

Indic NLP: A Multilinguality and Language Relatedness Perspective

Anoop Kunchukuttan

Machine Translation Group, Microsoft, Hyderabad

ankunchu@microsoft.com



Joint work with



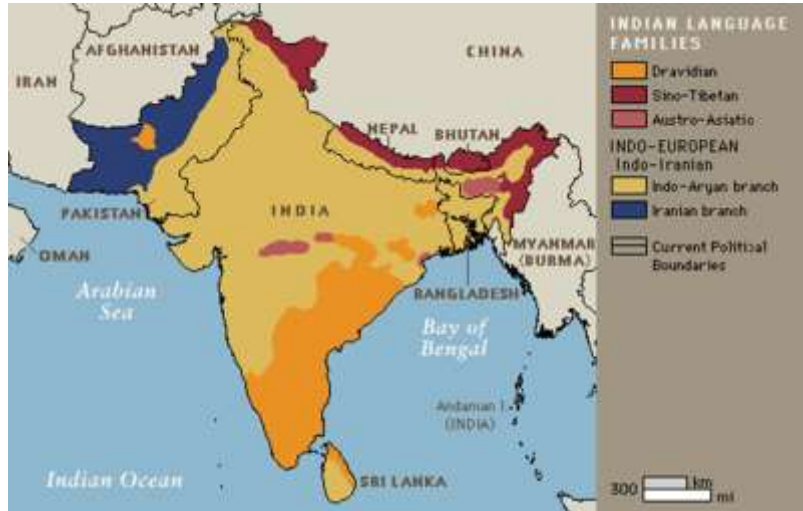
AI4Bharat



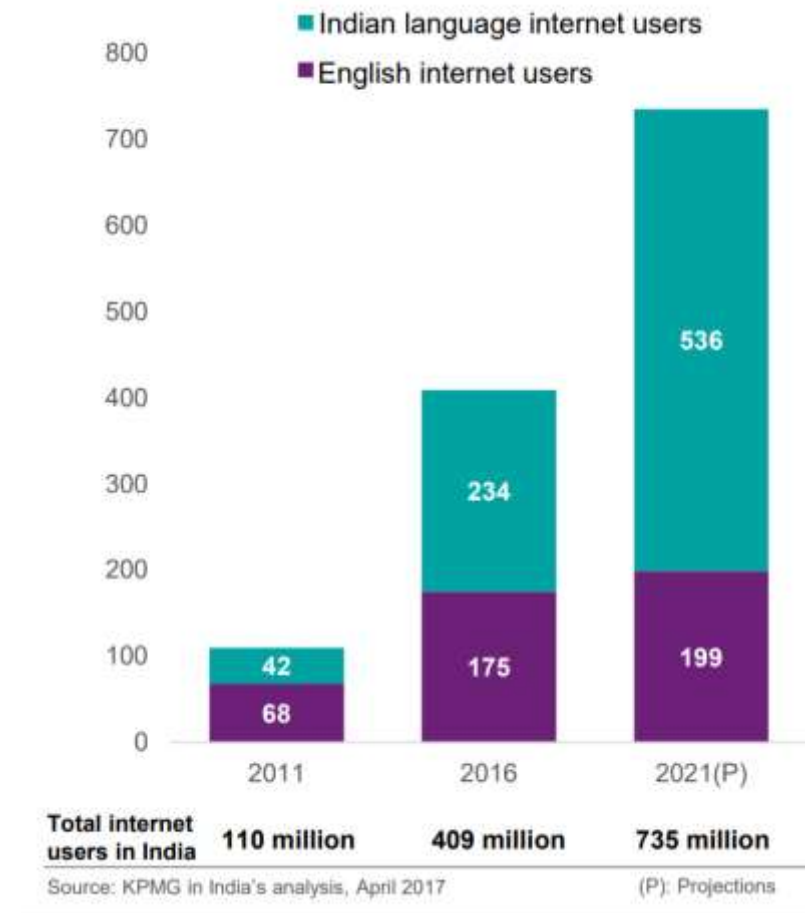
CFILT

*Vaibhav Summit, NLP Session
19th October 2020*

Multi-linguality – A matter of fact in India

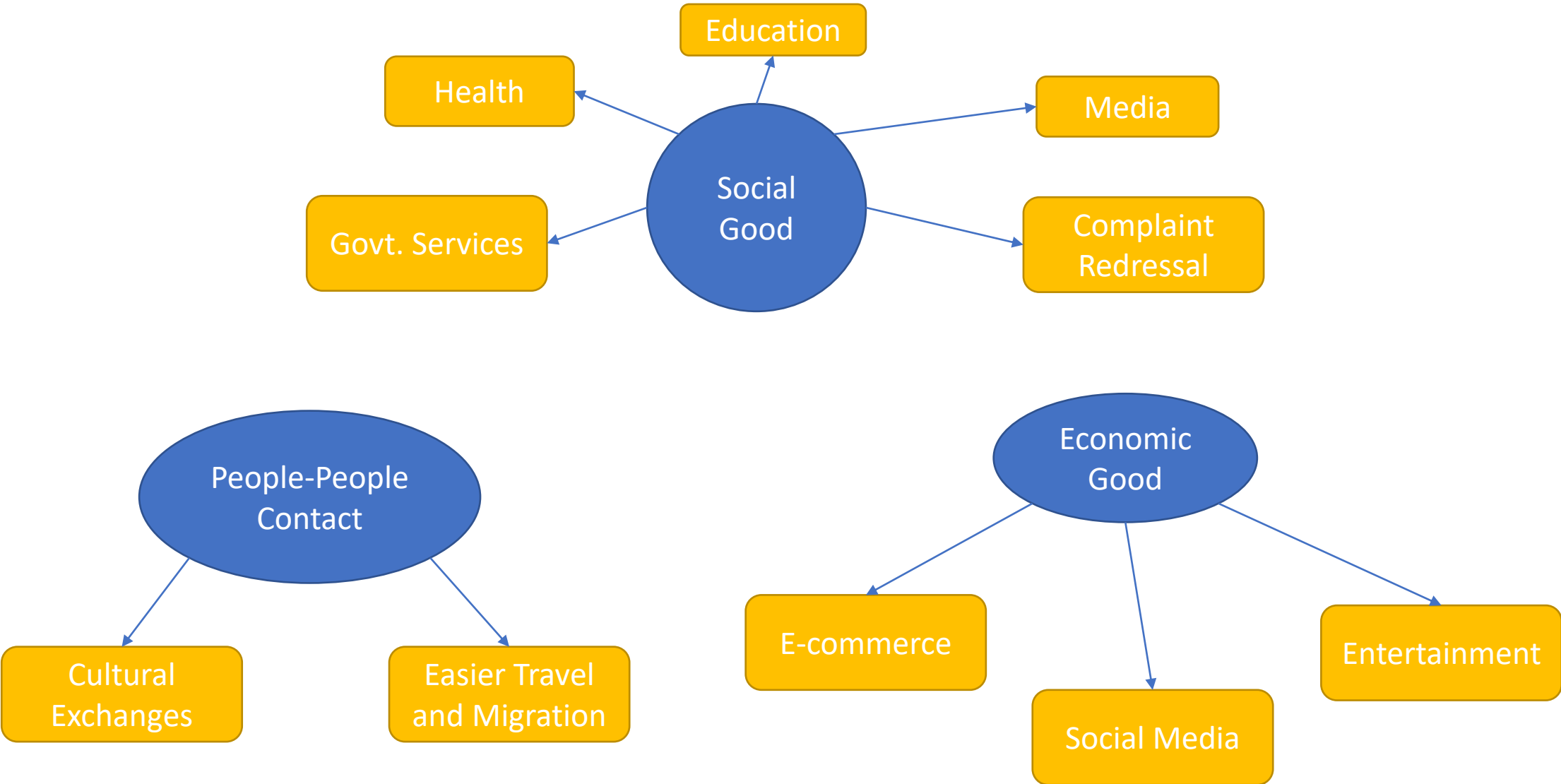


- 22 scheduled languages
- English as lingua franca: 125 million speakers
- 8 languages in the world's top 20 languages
- 30 languages with more than 1 million speakers

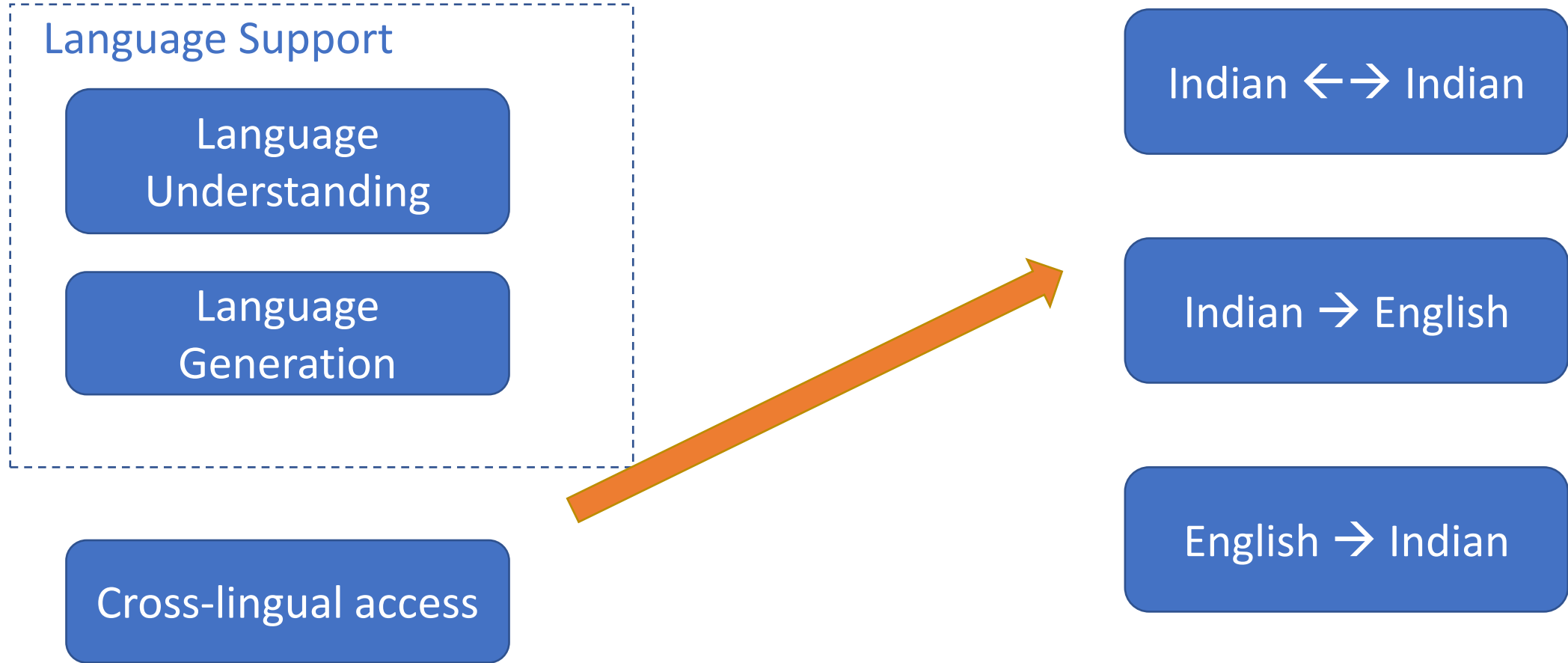


Internet User Base in India (in million)

Addressing Multilinguality is important to maximizing impact of language technologies

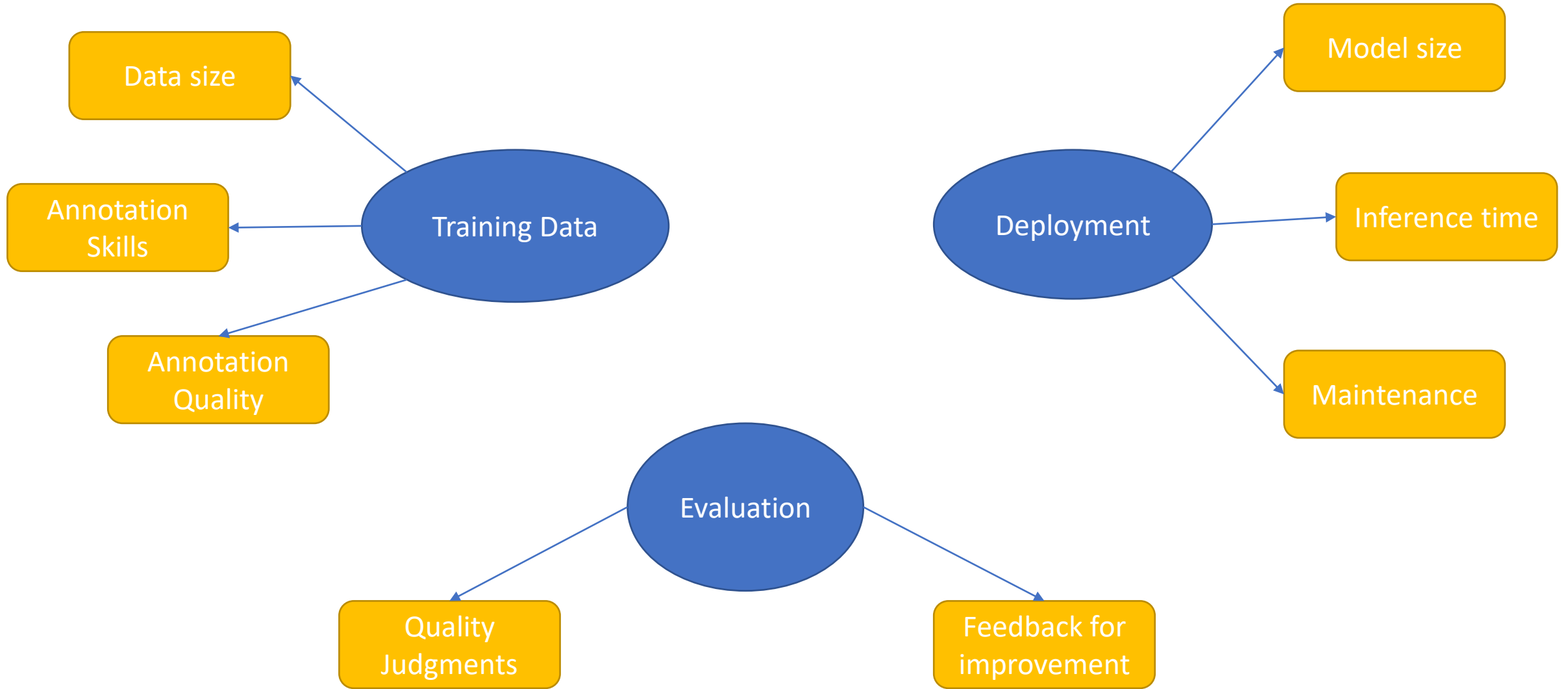


Multilingual NLP Scenarios



*Translation, transliteration, code-mixing
cross-lingual search/QA*

Scalability Challenges for NLP solutions



Effort and cost increase as languages increase

Need for a Unified Approach for Indic NLP

- *Can we share resources across languages?*
- *Can that also reduce effort & cost for deployment and maintenance?*
- *Can diversity of languages lead to better generalization?*

Can we utilize relatedness between Indian languages?

Multilingual challenges are not uniquely Indian → India is a microcosm of the world

Slavic

Germanic

South-East
Asian

Bantu

Global Language
Technology
Products

European
experience with
multi-linguality

How are Indian languages related?

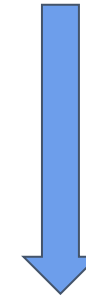
Related Languages

Related by Genealogy



Dravidian, Indo-Aryan,
Tibeto-Burman, Austro-Asiatic

Related by Contact



Indian Subcontinent as a linguistic area

English	Vedic Sanskrit	Hindi	Punjabi	Gujarati	Marathi	Odia	Bengali
<i>bread</i>	Rotika	chapātī, roṭī	roṭī	paū, roṭlā	chapāti, poli, bhākari	pauruṭī	(pau-)ruṭī

Sanskrit word	Language	Loanword	English
cakram	Tamil	cakkaram	wheel
matsyah	Telugu	matsyalu	fish

Key Similarities between related languages

भारताच्या स्वातंत्र्यदिनानिमित्त अमेरिकेतील लॉस एन्जल्स शहरात कार्यक्रम आयोजित करण्यात आला

bhAratAcyA svAta.ntryadinAnimitta ameriketIla lOsA enjalsA shaharAta kAryakrama Ayojita karaNyAta AIA

Marathi

भारता च्या स्वातंत्र्य दिना निमित्त अमेरिके तील लॉस एन्जल्स शहरा त कार्यक्रम आयोजित करण्यात आला

bhAratA cyA svAta.ntrya dinA nimitta amerike tIla lOsA enjalsA shaharA ta kAryakrama Ayojita karaNyAta AIA

Marathi
segmented

भारत के स्वतंत्रता दिवस के अवसर पर अमरीक के लॉस एन्जल्स शहर में कार्यक्रम आयोजित किया गया

bhArata ke svata.ntratA divasa ke avasara para amarIkA ke losa enjalsA shahara me.n kAryakrama Ayojita kiyA gayA

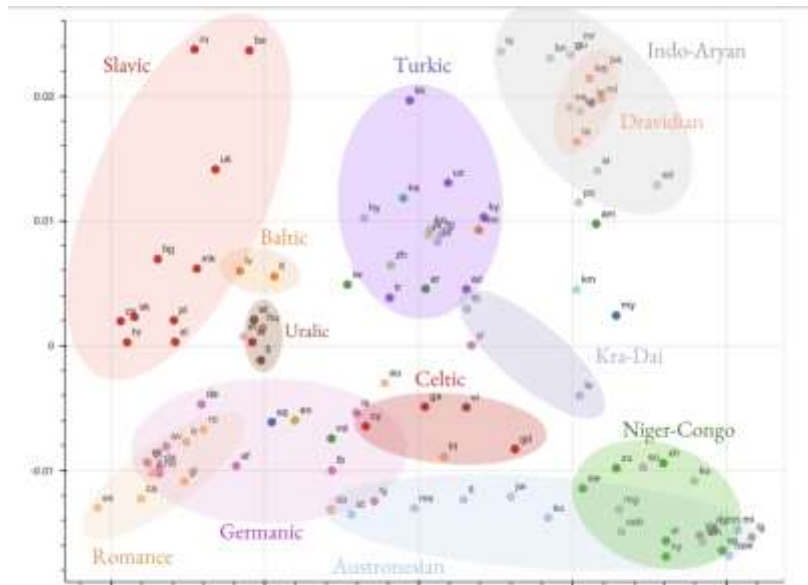
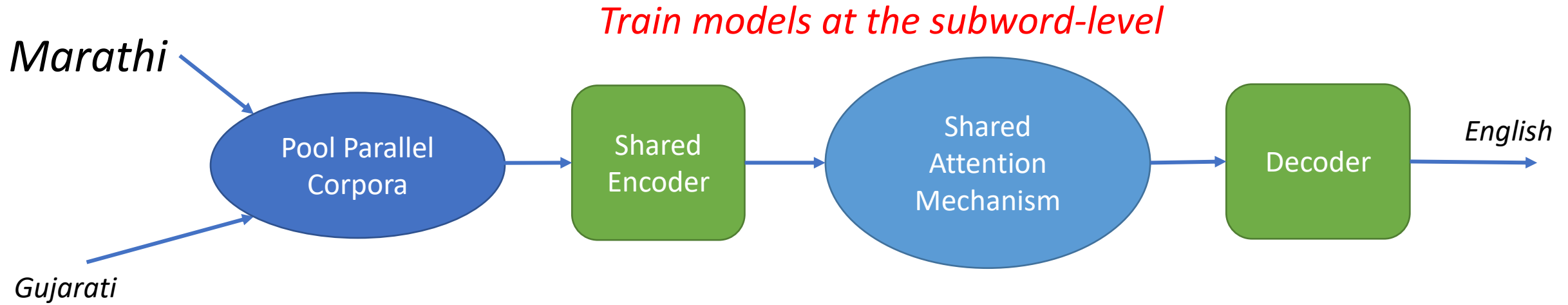
Hindi

Lexical: share significant vocabulary (cognates & loanwords)

Morphological: correspondence between suffixes/post-positions

Syntactic: share the same basic word order

Transfer Learning Recipe



(Kudungta et al, 2019)

Encoder Representations cluster by language family

Transfer Learning works best for related languages

Moving beyond the simple transfer learning paradigm

Can we better utilize the similarities between Indian languages?

Similarity in Scripts

Devanagari	अ आ इ ई उ ऊ ऋ ॠ एँ ऐ ए ऐ आँ ओ औ औ क ख ग घ ङ च छ ज झ
Bengali	অ আ ই ঐ উ ঊ ঋ ৠ এ ঐ ও ঔ ক খ গ ঘ ঙ চ ছ জ ঝ ঞ ট ঠ ড
Gurmukhi	ਅ ਆ ਇ ਈ ਉ ਊ ਏ ਓ ਅੰ ਕ ਖ ਗ ਘ ਙ ਚ ਛ ਜ ਝ ਟ ਠ ਡ ਢ ਤ ਥ
Gujarati	અ આ ઇ ઈ ઉ ઊ ઋ એ એ ઐ ઔ ઑ ઒ ઓ ઐ ક ખ ગ ઘ ઙ ચ છ જ ઝ ઞ ટ ઠ
Oriya	ଅ ଆ ଇ ଈ ଉ ଊ ଋ ଌ ଏ ଐ ଓ ଔ କ ଖ ଗ ଘ ଙ ଚ ଛ ଜ ଝ ଞ ଟ ଠ ଡ ଢ
Tamil	அ ஆ இ ஐ உ ஊ எ ஏ ஐ ஒ ஓ ஔ க ங ச ங ஞ ட ண த ந
Telugu	అ ఆ ఇ ఈ ఉ ఊ ఋ ఎ ఏ ఐ ఒ ఓ ఔ క ఖ గ ఘ ఙ చ ఛ జ ఝ ఞ
Kannada	ಅ ಆ ಇ ಈ ಉ ಊ ಋ ಌ ಎ ಏ ಐ ಒ ಓ ಔ ಕ ಖ ಗ ಘ ಙ ಚ ಛ ಜ ಝ ಞ
Malayalam	അ അ ഇ ഇയ ഉ ഉയ ള ള ള ള ള ള ള ള ള ള ള ള ള ള ള ള ള ള ള

Multilingual transliteration

Convert to a common script & Pool

Malayalam	കോഴിക്കോട്	कोळिकोट्	kozhikode
Hindi	केरल	केरल	kerala
Kannada	ಬೆಂಗಳೂರು	बेंगळूरु	bengaluru

Make inputs more similar, reduce vocabulary size

		Voiceless plosives		Voiced plosives		Nasals	
		unaspirated	aspirated	unaspirated	aspirated		
Velar	क ka	ख kha	ग ga	घ gha	ङ ṅa		
Palatal	च ca	छ cha	ज ja	झ jha	ञ ña		
Retroflex	ट ṭa	ठ ṭha	ड ḍa	ढ ḍha	ण ṇa		
Dental	त ta	थ tha	द da	ध dha	न na		
Labial	प pa	फ pha	ब ba	भ bha	म ma		
Sonorants and fricatives							
	Palatal	Retroflex	Dental	Labial			
Sonorants	य ya	र ra	ल la	व va			
Sibilants	श śa	ष ṣa	स sa				
Other letters							
	ह ha	ळ la					

Place of Articulation	velar, palatal, retroflex, dental, labial
Manner of Articulation	plosive, fricative, flap, approximant (central or lateral)
Aspiration	True, False
Voicing	True, False
Nasalization	True, False

Scientific Design of scripts enables feature representation of characters/sounds

- *Unsupervised Transliteration*
 - *Initialization and priors on Character transliteration probabilities*
- *Cognate Identification*

Word Segmentation using Aksharas

akshara, the fundamental organizing principle of Indian scripts

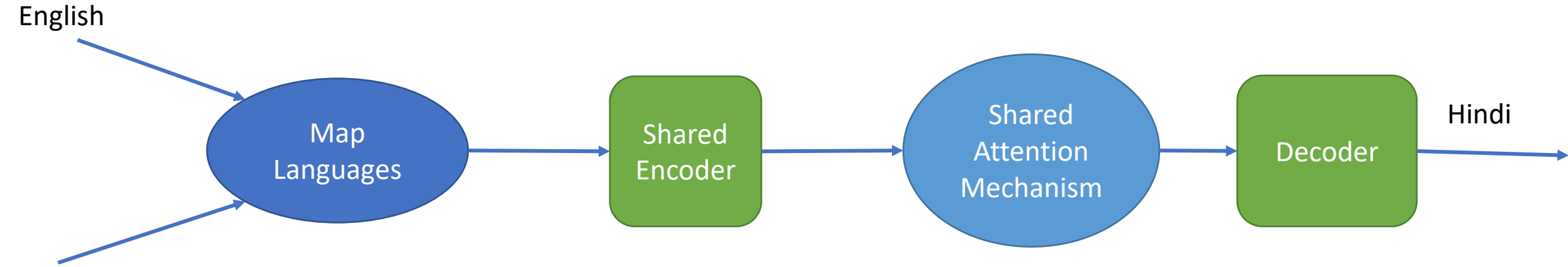
(CONSONANT) + VOWEL

Examples: की (ki), प्रे (pre)

Hindi	Kannada	English
वि द्या ल य	ವಿ ದ್ಯಾ ಲ ಯ	vi dya lay
अ र्जु न	ಅ ರ್ಜು ನ	a rju n

Useful as basic units for transliteration and translation

Transfer from English to Indian languages



Gujarati

Little G → H corpus

Cannot ensure similar Gujarat and English words have similar representations

Solution: Pre-order English sentence to match Gujarati word-order

Syntactic divergence can be overcome with shared rules

1. Rudra Murthy V, Anoop Kunchukuttan, Pushpak Bhattacharyya. *Addressing word-order Divergence in Multilingual Neural Machine Translation for extremely Low Resource Languages*. NAACL. 2019.
2. Anoop Kunchukuttan, Abhijit Mishra, Rajen Chatterjee, Ritesh Shah, Pushpak Bhattacharyya. *Shata-Anuvadak: Tackling Multiway Translation of Indian Languages*. Language and Resources and Evaluation Conference. 2014.

Transfer between contact languages



Tamil to English NMT with transfer-learning using Hindi
Addressing syntactic divergence in NMT using Hindi-driven rules

Language Relatedness can be successfully utilized between languages where contact relation exists

1. Vikrant Goyal, Anoop Kunchukuttan, Rahul Kejriwal, Siddharth Jain, Amit Bhagwat. *Contact Relatedness can help improve multilingual NMT: Microsoft STCI-MT @ WMT20*. WMT 2020. 2020.
2. Rudra Murthy V, Anoop Kunchukuttan, Pushpak Bhattacharyya. *Addressing word-order Divergence in Multilingual Neural Machine Translation for extremely Low Resource Languages*. NAACL. 2019.

Putting these ideas together into usable systems ...

Indic NLP Library

https://github.com/anoopkunchukuttan/indic_nlp_library

- **Utilize similarity** between Indian languages for scaling to multiple Indian languages
- Design to **support maximum number of Indian languages**
- Modular and Extensible
- Easy of use:
 - Installation `pip install indic-nlp-library`
 - Consistent Use
 - Separation between code and data resources

MIT License

Anoop Kunchukuttan. *The IndicNLP Library*. https://github.com/anoopkunchukuttan/indic_nlp_library/blob/master/docs/indicnlp.pdf. 2020.

Capabilities

Text Processing

- Text Normalizer
- Sentence Splitter
- Word Tokenizer
- Word Detokenizer

Word Segmentation

- Morphological Segmentation
- Syllabification

Script Processing

- Query Script Information
- Script Converter
- Romanization
- Indicization
- Acronym Transliterator
- Phonetic Similarity
- Lexical Similarity

Language Support

Indo-Aryan			Dravidian
Assamese (as)	Marathi (mr)	Sindhi (sd)	Kannada (kn)
Bengali (bn)	Nepali (ne)	Sinhala (si)	Malayalam (ml)
Gujarati (gu)	Odia (or)	Sanskrit (sa)	Telugu (te)
Hindi (hi)	Punjabi (pa)	Konkani (kok/kK)	Tamil (ta)

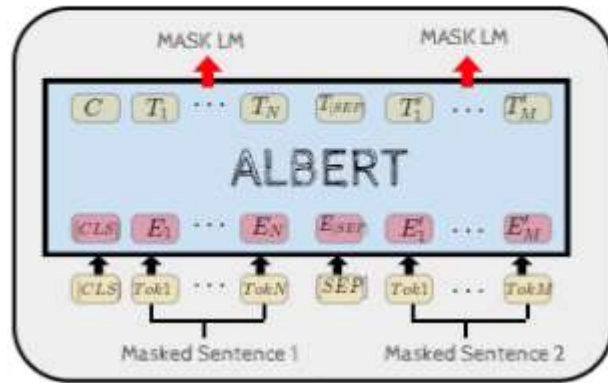
	as	bn	gu	hi	mr	ne	or	pa	sd	si	sa	kok	kn	ml	te	ta
Text Processing	✓	✓	✓	✓	✓	✓	✓	✓	✓	✗	✓	✓	✓	✓	✓	✓
Morphological Segmentation	✓	✓	✓	✓	✓	✗	✓	✓	✗	✗	✗	✓	✓	✓	✓	✓
Syllabification	✓	✓	✓	✓	✓	✓	✓	✓	✓	✗	✓	✓	✓	✓	✓	✓
Script Processing	✓	✓	✓	✓	✓	✓	✓	✓	✓	✗	✓	✓	✓	✓	✓	✓

Future: add support for multilingual pre-trained embeddings, fundamental tools like POS, Dependency parsing, NER, etc.

IndicBERT

<https://indicnlp.ai4bharat.org/indic-bert>

<https://huggingface.co/ai4bharat/indic-bert>



यं हि वा ॐ अ^A
गु म ष ष ड ङ
Joint Pre-training

- Large Indian language content (8B tokens)
 - 11 Indian languages
 - + Indian English content
- Multilingual Model
- Compact Model (~20m params)
- Competitive/better than mBERT/XLM-R
- Simplify **fine-tune** for your application

Multilingual Approaches are important for language technologies to scale and make social impact

The field is nascent, there are many directions to explore

- *Better representation methods to utilize relatedness*
- *Bridging typological divergence between English and Indian languages*
- *Utilizing relatedness for generation tasks*
- *Cross-lingual Evaluation benchmarks*

Thank You!

<http://anoopk.in>