



PennState



# Does Transliteration Help Multilingual Language Modeling?

Ibraheem Muhammad Moosa, Mahmud Elahi Akhter, Ashfia Binte Habib

EACL 2023 Paper ID: 433

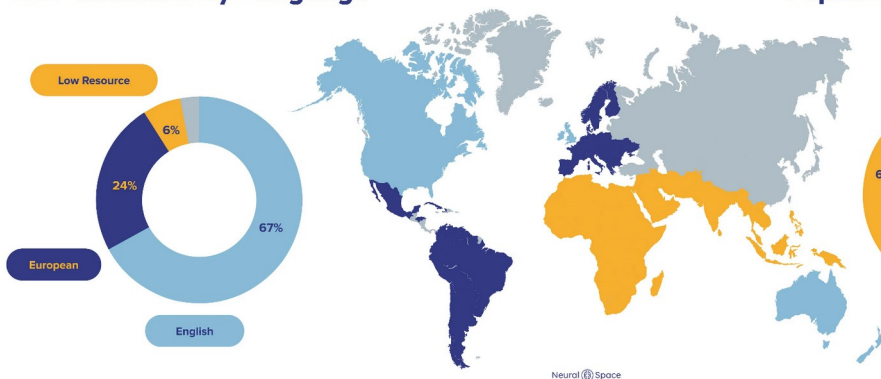


# Motivation & Background

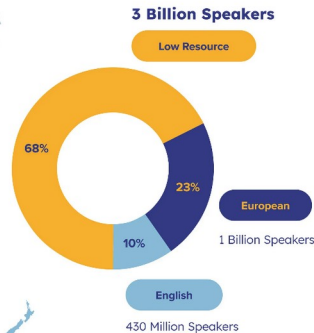
## Open Problems

- Resource disparity between Low Resource Languages and High Resource Languages
- Diversity and Inclusivity of Low Resource Language tasks

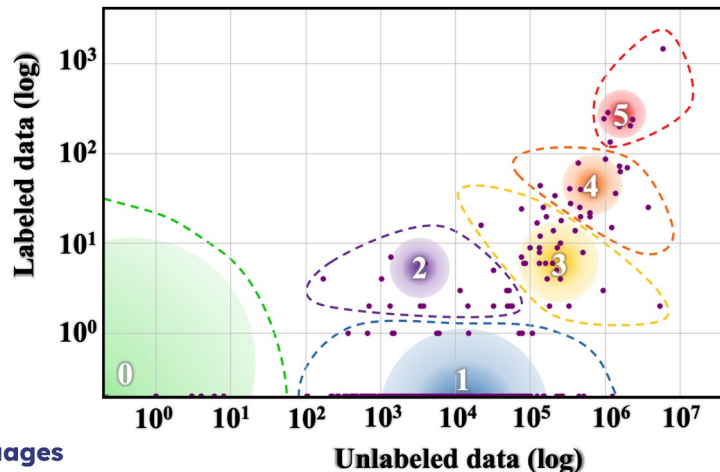
### NLP Solutions by Language



### Population Size of Languages



Diversity and inclusivity disparity<sup>2</sup>



## Issues for Low Resource Languages

- Lack of large pretraining corpus
- Lack of diverse evaluation dataset

1. The State and Fate of Linguistic Diversity and Inclusion in the NLP World, Joshi et al. (2020)  
2. Challenges in using NLP for low-resource languages and how NeuralSpace solves them, Felix Laumann (2022)

# Motivation & Background

## Additional obstacles for Low Resource Language

- ❌ Script barrier leads to poor lexical overlap <sup>1</sup>
- ❌ Poor Tokenization leads to increase in unknown tokens <sup>2</sup>

## Proposed Solution

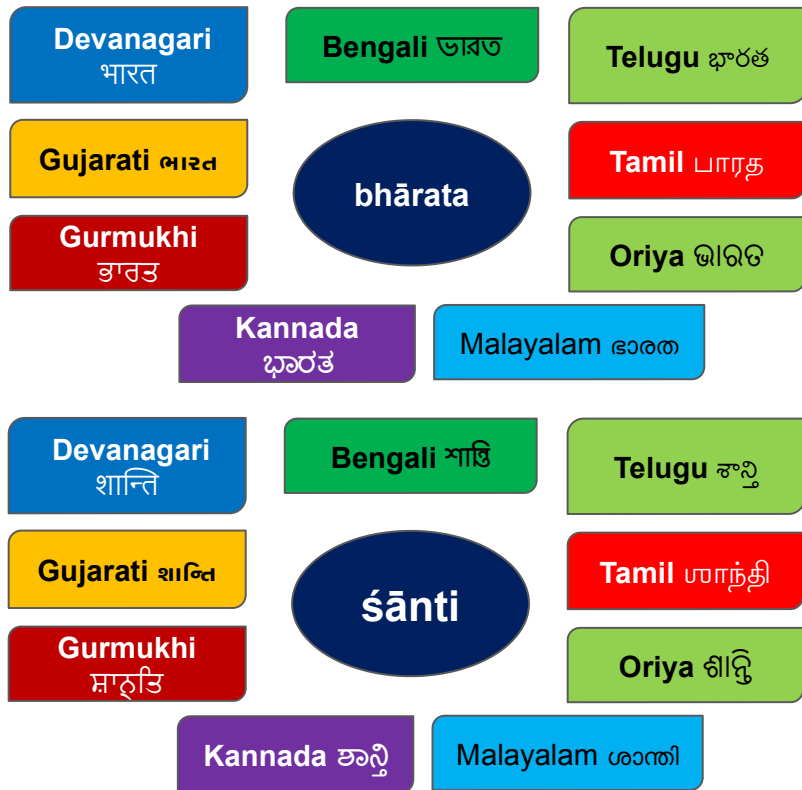
- ✅ Grapheme-to-phoneme (G2P)
- ✅ Transliteration



Script diversity in South Asia<sup>3</sup>.

1. Pushing the Limits of Low-Resource Morphological Inflection, Anastasopoulos and Neubig (2019):
2. UNKs Everywhere: Adapting Multilingual Language Models to New Scripts, Pfeiffer et al.(2021)
3. [en.wikipedia.org/wiki/Languages\\_of\\_South\\_Asia#/media/File:States\\_of\\_South\\_Asia.png](https://en.wikipedia.org/wiki/Languages_of_South_Asia#/media/File:States_of_South_Asia.png)

# Case for Transliteration



## Benefits of Transliteration?

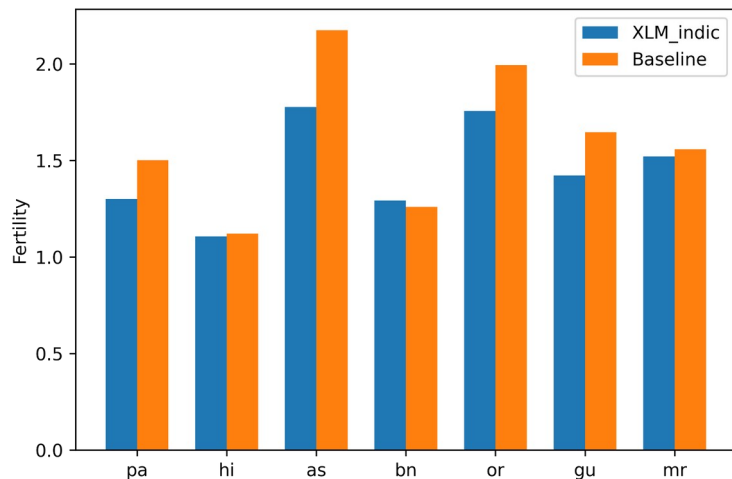
- ✓ Transliteration collapses multiple scripts into a single script
- ✓ Improves lexical overlap
- ✓ Reduces number of token

## However..

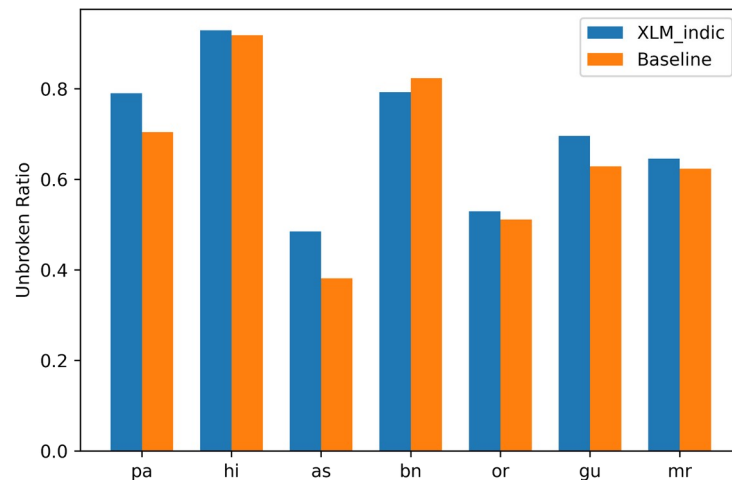
- Does transliterating input scripts improve performance of multilingual language models?
- Given it does, is this improvement statistically significant?
- Does a model trained on transliterated script learn better representation?

# Results: Tokenization Quality

✓ Transliteration Reduces number of token



Average number of tokens per word



Ratio of words unbroken by the tokenizer

# Results: Performance analysis

Model	pa	hi	bn	or	as	gu	mr	kn	te	ml	ta	avg
<b>Wikipedia Section Title Prediction</b>												
RemBERT <sub>MS</sub>	68.42±0.92	70.90±0.39	72.58±0.45	69.92±0.90	68.37±1.37	72.93±0.58	73.23±0.61	71.67±0.41	92.98±0.19	69.03±0.57	69.77±0.45	73.00
RemBERT <sub>US</sub>	71.01±0.22	72.45±0.29	73.65±0.21	75.37±0.69	72.50±0.91	76.35±0.29	74.58±0.72	74.21±0.29	93.66±0.09	69.33±0.35	70.63±0.22	74.89
$\delta$	2.59	1.55	1.07	5.45	4.13	3.42	1.34	2.54	0.68	0.31	0.86	1.89
$p - value$	0.0004	0.0004	0.0004	0.0004	0.0004	0.0004	0.0035	0.0004	0.0004	0.2505	0.0006	-
ALBERT <sub>MS</sub>	74.33±0.83	78.18±0.33	81.18±0.28	74.35±1.2	76.70±0.83	76.37±0.53	79.10±0.84	-	-	-	-	77.17
ALBERT <sub>US</sub>	77.55±0.61	82.24±0.18	84.38±0.29	81.47±0.99	81.74±0.82	82.39±0.27	82.74±0.52	-	-	-	-	81.78
$\delta$	3.22	4.06	3.20	7.12	5.04	6.02	3.64	-	-	-	-	4.61
$p - value$	0.0004	0.0004	0.0004	0.0004	0.0004	0.0004	0.0004	-	-	-	-	-
<b>Named Entity Recognition (F1-Score)</b>												
RemBERT <sub>MS</sub>	69.47±1.72	90.95±0.33	95.51±0.18	87.92±1.26	79±0.22	69±0.94	90.72±0.17	72.65±1.81	81.82±1.81	89.17±0.25	90.07±0.33	83.40
RemBERT <sub>US</sub>	81.91±1.93	91.73±0.39	96.19±0.21	88.92±2.88	83.50±2.75	80.25±1.42	90.75±0.35	78.98±1.50	84.97±0.45	89.26±0.46	90.18±0.27	86.97
$\delta$	12.44	0.78	0.68	1.00	4.28	10.31	0.02	6.33	3.15	0.01	0.12	3.56
$p - value$	0.00004	0.0005	0.00001	0.1615	0.0019	0.00004	0.6665	0.00004	0.00004	0.7304	0.2973	-
ALBERT <sub>MS</sub>	76.69±1.5	91.80±0.42	96.39±0.19	84.18±1.8	75.45±1.8	69.10±2.9	88.72±0.40	-	-	-	-	83.19
ALBERT <sub>US</sub>	85.42±1.9	92.93±0.21	97.31±0.22	93.54±0.58	89.06±2.2	80.16±0.15	90.56±0.44	-	-	-	-	89.85
$\delta$	8.73	1.13	0.92	9.36	13.61	11.06	1.84	-	-	-	-	6.66
$p - value$	0.0004066	0.0004066	0.0003983	0.0004038	0.000401	0.0004066	0.0004095	-	-	-	-	-

Orange indicates the multi-script and uni-script models are equal and blue indicates the uni-script model is better

- ✓ On average transliteration improves performance of multilingual language models.
- ✓ It improves the performance of low resource languages more compared to high resource languages
- ✓ It does not deteriorate the results of low resource languages

# Results: Performance analysis

Language	Dataset	RemBERT <sub>MS</sub>	RemBERT <sub>US</sub>	$\delta$	$p - value$	ALBERT <sub>MS</sub>	ALBERT <sub>US</sub>	$\delta$	$p - value$
<b>Article Genre Classification</b>									
hi	BBC News	76.80±0.84	77.78±0.92	<b>0.98</b>	0.0466	77.28±1.51	79.14±0.60	<b>1.86</b>	0.0088
bn	Soham News Article Classification	92.86±0.10	93.69±0.20	<b>0.83</b>	0.0004	93.22±0.49	93.89±0.48	<b>0.67</b>	0.0090
gu	INLTK Headlines	90.27±0.47	91.60±0.28	<b>1.33</b>	0.0004	90.41±0.69	90.73±0.75	<b>0.32</b>	0.6249
mr	INLTK Headlines	91.24±0.50	92.27±0.39	<b>1.03</b>	0.0008	92.21±0.23	92.04±0.47	<b>-0.17</b>	0.3503
ml	INLTK Headlines	94.11±0.49	93.33±0.22	<b>-0.78</b>	0.003	-	-	-	-
ta	INLTK Headlines	95.59±0.70	94.93±0.30	<b>-0.65</b>	0.013	-	-	-	-
<b>Sentiment Analysis</b>									
hi	IITP Product Reviews	72.17±1.98	72.85±0.63	<b>0.68</b>	0.9646	76.33±0.84	77.18±0.77	<b>0.85</b>	0.04099
hi	IITP Movie Reviews	58.66±1.09	62.65±2.74	<b>3.99</b>	0.0023	65.91±2.2	66.34±0.16	<b>0.15</b>	0.8941
te	ACTSA	61.18±1.38	60.53±0.85	<b>-0.66</b>	0.1981	-	-	-	-
<b>Discourse Mode Classification</b>									
hi	MIDAS Discourse	78.07±0.83	79.46±0.67	<b>1.39</b>	0.0415	78.39±0.33	78.54±0.91	<b>0.15</b>	0.7561

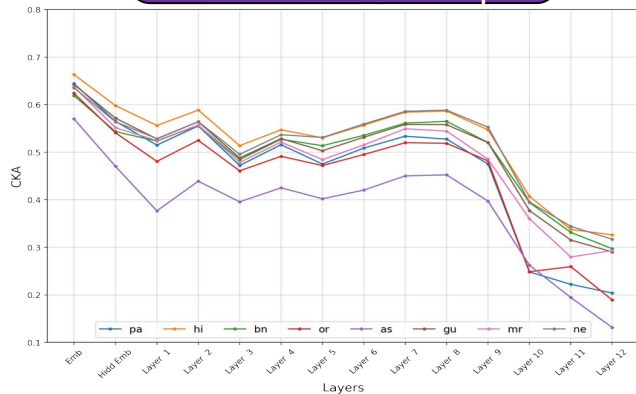
  

Model	pa	hi	bn	or	as	gu	mr	ta	te	ml	kn	avg
<b>Cloze-style QA (Zero Shot)</b>												
RemBERT <sub>MS</sub>	33.93	39.06	38.93	37.32	37.66	84.21	46.15	37.02	34.42	38.45	40.75	42.53
RemBERT <sub>US</sub>	33.92	40.10	39.62	38.28	39.26	85.37	45.92	36.68	34.36	37.16	44.29	43.17
$\delta$	<b>-0.01</b>	<b>1.04</b>	<b>0.69</b>	<b>0.96</b>	<b>1.6</b>	<b>1.16</b>	<b>-0.23</b>	<b>-0.34</b>	<b>-0.06</b>	<b>-1.29</b>	<b>3.54</b>	<b>0.64</b>
ALBERT <sub>MS</sub>	31.04	36.72	35.19	34.63	33.92	59.86	36.14	-	-	-	-	38.21
ALBERT <sub>US</sub>	32.77	38.52	36.38	36.00	37.36	70.22	39.53	-	-	-	-	41.54
$\delta$	<b>1.73</b>	<b>1.8</b>	<b>1.19</b>	<b>1.37</b>	<b>3.44</b>	<b>10.36</b>	<b>3.39</b>	-	-	-	-	<b>3.33</b>

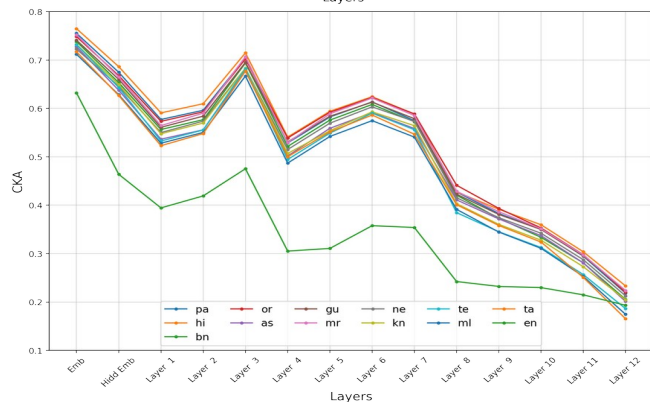
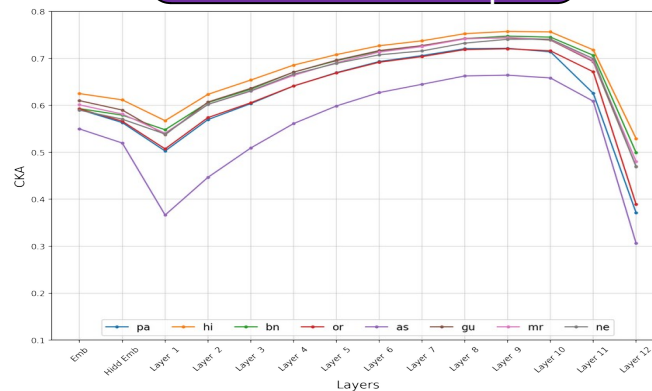
Orange indicates the multi-script and uni-script models are equal, cyan indicates multi-script is better than uni-script and blue indicates vice versa

# Results: Cross lingual representation

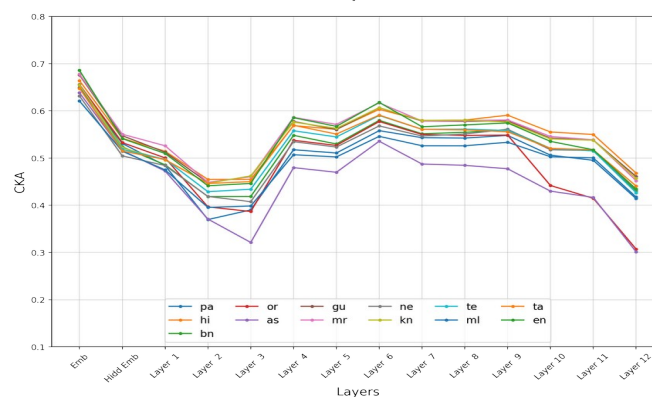
ALBERT Multi-script



ALBERT Uni-script



RemBERT Multi-script



RemBERT Uni-script



PennState



# Does Transliteration Help Multilingual Language Modeling?

Ibraheem Muhammad Moosa, Mahmud Elahi Akhter, Ashfia Binte Habib

EACL 2023 Paper ID: 433

