

Challenges of Processing South Asian Languages (CPSAL)

Kengatharaiyer Sarveswaran
University of Konstanz, Germany
University of Jaffna, Sri Lanka.

Tafseer Ahmed
Senior NLP Scientist
Alexa Translations

Course outline

- **Topics (Tentative):**

- Day 01: Languages, Scripts, and Encoding of South Asian Languages.
- **Day 02: Phonology, Transliteration and Morphology of South Asian Languages.**
- Day 03: Part of Speech and Multiword tokenisation
- Day 04: Syntax, Morphosyntax, and Semantics of South Asian Languages.
- Day 05: Deep Learning for South Asian Languages and winding up the course.

Challenges of Processing South Asian Languages (CPSAL)

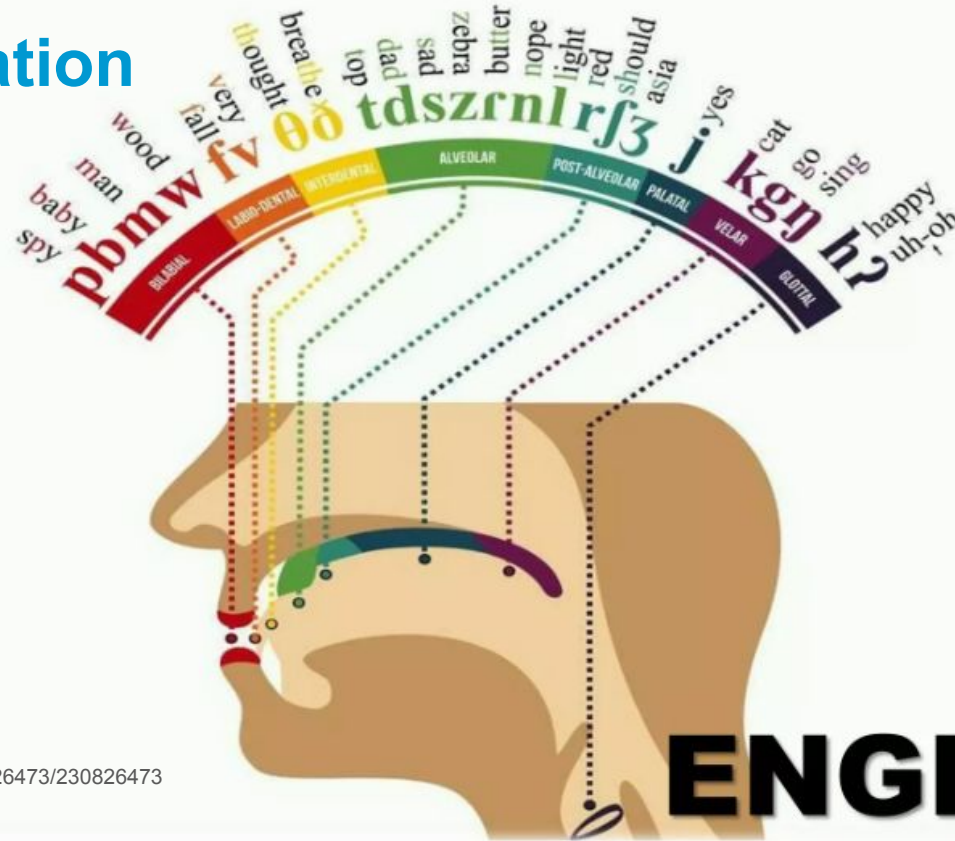
Day 02: Transliteration/Phonology and Morphology of South Asian Languages.

Kengatharaiyer Sarveswaran
University of Jaffna, Sri Lanka.
University of Konstanz, Germany

Tafseer Ahmed
Senior NLP Scientist
Alexa Translations³

Position of Articulation

- Bilabial
- Labio-Dental
- InterDental
- Alveolar
- Palatal
- Velar
- Glottal



<https://www.slideshare.net/slideshow/english-consonants-230826473/230826473>

ENGLISH CONSONANTS

Position and manner of articulation

| | bilabial | labio-dental | dental | alveolar | post-alveolar | palatal | velar | glottal |
|---------------------|----------|--------------|--------|----------|---------------|---------|-------|---------|
| stop | p b | | | t d | | | k g | ʔ |
| nasal | m | | | n | | | ŋ | |
| flap | | | | r | | | | |
| fricative | | f v | θ ð | s z | ʃ ʒ | | | h |
| approximant | | | | ɹ | | j | | |
| lateral approximant | | | | l | | | | |

Consonants in Sindhi

| PLACE \ MANNER | | Bilabial | Labio-dental | Dental | Alveolar | Retroflex | Alveo-Palatal | Palatal | Velar | Glottal |
|----------------|-----------|--------------------------------------|--------------|--|------------------|--|--|---------|--------------------------------------|---------|
| Stop | Plosive | p b p ^h b ^h | | t̪ d̪ t̪ ^h d̪ ^h | | t̠ d̠ t̠ ^h d̠ ^h | | | k g k ^h g ^h | ʔ |
| | implosive | ɓ | | | | ɗ | | f | ɠ | |
| Nasal | | m | | | n n ^h | ɳ ɳ ^h | | ɲ | ŋ | |
| Fricative | | | f v | | s z | | | ʃ | x y | h |
| Affricate | | | | | | | tʃ dʒ tʃ ^h dʒ ^h | | | |
| Trill | | | | | r r ^h | | | | | |
| Flap | | | | | | ɾ ɾ ^h | | | | |
| Lateral | | | | | l l ^h | | | | | |
| Approximant | | | | | | | | j | | |

Table 1: Consonantal Inventory of Sindhi

<https://uogenglish.wordpress.com/wp-content/uploads/2011/07/sindhi-phonemic-inventory.pdf>

Tamil consonants

| | Labial | Dental | Alveolar | Retroflex | (Alveolo-) palatal | Velar | Glottal |
|-------------------------------|------------------|--------|------------------------------------|---------------------|---------------------|------------------|---------------------|
| Nasal | m ம் | (n) ந் | n ன் | ɳ ண் | ɲ ஞ் | (ŋ) ங் | |
| Plosive/ Affricate | p ப் | t̪ த் | (tʀ) த்ற | ʈ ட் | t͡ʃ ச் ⁵ | k க் | |
| Fricative | (f) ¹ | | s ⁵ ஸ் (z) ¹ | (ʂ) ¹ ஷ் | (ɕ) ¹ ஸ் | (x) ² | (h) ² ஹ் |
| Tap | | | r ர் | | | | |
| Trill | | | r ற் | | | | |
| Approximant | ɸ வ் | | | ɭ ழ் | j ய் | | |
| Lateral approximant | | | l ல் | ɭ ள் | | | |

Interesting/Uncommon Features - Some examples

- Aspirated Consonants

| | | Sindhi | Hindi | Urdu |
|-------------|----------------|--------|-------|------|
| unaspirated | b | ب | ब | ب |
| aspirated | b ^h | پ | भ | بھ |

- Retroflex Consonants

| | | Sindhi | Hindi | Urdu |
|-----------|---|--|-------|--|
| Dental |  |  | त |  |
| Retroflex | t |  | ट |  |

Interesting/Uncommon Features

Implosive Consonants:

| | | Sindhi |
|-----------|---|--------|
| Plosive | b | ڀ |
| Implosive | ɓ | ɓ̥ |

Foreign Sounds

Consonants:

| Foreign | Origin | Nearest Native | Sindhi | Hindi |
|----------|----------------------------|----------------|--------|-------|
| f | Persian/Arabic, English | p ^h | ف | फ़ |
| x | Persian/Arabic | k ^h | خ | ख़ |

Vowels:

Urdu/Hindi: bε:l (bull), be:l (creeper plant)

People know how to pronounce the borrowed English words *bell* (bεl), and *cell* (sεl) etc. However, there is no symbol to represent these in traditional orthography.

Sandhi - Example Tamil

- Internal Sandhi / External Sandhi
 - Phonological change that occurs when two morphs are concatenated together
- External Sandhi (Gemination of stops)
 - Word initial stops are geminated in Tamil when preceded by either a short vowel or a glide - this rule does not always apply even when these phonological conditions are met!
 - The phonological phenomenon is shown to be 'directly syntax-driven??'
 - If word initial k, t, p is preceded by a noun in the accusative case, a noun in the dative case, a demonstrative or interrogative adjective (inta, anta, enta), an infinitive or one of the adverbs ippaDi 'like this', appaDi 'like that', eppaDi 'how', the realisation is with a voiceless plosive....

Sandhi - A few example in Tamil

anta = ppustakattay = kkuDu
that book (acc) give imp.
Give (me) that book.
(Dem=N,Obj.NP=V)

avanukku = ppaNam kuDutteen
him (dat) money give past 1s.
(I) gave him money.
(I.O=D.O)

raamanukku = ppacikkiradu
ram (dat) hunger pres 3p.
Rama is hungry.
(Dative subj.=V)

kaar-ootti centran
kaar-driver go past 1sm
Car driver went
(kaarootti - the one who drives car)

kaar-ootti =(c) centran
Car-drive.Inf =(c) go past 1sm
(he) drove a car.

Transliteration standards

- ISO_15919: Transliteration of Devanagari and related Indic scripts into Latin characters
 - https://en.wikipedia.org/wiki/ISO_15919
- ISO 233-3:2023: Transliteration of Arabic characters into Latin characters
 - https://en.wikipedia.org/wiki/ISO_233

Transliteration standards

- Not all sounds have an orthographic character

| | | | |
|-----|--------|--------------|-------------------|
| p | paṭu | <i>pattu</i> | ‘ten’ |
| [b] | ṭambi | <i>tampi</i> | ‘younger brother’ |
| (b) | baḍil | <i>patil</i> | ‘answer’ |
| ṭ | ṭapu | <i>tappu</i> | ‘mistake’ |
| [ḍ] | paṇḍu | <i>pantu</i> | ‘ball’ |
| (ḍ) | ḍinam | <i>tinam</i> | ‘day’ |
| ṭ | | | |
| [ḍ] | vaṇḍi | <i>vaṇṭi</i> | ‘cart’ |
| k | kaḷ | <i>kaal</i> | ‘leg’ |
| [g] | aṅgeṛ | <i>aṅkee</i> | ‘there’ |
| (g) | gaṇam | <i>kanam</i> | ‘heaviness’ |
| ṭṣ | ṭṣinṇə | <i>cinna</i> | ‘small’ |

Morphology

What can you find out about morphology?

| | | |
|---|-------------------|---------------|
| 1 | அவன் அழுதான் | avan alutān |
| | He cried | |
| 2 | அவள் அழுதாள் | aval alutāḷ |
| | She Cried | |
| 3 | அவன் அழுகிறான் | avan alukirān |
| | He is crying | |
| 4 | அவள் அழுகிறாள் | aval alukirāḷ |
| | She is crying | |
| 5 | அவன் அழுவான் | avan aluvān |
| | He will cry | |
| 6 | அவள் அழுவாள் | aval aluvāḷ |

| | | |
|---|------------------------|---------------------|
| 7 | அவர்கள் அழுதனர் | avarkaḷ alutaṇar |
| | They cried | |
| 8 | அவர்கள் அழுகின்றனர் | avarkaḷ alukinraṇar |
| | They are crying | |
| 9 | அவர்கள் அழுவார்கள் | avarkaḷ aluvārkaḷ |
| | They will cry | |