

# Unicode Normalization and Grapheme Parsing of Indic Languages

Nazmuddoha Ansary\*, Quazi Adibur Rahman Adib\*, Tahsin Reasat  
Asif Shahriyar Sushmit, Ahmed Imtiaz Humayun, Sazia Mehnaz Kanij  
Fatema, Mohammad Mamun Or Rashid, Farig Sadeque

Presented By:  
Quazi Adibur Rahman Adib

\* These authors share first authorship

# Introduction

- **Motivation**
  - Fixing Malformed Words
  - Encapsulating feature without disregarding writing forms of Indic Languages
- **Core Objectives**
  - Developing a Unicode normalizer for Indic languages to fix malformed words.
  - Developing a Grapheme Parser

# Unicode Normalizer

- Broken Diacritics (BD)
- Broken Nukta (BN)
- Invalid Unicode (IU)
- Invalid Connector (IC)
- Fix Diacritics (FD)
- Vowel-Vowel Diacritic (VDV)
- Legacy Symbols Handling (LS)
- Language Specific Treatment

# Broken Diacritics (BD)

## Incorrect Form

संस्कृति

[स, ं, स, ्र, क, ्र, त, ि]

[U+09B8, U+0982, U+09B8, U+09CD,  
U+0995, **U+09C4**, U+09A4, U+09BF]

## Correct Form

संस्कृति [ʃɒnskɾiti]

[स, ं, स, ्र, क, ्र, त, ि]

[U+09B8, U+0982, U+09B8, U+09CD,  
U+0995, **U+09C3**, U+09A4, U+09BF]

Potential Source: Cognitive

# Broken Nukta (BN)

## Incorrect Form

কেন্দ্রীয়

[ক, ে, ন, ্, দ, ্র, ী, য, ়]

[U+0995, U+09C7, U+09A8, U+09CD,  
U+09A6, U+09CD, U+09B0, U+09C0, **U+09AF**,  
**U+09BC**]

## Correct Form

কেন্দ্রীয় [kendriio]

[ক, ে, ন, ্, দ, ্র, ী, য]

[U+0995, U+09C7, U+09A8, U+09CD,  
U+09A6, U+09CD, U+09B0, U+09C0,  
**U+09DF**]

Potential Source: Keyboard

# Invalid Unicode (IU)

Incorrect Form	Correct Form
<p data-bbox="434 565 629 612">াটোবাকো</p> <p data-bbox="272 663 792 714">[ া, ট, ো, ব, া, ক, ো]</p> <p data-bbox="141 765 923 877">[U+09BE, U+099F, U+09CB, U+09AC, U+09BE, U+0995, U+0995, U+09CB]</p>	<p data-bbox="1232 565 1580 612">টোবাকো [tobeko]</p> <p data-bbox="1182 656 1630 706">[ট, ো, ব, া, ক, ো]</p> <p data-bbox="1006 743 1808 849">[U+099F, U+09CB, U+09AC, U+09BE, U+0995, U+0995, U+09CB]</p>

Potential Source: Typing Mistake

# Invalid Connector (IC)

Incorrect Form	Correct Form
<p data-bbox="484 561 581 620">দুইটি</p> <p data-bbox="345 659 716 718">[দ, ু, ই, ্, ট, ি]</p> <p data-bbox="137 762 925 877">[ U+09A6, U+09C1, U+0987, <b>U+09CD</b>, U+099F, U+09BF]</p>	<p data-bbox="1286 561 1522 620">দুইটি [duiti]</p> <p data-bbox="1257 648 1557 707">[দ, ু, ই, ট, ি]</p> <p data-bbox="1020 746 1792 860">[ U+09A6, U+09C1, U+0987, U+099F, U+09BF]</p>

Potential Source: Cognitive

# Fix Diacritics (FD)

Incorrect Form	Correct Form
<p data-bbox="504 612 562 672">दूई</p> <p data-bbox="417 713 649 773">[द, ू, ू, ई]</p> <p data-bbox="137 812 925 863">[ U+09A6, U+09C1, <b>U+09C1</b>, U+0987]</p>	<p data-bbox="1325 612 1489 672">दूई [दूई]</p> <p data-bbox="1325 699 1489 760">[द, ू, ई]</p> <p data-bbox="1112 792 1707 844">[ U+09A6, U+09C1, U+0987]</p>

Potential Source: Keyboard



# Vowel-Vowel Diacritic (VDV)

## Incorrect Form

একএ  
[এ, ক, এ, ে]

[U+098F, U+0995, **U+098F**, U+09C7]

## Correct Form

একত্রে [ekot̪re]  
[এ, ক, ত, ্, র, ে]

[U+098F, U+0995, **U+09A4**, **U+09CD**,  
**U+09B0**, U+09C7]

Potential Source: Typing Mistake and Cognitive

# Legacy Symbols Handling (LS)

Incorrect Form

□

[U+0980]

Correct Form

१ [ɟɛ̃]

[U+09ED]

Potential Source: Typing Mistake

# Language Specific Treatment [Assamese Replacement (AR)]

Incorrect Form	Correct Form
<p>ব্যৱহাৰ</p> <p>[ব, ্, য, ব্ৰ, হ, া, ৰ]</p> <p>[U+09AC, U+09CD, U+09AF, <b>U+09F1</b>, U+09B9, U+09BE, <b>U+09F0</b> ]</p>	<p>ব্যবহার [bæbohar]</p> <p>[ব, ্, য, ব, হ, া, ৰ]</p> <p>[U+09AC, U+09CD, U+09AF, <b>U+09AC</b>, U+09B9, U+09BE, <b>U+09B0</b> ]</p>

Potential Source: Cognitive

# Language Specific Treatment [To-hosonto Normalize (THN)]

Incorrect Form	Correct Form
<p data-bbox="479 598 585 653">উত্‌স</p> <p data-bbox="421 699 643 754">[উ, ত, ্, স]</p> <p data-bbox="141 798 923 852">[U+0989, <b>U+09A4</b>, <b>U+09CD</b>, U+09B8]</p>	<p data-bbox="1300 603 1512 658">উৎস [ut̪ʃo]</p> <p data-bbox="1329 696 1483 751">[উ, ৎ, স]</p> <p data-bbox="1116 789 1696 844">[U+0989, <b>U+09CE</b>, U+09B8]</p>

Potential Source: Cognitive

# Language Specific Treatment [Complex Root Normalization (CRN)]

Incorrect Form	Correct Form
<p>बिष्पद [ब, ि, ्ष, ्, प, ्, द]</p> <p>[U+09AC, U+09BF, <b>U+09CD</b>, U+09B7, U+09CD, U+09AA, <b>U+09CD</b>, U+09A6]</p>	<p>बिष्पद [biʃpod̪] [ब, ि, ष, ्, प, द]</p> <p>[U+09AC, U+09BF, U+09B7, U+09CD, U+09AA, U+09A6]</p>

Potential Source: Cognitive

Language	Total	Affected	%
Bangla	2,883,731	369,348	12.81
Devanagari	2,887,725	257,615	8.92
Gujarati	1,119,927	28,512	2.55
Odiya	414,483	36,267	8.69
Tamil	6,885,008	214,242	3.11
Punjabi	421,537	132,993	<b>31.55</b>
Malayalam	6,021,715	932,314	15.48

Table 1: Summary of the OSCAR Abugida language corpus. It shows how many, out of the total unique words, our pipeline was able to capture as unnormalized words

## Word Level Experiments

# NLP Experiments

[Normalization Performance]

Train Norm.	Test Norm.	F1 (NER)	F1 (Sentiment Analysis)
None	IndicNLP	$0.89 \pm 0.012$	$0.79 \pm 0.014$
	Ours	$0.89 \pm 0.011$	$0.79 \pm 0.013$
IndicNLP	None	$0.90 \pm 0.011$	$0.78 \pm 0.002$
	IndicNLP	$0.91 \pm 0.011$	$0.78 \pm 0.001$
Ours	None	$0.90 \pm 0.013$	$0.78 \pm 0.008$
	Ours	$0.90 \pm 0.007$	$0.78 \pm 0.008$

Table 2: Normalization Performance of IndicNLP (I) and our proposed (O) Normalizer.

# NLP Experiments

## [Robustness]

Train Norm.	Test Norm.	F1 (NER)	F1 (Sentiment Analysis)
None	att-1°	0.79 ± 0.039	0.77 ± 0.013
	att-5°	0.73 ± 0.011	0.74 ± 0.006
IndicNLP	att-1°	0.81 ± 0.003	0.76 ± 0.012
	att-5°	0.74 ± 0.013	0.73 ± 0.004
	att-1° + I	0.84 ± 0.019	0.76 ± 0.012
	att-5° + I	0.75 ± 0.015	0.73 ± 0.004
Ours	att-1°	0.80 ± 0.013	0.77 ± 0.004
	att-5°	0.73 ± 0.021	0.73 ± 0.013
	att-1° + O	0.89 ± 0.002	0.78 ± 0.008
	att-5° + O	0.89 ± 0.007	0.76 ± 0.007

Table 3: Robustness of IndicNLP (I) and our proposed (O) Normalizer. Here, att- $\{1, 5\}^\circ$  means attack degree.



# Grapheme Parser

$t = \text{"উজ্জ্বল নক্ষত্র অন্তোনীয় গ্রামসি"}$

$w = [\text{'উ', 'জ', '্', 'জ', '্', 'ব', 'ল', ' ', 'ন', 'ক', '্', 'ষ', 'ত', '্', 'র', ' ', 'অ', 'ন', '্', 'ত', 'ো', 'ন', 'ী', 'য়', ' ', 'গ', '্', 'র', 'া', 'ম', 'স', 'ি'}]$   $n=32$

$[0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31]$

$C = [2, 4, 10, 13, 18, 26]$

Index:  $i, 0 < i < n-1$

Values:  $C[i=0]=2, C[i=1]=4$

$D = [(1, 2, 3), (3, 4, 5), (9, 10, 11), (12, 13, 14), (17, 18, 19), (25, 26, 27)]$   $k = 6$

Index:  $j, 0 \leq j \leq k-1$

Values:

$D[j=0] = (1, 2, 3) \rightarrow (C[i=0]-1, C[i=0], C[i=0]+1)$

$D[j=1] = (3, 4, 5) \rightarrow (C[i=1]-1, C[i=1], C[i=1]+1)$

Merging recursively

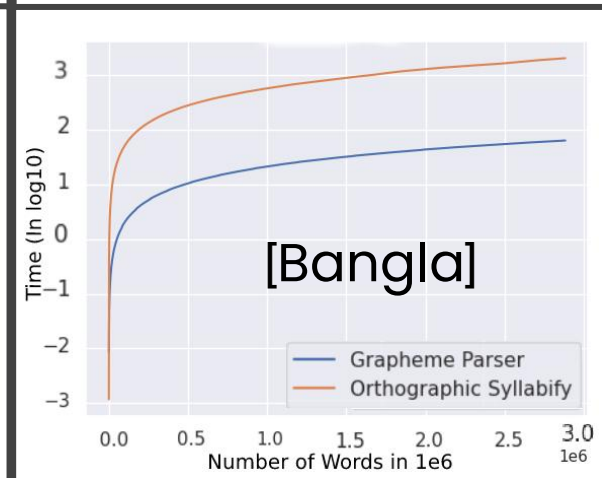
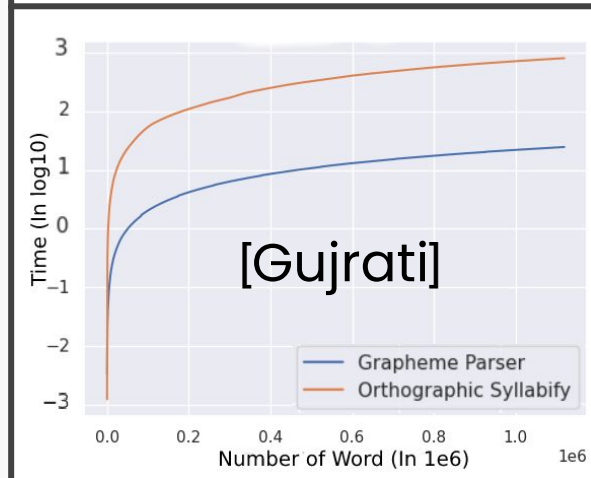
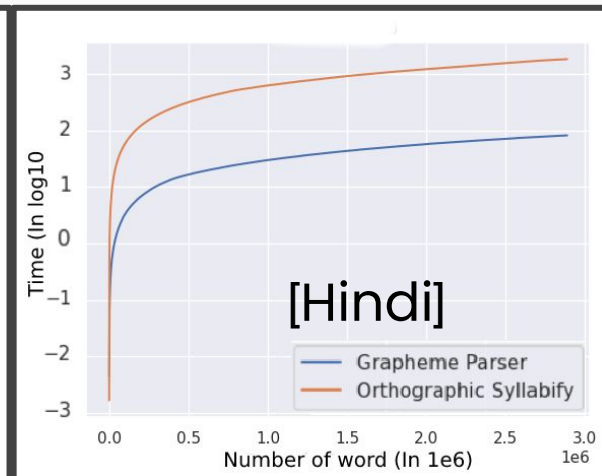
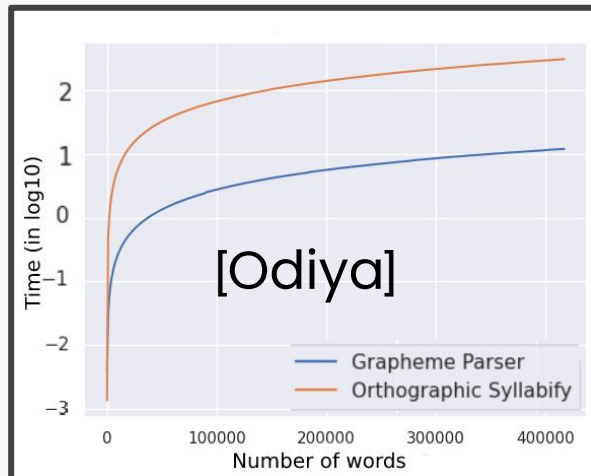
$D' = [(1, 2, 3, 4, 5), (9, 10, 11), (12, 13, 14), (17, 18, 19), (25, 26, 27)]$   $p = 5$

$w' = [\text{'উ', 'জ্জ', 'ল', ' ', 'ন', 'ক্ষ', 'ত্র', ' ', 'অ', 'ন্ত', 'ো', 'ন', 'ী', 'য়', ' ', 'গ্র', 'া', 'ম', 'স', 'ি'}]$

Merging gc, vd and cd

$G = [\text{'উ', 'জ্জ', 'ল', ' ', 'ন', 'ক্ষ', 'ত্র', ' ', 'অ', 'ন্তো', 'নী', 'য়', ' ', 'গ্রা', 'ম', 'সি'}]$

Time taken for our parser to process examples vs. time taken by IndicNLP syllabifier. Time is shown as a logarithmic function, and the example count is in scale of  $10^6$ .



Thank You

# Grapheme Parser [Algorithm]

1. For a given text  $t$ , get the list of Unicode characters:  
 $w = u_0, u_1, u_2, \dots, u_{(n-1)}$
2. Create a list  $C$  with  $k$  elements having positions of connectors present in  $w$  in ascending order.  
Here  $k$  is the total number of connectors in  $w$ .  
For every element  $i$  in  $C$ :  $0 < i < n-1$  and  $u_{(i-1)}, u_{(i+1)}$  are Consonants.
3. Create  $k$  lists  $D$ , maintaining the ascending order of  $C$ .  
Each list has 3 consecutive numbers.  
A list  $d$  that is an element of  $D$  is formed as:  $(i_{j-1}, i_j, i_{j+1})$  where  $0 <= j <= k-1$ .
4. Two lists are mergeable if and only if  $i_{j+1} = i_{(j+1)-1}$ .  
The new merged list  $d$  in  $D$  will have 5 elements:  $(i_{j-1}, i_j, i_{j+1}, i_{(j+1)}, i_{(j+1)+1})$ .  
Merge lists in  $D$  recursively until lists are not mergeable anymore, maintaining order in  $C$ .  
At the end, we have a new list  $D'$  with  $p$  lists, where each  $d'$  in  $D'$  has varying odd lengths  $>= 3$ .
5. For every  $d'$  with consecutive ascending numbers, merge the Unicode characters in  $w$  at those positions and construct  $w'$ .  
Elements in  $w'$  can only be: pure vowel diacritics  $vd$ , pure consonant diacritics  $cd$ , or grapheme root+conjunct diacritic  $gc$ .
6. To get the list of graphemes  $G$ , add a  $gc$  with the next element in  $w'$  if and only if the next element is  $vd$  or  $cd$ .