenges in Indian languages

**Use case Scenarios: Phonetic Typing/Transliteration**



## Challenges in Indian languages

### Use case Scenarios : Letter spacing on Web browsers

आर्थिक सुस्ती और इसकी वजह से मांग घटने के बावजूद देश की बड़ी कंपनियां इस साल अपने कर्मचारियों को बोनस दे रही हैं। इन कंपनियों का कहना है कि पिछले साल उन्हें लाभ कमाने की अपनी क्षमता बढ़ाई थी, लिहाजा कर्मचारियों को उसका इनाम दिया जा रहा है।

## Challenges in Indian languages

### Use case Scenarios: Line breaking on applying word wrap

ना 50 से 100 प्रतिशत तक बढ़ाते जा रहे हैं. कपिल गुप्ता कहते हैं, "डिजिटल एडवर्टाइजिंग का सबसे बड़ा आकर्षण यह है कि आप इन विज्ञापनों से इंटरैक्ट कर सकते हैं, यानी प्रतिक्रिया दे सकते हैं. पुराने प्रचार माध्यम - जैसे टी.वी. - के विज्ञापन धीरे-धीरे उबाऊ होते जा रहे हैं." ऐप्स और सोशल मीडिया वेबसाइटों पर बैनर के रूप में विज्ञापन चलाए जाते हैं.

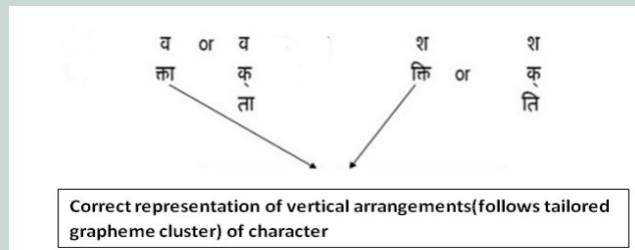
Words not break at Indic syllable base

आकर्षण

विज्ञापन

## Challenges in Indian languages

### □ Vertical arrangements of characters



## Grapheme cluster boundaries defined in UAX#29

### □ **legacy grapheme cluster** :

It is defined as a base followed by zero or more continuing characters.

### □ **Extended grapheme cluster**

It is the same as a legacy grapheme cluster, with the addition of some other characters.

### □ **Tailored Grapheme cluster**

Tailoring of Grapheme cluster to meet further requirements

## Approach to be taken for Possible Solution

- Due to high complexities of Indian languages , it is required to tailored the grapheme cluster for Indian languages
- Indian languages Orthographic syllable should be based on tailored Grapheme Cluster as defined in UAX#29
- Rules for wrapping of Indian languages characters and identification of syllable boundaries needs to be evolved for tailoring of grapheme cluster so that segmentation in Indian languages seems logically.

## Indic Orthographic syllable

- An Orthographic syllable includes Independent vowel or a base consonant and/or any combination of the following characters in the text stream:
  - Consonant/s and consonant + virama sequences
  - vowel signs
  - Modifiers

The above definition of Orthographic syllable is based on the tailored grapheme cluster discussed in section 3 of UAX#29 report.



## Sample tailored Grapheme Cluster Boundaries for Indian languages

- Examples of Indic Orthographic syllable based on tailored grapheme cluster boundaries

क्या	0915 (क)DEVANAGARI LETTER KA	Devanagari kya	स्तः	0938 (स) DEVANAGARI LETTER SA	Devana gari sth
	094D (ः) DEVANAGARI SIGN VIRAMA			0924 (त) DEVANAGARI LETTER TA	
	092F (य)DEVANAGARI LETTER SSA			0903 (ः) DEVANAGARI Sign Visarga	
स्थि	093E (ा)DEVANAGARI SIGN AA		क्ल	0924 (त) DEVANAGARI LETTER TA	Devana gari tkl
	0938 (स)DEVANAGARI LETTER SA	Devanagari sthi		094D (ः) DEVANAGARI SIGN VIRAMA	
	094D (ः)DEVANAGARI SIGN VIRAMA			0915 (क) DEVANAGARI LETTER KA	
	0925 (थ)DEVANAGARI LETTER THA			094D (ः) DEVANAGARI SIGN VIRAMA	
	091C (ि)DEVANAGARI LETTER I			0932 (ल) DEVANAGARI LETTER LA	

## Improving Indic text segmentation....

### Formulation of ABNF based Indic Orthographic syllable definition for defining rules

- ABNF Valid Segmentation based Indic orthographic syllable definition is provided for correct and standardized representation of Indian languages text segmentation
- Augmented Backus–Naur Form (ABNF) is a meta-language based on Backus–Naur Form (BNF), but consisting of its own syntax and derivation rules. The motive principle for ABNF is to describe a formal system of a language to be used as a bidirectional communications protocol.

## Indic Orthographic syllable definition

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

- The linguistic definition of Indic orthographic syllable has been mapped to ABNF(Augmented Backus–Naur Form) for the purpose of text segmentation, line breaking , drop letter, letter spacing in horizontal text and vertical text representation.

## Indic Orthographic syllable definition

**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 [**V**<sub>vs</sub> has been used as symbol in Unicode for dependent vowel of full vowel  
V e.g **AA**<sub>vs</sub>]
- H is Virama/ halant
- | is a rule separator
- [ ] - The enclosed items is optional under this bracket
- { } - The enclosed item/items occurs zero or repeated multiple times

## Indic syllable boundary determination

### No break rules for Indian languages

Rules	Do not break between
V[m]	Independent vowel and Modifier
{CH}C[v][m]	one or more consonant(N) + virama sequences and Consonant  zero or more consonant(N) + virama sequences , Consonant and dependent vowel sign  zero or more consonant(N) + virama sequences , Consonant and modifier  zero or more consonant(N) + virama sequences, Consonant ,dependent vowel sign and modifier
CH	Consonant(N) with virama (applicable only for those Indian languages where pure consonant appears at the end of the word)

Note : Consonant may or may not include Nukta(N)

## Categories values of Indic Orthographic syllable

The precise list of characters with their Unicode code points of all the categories i.e C, H, V etc defined in Indic syllable definition are enclosed as appendix 1 on the following link :

<http://www.unicode.org/L2/L2016/16161-indic-text-seg.pdf>

## Boundary determination for line breaking

- In Indic writing system , it is preferred that line breaks at word boundaries ,if required following principle may be adhered : New line cannot begin with following symbols/Punctuation marks. Also these should be retain with the associated text :

Symbol	Character Name	Unicode code-point
	DEVANAGARI DANDA	U + 0964
	DEVANAGARI DOUBLE DANDA	U + 0965
)	RIGHT PARENTHESIS	U + 0029
+	PLUS SIGN	U + 002B
*	ASTERISK	U + 002A
-	HYPHENATIONPOINT-VISIBLE HYPHEN	U + 2027
	HYPHENATION-SOFT HYPHEN	U+ 00AD
/	SOLIDUS	U + 002F
,	COMMA	U + 002C
.	FULL STOP	U + 002E
:	COLON	U + 003A
;	SEMICOLON	U + 003B
=	EQUALS SIGN	U + 003D
>	GREATER-THAN SIGN	U + 003E
]	RIGHT SQUARE BRACKET	U + 005D
_	LOW LINE	U + 005F
	VERTICAL LINE	U + 007C
}	RIGHT CURLY BRACKET	U + 007D
~	TILDE	U + 007E
%	PERCENT SIGN	U + 0025

## Hyphenation at line boundary

- The definition of Indic orthographic syllable may be used to break the line and a hyphen should be at the breaking point so that word can be read intuitively.
- However the language specific morpho-phonemic rules and industry practices (from media, publishing and grammar books) could be used for hyphenation. U+ 00AD (soft hyphen) is used in some languages such as Tamil and Malayalam.
- The hyphenated words can be broken at the hyphenation point (U + 2027) e.g.:  
नर-नारी should be treated as:  
नर- on the first line and नारी on the next line

# Hyphenation used in printed documents

## Hindi

अन्य भाषाओं से भी लिये जाते हैं। वे सब बने-बनाए लिए जाते हैं। सन्धि-समास विधेयी शब्द भी लिये जाते हैं। उनका भी विकल्पण अनावश्यक है। 'योगाक्षर' 'हासिक' दोनों संस्कृत से लिये गए हैं पर 'हृद्य' 'दृक्' और 'आदि' 'वृद्धि' के नियम को जानकर नहीं। 'उपदेश' में 'उप' उपसर्ग है यह जानना केवल संस्कृत के लिए उपयोगी है, हिन्दी-मुजरती आदि के छात्र के लिए नहीं। फिर भी सभी भारतीय भाषाओं के व्याकरणों में संस्कृत का व्याकरण मरा होता है। 'प्रयोग और प्रयोग' में दिया हुआ 'उपसर्ग' का विवरण देखें :

21.1.0. 'उपसर्ग'—हिन्दी में संस्कृत और उर्दू से आए पूर्व प्रत्ययों के अलावा हिन्दी के अपने उपसर्ग भी हैं। 'फिर संस्कृत के सभी उपसर्गों के उदाहरण दिए हैं। उर्दू के उपसर्गों के उदाहरणरूप 'कमउन्न, गैरकानूनी, दरअसल, नापसन्द, बदनाम' आदि गिनाए गए हैं। बास्तव में हिन्दी में ये शब्द सोपसर्ग अपनाए गए हैं। पर उपसर्गों को हिन्दी के उपसर्ग मानकर नहीं। यदि माना होता तो उपपदाना, अनुसमझना, गैरआवश्यक जैसे प्रयोग भी बराबर चलते। हिन्दी के उपसर्गों के भी उदाहरण प्रायः ठीक नहीं हैं। उनसठ, औपुन, कपूल, सपूल आदि में 'उन, औ, क, स, उपसर्ग नहीं हैं। पूरे शब्द तद्भव रूप में गृहीत हुए हैं। तत्सम रूप ये ऊनघट्टि, अवगुण, सुपुत्र, सुपुत्र। 'उनीदा' के 'उ' को भी उपसर्ग मानने का सुझाव दिया गया है जो ठीक नहीं। यह भी 'उनिद्र' का 'तद्-भव रूप है। सार यह है कि हिन्दी आदि आधुनिक भाषाओं में गृहीत तत्सम, तद्भव और विदेशी शब्दों में सन्धि-समास, उपसर्ग-परसर्ग का विचार अनावश्यक

## Punjabi

ਵਿਗਿਆਨ ਤੇ ਹੋਰ ਸਾਸਤ੍ਰਾਂ ਆਦਿ ਦਾ ਹੋਣਾ ਅੰਬਵ ਜਿਹਾ ਪ੍ਰਬੰਧ ਹੋਣਾ ਹੈ, ਪਰ ਵਾਸਤਵ ਵਿਚ ਇਹ ਗੱਲ ਨਹੀਂ ਹੈ। ਲਿਪੀ ਦੇ ਅਭਾਵ ਵਿੱਚ ਵੀ ਸਾਹਿਤ ਤੇ ਇਤਿਹਾਸ ਆਦਿ ਹੋ ਸਕਦੇ ਹਨ ਅਤੇ ਪਹਿਲਾਂ ਵੀ ਸਨ, ਕੋਵਲ ਇਹ ਅੰਤਰ ਹੋ ਜਾਂਦਾ ਹੈ ਕਿ ਬਾਣੀ ਇਹ ਸਾਹਿਤ ਅਨਿਸ਼ਚਤ ਜਿਹਾ ਰਹਿੰਦਾ ਹੈ—ਧਰਮ, ਜੰਤੂ-ਮੰਤ੍ਰ ਦਾ ਸਾਹਿਤ, ਵਿਦਾ-ਬਧ ਸਾਹਿਤ, ਇਤਿਹਾਸ ਲੋਕ-ਕਥਾਵਾਂ ਦਾ ਰੂਪ ਧਾਰਨ ਕਰ ਲੈਂਦਾ ਹੈ। ਸਾਡੇ ਇਕ ਮੰਤ੍ਰ, ਰਾਮਾਇਣ ਤੇ ਮਹਾਂਭਾਰਤ ਆਦਿ ਦੀਆਂ ਕਥਾਵਾਂ, ਯੂਨਾਨੀ ਸਾਹਿਤ ਦੀਆਂ ਡੀਸੀ ਕਹਾਣੀਆਂ ਤੇ ਵੱਖ ਵੱਖ ਈਸਾਈ ਪਰੰਪਰਾ-ਗਣ ਕਥਾਵਾਂ ਇਸ ਦੇ ਉਦਾਹਰਣ ਰੂਪ ਸਬਤ ਹਨ। ਅਤੇ ਲੇਖਨ-ਕਲਾ ਦੀ ਅਣਹੋਂਦ ਵਿਚ ਧਰਮ, ਸਾਹਿਤ ਤੇ ਇਤਿਹਾਸ ਦਿ ਦਾ ਹੋਣਾ ਸੰਭਵ ਹੈ। ਜਿਵੇਂ ਗੱਲ ਇਹ ਹੈ ਕਿ ਲਿਪੀ ਤੋਂ ਭਾਵ ਕੋਵਲ ਵਰਣ ਪੀ ਤੋਂ ਹੀ ਨਹੀਂ ਹੈ। ਜਿਵੇਂ ਲੇਖਨ-ਕਲਾ ਦੇ ਅਭਾਵ ਵਿਚ ਸਾਹਿਤ ਦਾ ਹੋਣਾ ਸੰਭਵ ਉੱਚ ਹੀ ਵਰਣਮਾਲਾ ਦੇ ਅਭਾਵ ਵਿਚ ਲਿਪੀ ਦਾ ਹੋਣਾ ਵੀ ਸੰਭਵ ਹੈ। ਵਰਣਮਾਲਾ ਅਭਾਵ ਵਿਚ ਮਨੁੱਖ ਰਸੂ (ਗੋਵ), ਰੇਖਾ (ਲਕੀਰ) ਤੇ ਹਿੰਦੂ ਆਦਿ ਰਾਹੀਂ ਅਪ-ਮਨੁੱਖਾਂ, ਮਨੁੱਖਾਂ, ਸੱਧਰਾਂ, ਗੋਬਾਂ, ਵਲੰਕਿਅਾਂ ਅਤੇ ਜੰਗਲੀ ਸ਼ੁਕੁਪੁਰ ਮਾਮੂਲੀ ਤੇ ਗਾਂ ਨੂੰ ਲਿਪੀ-ਬਧ ਕਰਦਾ ਸੀ। ਇਸ ਲਈ ਲਿਪੀ ਦੇ ਅੰਤਰਗਤ ਵਰਣ-ਲਿਪੀ ਬਲਾਦਾ, ਗੋਵ-ਲਿਪੀ, ਲਕੀਰ-ਲਿਪੀ ਤਸਵੀਰ-ਲਿਪੀ ਆਦਿ ਵੀ ਆ ਜਾਂਦੀਆਂ ਹਨ। ਲਿਪੀ

# Word-break at line boundary in south Indian language

## Malayalam

എടുത്തതിനു ശേഷം ഉറപ്പിച്ചു കുടി വയ്ക്കില്ല എന്നതായിരുന്നു കാടണം.

**കാർക്കിളിൻ അവാടും ഭാണക്കളികളും**  
 പുക്കളുവു, ഭാണക്കളികളുമായി വലിയൊരു അടുത്താണ് ഒരതിൽ തറവാട്ടിൽ ധാരാളം പണിക്കാരും വാങ്ങിക്കൊടുക്കേണ്ടുന്ന ഒരു വലിയ യോഗം. സാഹിത്യവും കവിയായകളും അവിടെ ചർച്ചകളിൽ സജീവമാണ്. ഭാണക്കാരന്മാരുടെയും നിരവധി കളികളും തറവാട്ടു മുറ്റത്ത് അടുത്തും.

ഭാണക്കളികളുടെ എതിർപ്പും പൊലെ സമതല നിറഞ്ഞു നിൽക്കുന്നത് ഇത്തരം കളികളിലാണ് പലപ്പോഴും. തോന്നിയിട്ടുണ്ട്. ഇപ്പോൾ ഇത്തരം ഭാണക്കളികൾ എങ്ങനെയോ അപമാനം കഴിഞ്ഞിരിക്കുന്നു. ഭാണക്കളികളുടെ എതിർപ്പും അതിൽ എതിർപ്പും വലിയ വികാസമില്ല. മറ്റൊരു വിഷയത്തിൽ വാടനാവതാം. ജനപ്രിയനായ മഹാബലി രാജാവിനെ ചുമട്ടിത്താഴ്ത്തി എന്ന് വിശ്വസിക്കാൻ കഴിയില്ല. ഭാണക്കളികൾ വീണ്ടും കൊണ്ടു വരാനു വേണ്ടി സൂക്ഷിച്ച് താക്കും ഈ എതിർപ്പും.

ഭാണക്കളികളും പാട്ടുകളും നടുക്ക് അന്യമായിരിക്കുകയാണ്. എല്ലാവരും സൂര്യയിൽ ഇരുത്താൻ ശ്രമിക്കുന്നത്. ഭാണക്കളികളുടെ എതിർപ്പും വേഗം തടഞ്ഞുവെക്കൽ ഈ കളികളും പാട്ടുകളും കുടി മടങ്ങിവരണം. ഭാണക്കളികളിലാണ് പുർത്തമായ സമതലം രൂപമാകുന്നത്. ജാതിമതഭേദമന്യം എല്ലാവരും ഈ കളികളിൽ പങ്കെടുക്കേണ്ടവരാണ് ഭാണം. പുർത്തമാകുന്നത്. പാവപ്പെട്ടപ്പോൾ എന്റെ നടുക്കളിയിലെ ഭാണക്കാരന്മാരുടെ ഇത്തരത്തിലുള്ള സന്ദേശം നൽകാനാണ് ശ്രമിക്കുന്നത്. മൂന്നു വർഷമായി എന്റെ ശിഷ്യർ ഈ ആഘോഷം സംഘടിപ്പിക്കുന്നുണ്ട്. അതിൽ മരണശോചനയും പങ്കെടുക്കുന്നവർ അറിയാതെ ഞാൻ മേതിൽ തറവാട്ടിലെ മുറ്റത്ത് കളിച്ച് ചിരിച്ചു നടക്കുന്ന ആ കൊച്ചു കുട്ടിയായി മാറും.

# Indic text segmentation results based on Indic syllable definition

0938: स DEVANAGARI LETTER SA  
 094D: ् DEVANAGARI SIGN VIRAMA  
 0925: थ DEVANAGARI LETTER THA  
 093F: ि DEVANAGARI VOWEL SIGN I

→ syllable

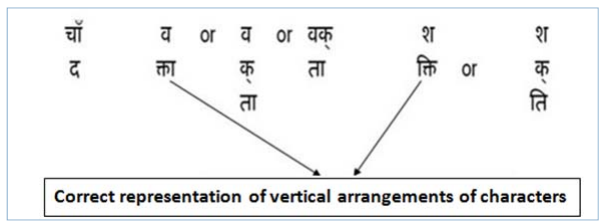
কী লামবাও চেঙ্গকলকলো  
 ঙ্গা চরোই-তেতু খোংলকলো যুগ-  
 যুগকী মীংয়েংলা অভিসা  
 লস্বী কুপখোো নোংগী নাউপোমদা সনাগী  
 লিকলাং ভারকলো।

শ হ दि मू क र

क्रि केट के लिए मशहूर व  
 इसके लिए दीवानगी  
 दिखाने वाले खेल प्रेमियों के देश  
 भारत में अगर कोई अन्य खिलाड़ी  
 प्रशंसा हासिल करता है तो यह  
 किसी उपलब्धि से कम नहीं है।  
 फिर हम जिस खिलाड़ी की बात  
 कर रहे हैं उसने न केवल प्रशंसा

ഈ ലേഖകനെ വളരെയൊരായ  
 കർഷിച്ഛുള്ള വിഭാഗകാ  
 വ്യയാണു കണ്ണൂരിത്തുള്ളി. നാല്പാട്ട്  
 നാരായണമേനോന്റെ ഈ വിഭാഗകാ  
 വ്യവ്യാ വിക്റ്റർ യുഗോവീന്റെ നോവൽ

# Indic text segmentation results based on Indic syllable definition



क  
ल  
क  
त

धा रा वा हि क उ प ना स  
 पा ल्ठा हा उ या

## Proposal to incorporate Indian languages requirements in UAX#29

It is proposed to incorporate following Indian languages text segmentation requirements in UAX#29

- Additional information on Indic orthographic syllable boundaries based on tailored grapheme cluster define in UAX#29
- ABNF valid segmentation definition to define Indian languages orthographic syllable
- No break rules for determination of Indic syllable boundary
- Information for identification of boundaries of first letter styling, Guiding principles of line breaking at syllable level for Indian languages.
- Detailed report at [L2/16-161](#)

Thanks