

# Language Resources and Technology Development Efforts for some Lesser-known Indian Languages

Ritesh Kumar, KMI, Dr. Bhimrao Ambedkar University, India; Atul Kr. Ojha, Panlingua Language Processing LLP, India  
Bornini Lahiri, IIT-Kharagpur, India; Mayank, Kendriya Hindi Sansthan, India; Deepak Alok, Rutgers University, USA  
All present and past **students of KMI, Agra** and **JNU, New Delhi** involved in these projects

## Abstract

For the last few years, we have been involved in the development of language technologies and resources for some of the lesser-known Indian languages viz. Magahi, Bhojpuri, Awadhi and Braj Bhasha. These languages (among others) have been largely marginalised and ignored and have low prestige and negative attitude towards them because of these being considered ‘illiterate’ and ‘rural’ varieties of Hindi. Our poster will showcase different kinds of corpora as well as basic technologies like part-of-speech tagger, morphological analyser as well as some applications like machine translation systems that have been developed for these languages.

पिछला कै साल से, हमनी मगही, भोजपुरी, अवधी आउ ब्रज भासा जईसन भासा, जेकरा पर बहुते कम काम होल हई, ओकर प्रौद्योगिकि आउ संसाधन के विकास में लगल ही। इ सब भासा (औ अइसन बड़ीमनी) के नजर अंदाज और हासिया पर कर देल गेल हे आउ कम परतिस्था आउ नीचा दृस्टि से देखल जा हे काहे की इ सब हिंदी भासा के देहाती आउ असीछित बोली के प्रकार समझल जा हे। हमनी के पोस्टर बिभिन्न प्रकार के कॉर्पोरा के साथ-साथ बुनयादी प्रौद्योगिकि जैसे की सब्द भेद टैगर (पार्ट ऑफ़ स्पीच टैगर), रूप सम्बन्धी बिश्लेसक (मोरफोलॉजिकल अनलाइज़र) आउ कुछ अनुप्रयोग सॉफ्टवेयर जैसे की मसीनी अनुवाद पद्धति/सिस्टम जे इसब भासा ला विकिसित कइल गइल हे ओकरा परदरसीत करत।

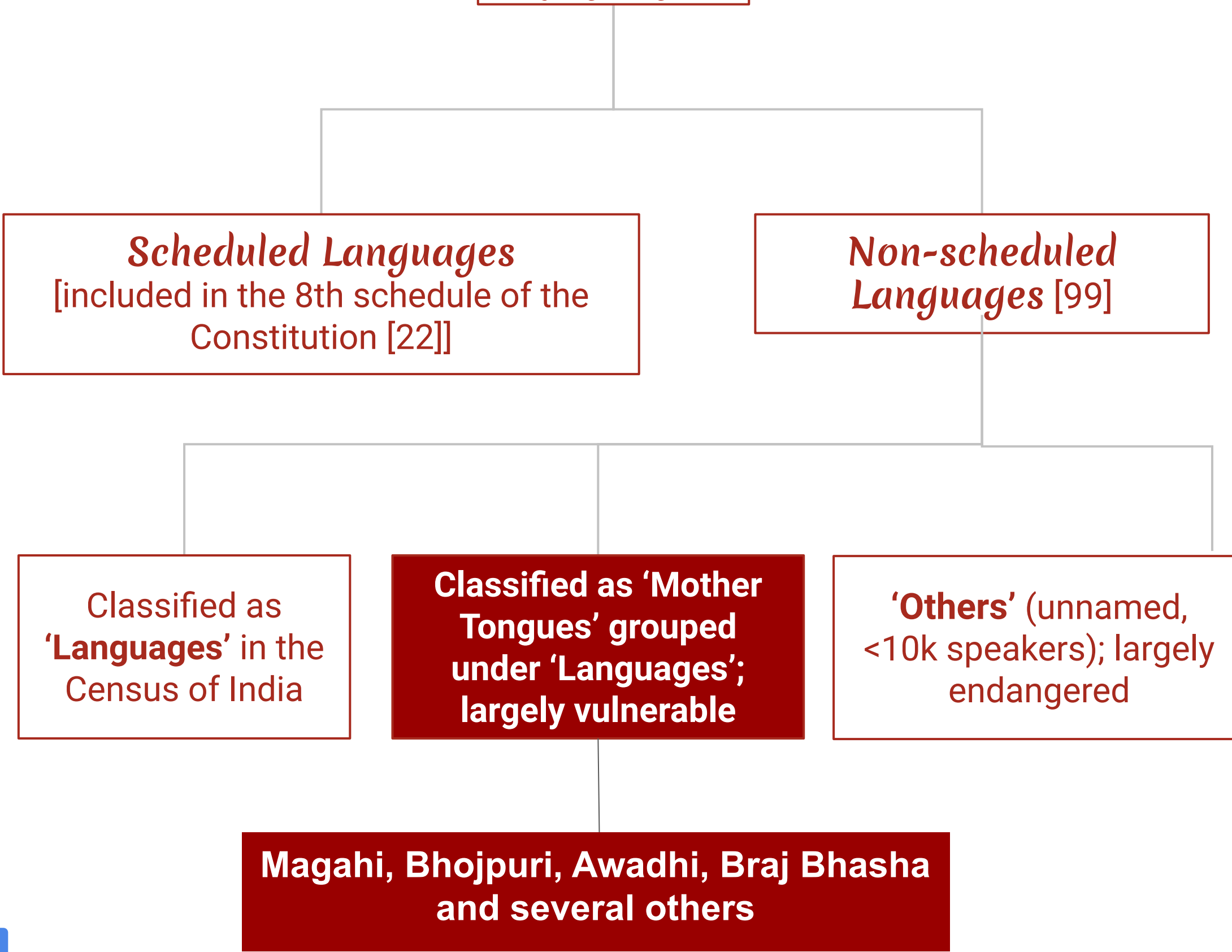
## Magahi Resources and Technologies

>200k sentences	Raw Corpus	<ul style="list-style-type: none"><li>Collection of Folk Tales</li><li>Blogs</li><li>Modern Fiction &amp; Non-fiction</li></ul>
55k sentences	POS-annotated Corpus	<ul style="list-style-type: none"><li>Annotated with UD [manual and automatically]</li></ul>
10k sentences	Parallel Corpus	<ul style="list-style-type: none"><li>English - Magahi</li><li>Russian - Magahi [non-aligned; auto created]</li></ul>
1.3k sentences	Morph-annotated Corpus	<ul style="list-style-type: none"><li>Annotated with UD Feature Sets</li></ul>
1k sentences	UD Treebank	<ul style="list-style-type: none"><li>Annotated with UD relations</li></ul>
0.99 / 0.94 F-score	POS Tagger	<ul style="list-style-type: none"><li>SVM - based POS tagger</li></ul>
0.94 F-score	Morph Analyser	<ul style="list-style-type: none"><li>Rule-based system</li><li>Partially implemented</li></ul>
	MAG - EN MT	<ul style="list-style-type: none"><li>Phrase-based MT System</li></ul>

## Braj Bhasha Resources and Technologies

9k sentences	Raw Corpus	<ul style="list-style-type: none"><li>Modern Fiction and Non-fiction</li></ul>
9k sentences	POS-annotated Corpus	<ul style="list-style-type: none"><li>Annotated with BIS POS Tagset</li></ul>
1.6k sentences	Morph-annotated Corpus	<ul style="list-style-type: none"><li>Annotated with UD Morph Features</li></ul>
88.27%	POS Tagger	<ul style="list-style-type: none"><li>SVM-based POS Tagger</li></ul>

## Indian Languages



## Bhojpuri Resources and Technologies

>100k sentences	Raw Corpus	<ul style="list-style-type: none"><li>News and Blogs</li><li>Collection of Folk Tales</li><li>Modern Fiction and Non-fiction</li></ul>
65k sentences	POS-annotated Corpus	<ul style="list-style-type: none"><li>Annotated with UD as well as BIS POS tags</li></ul>
65k sentences	Parallel Corpus	<ul style="list-style-type: none"><li>English - Bhojpuri</li><li>Hindi-Bhojpuri (&gt;65k)</li></ul>
4k tokens	UD Treebank	<ul style="list-style-type: none"><li>Annotated with UD Tagset (including Morph, POS and relations)</li></ul>
86.78%	POS Tagger	<ul style="list-style-type: none"><li>CRF - based POS tagger</li></ul>
	Eng↔Bho MT	<ul style="list-style-type: none"><li>Phrase-based MT</li><li>Factor-based MT</li><li>Neural MT</li><li>Tree-based MT</li></ul>

## Awadhi Resources and Technologies

70k tokens	Raw Corpus	<ul style="list-style-type: none"><li>Fiction and Stories</li></ul>
20k tokens	POS-annotated Corpus	<ul style="list-style-type: none"><li>Annotated with UD POS Tagset</li></ul>
0.88 F-score	POS Tagger	<ul style="list-style-type: none"><li>SVM based POS tagged</li></ul>

**FOR COLLABORATION AND QUERIES:**  
**panlingua@outlook.com**