

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.

Introductions

- Name
- Where are you from
- Are you already working in SALs?
- Other

Challenges of
Processing South Asian Languages
(CPSAL)

Day 01: Languages and Scripts
of South Asian Languages

Topics

- **South Asia**
- **Languages**
- **Scripts**
- **Encodings**
- **Challenges**

Some Background

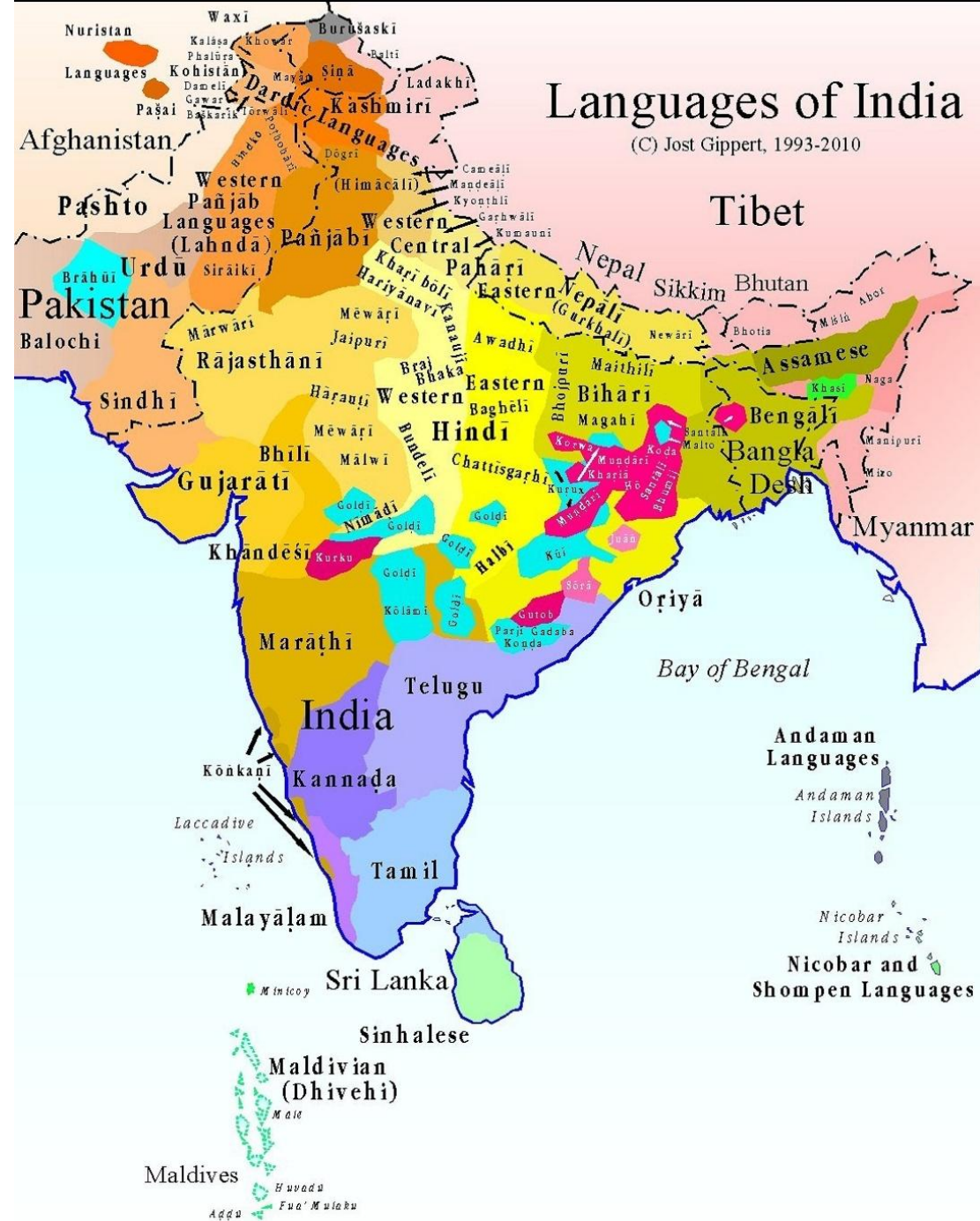
South Asia



South Asia

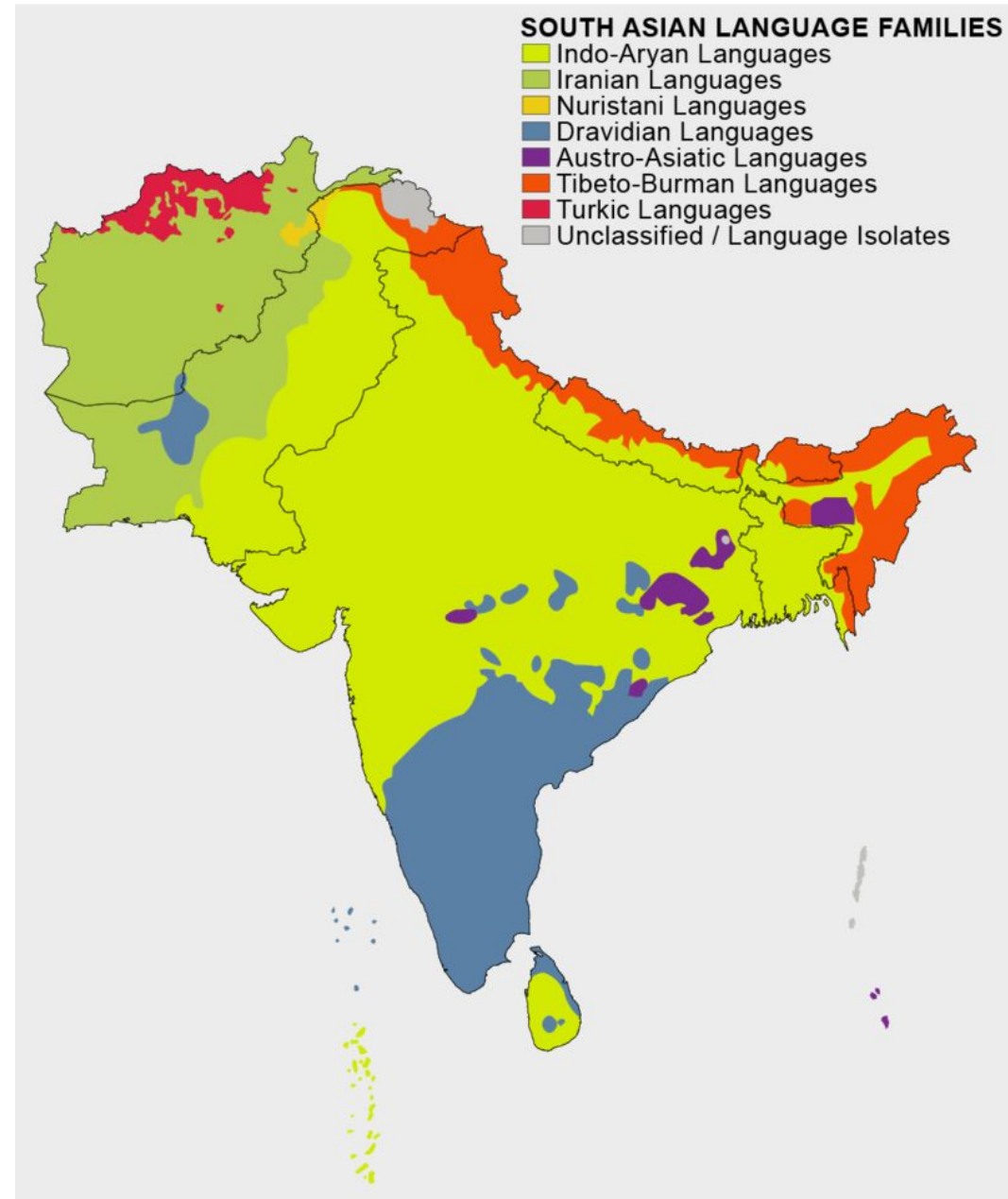
Country	Population	Languages (Primary)
Afghanistan	42+ millions	Pashto, Dari
Bangladesh	172+ millions	Bangla
Bhutan	787 thousands	Dzongkha
India	1.4+ billions	Assamese, Bengali, Bodo, Dogri, Gujarati, Hindi, Kannada, Kashmiri, Konkani, Maithili, Malayalam, Manipuri, Marathi, Nepali, Odia, Punjabi, Sanskrit, Santali, Sindhi, Tamil, Telugu, Urdu
Maldives	521 thousands	Divehi
Nepal	30+ millions	Nepali
Pakistan	240+ millions	Urdu, Punjabi, Saraiki, Sindhi, Baluchi and Pashto
Sri Lanka	21+ millions	Sinhala and Tamil
(3% of world's land area)	1.9+ billions (24+% of the world's population)	650 living languages
Migrants, Diaspora	?	New dialects?

Languages of South Asia

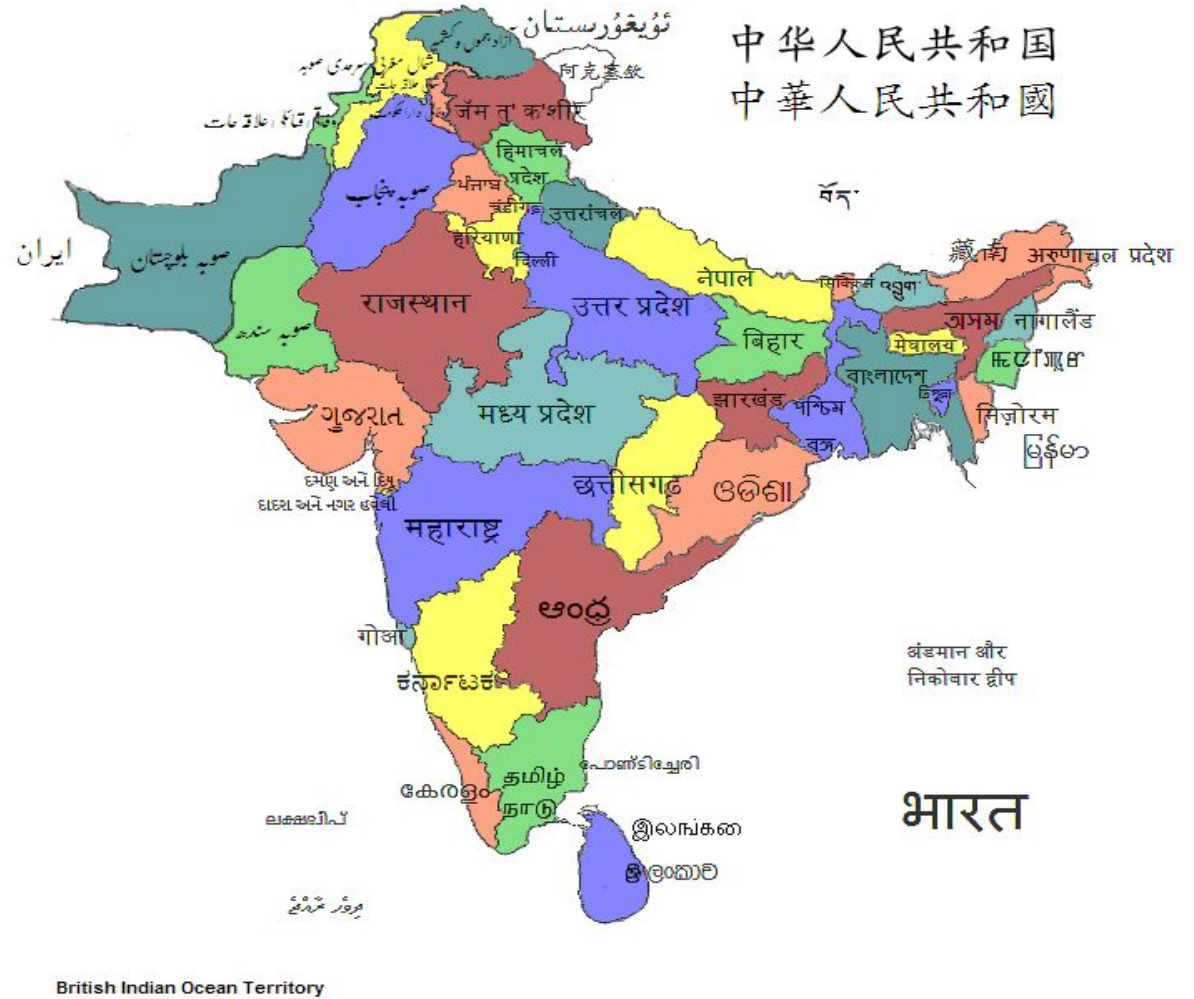


Language families

- **Indo-Aryan**
 - Urdu, Hindi, Rajasthani, Nepalese...
- **Iranian**
 - Pashto, Balochi, Dari...
- **Nuristani**
 - Kamkata-viri, Ashkun, Tregami...
- **Dravidian**
 - Tamil, Malayalam, Telugu, Kannada, Brahui...
- **Austro-Asiatic**
 - Munda, Santali, Khasi, Mundari...
- **Tibeto-Burman**
 - Bodo, Manipuri (Meitei)...
- **Turkic**
 - Uzbek, Turkmen...



Scripts of South Asia



Scripts of South Asia: Examples of scripts used to write Indo-Aryan languages

- कोन्स्टांज़ - Devanagari (Hindi)
 - কোনস্টোন্জ - Bengali
 - ਵੇਨਸਟਾਂਜ - Gurmukhi (Punjabi)
 - કોન્સ્ટાંઝ - Gujarati
 - කොන්ස්ටාන්ස් - Sinhala (Sinhala)
-

Scripts of South Asia: Examples of scripts used to write **Iranian** languages

• کونستانس - (Urdu) Naskh

• کونستانس - (Urdu) Nastaliq

• کونستانس - (Sindhi) Naskh

• کونستانس، - (Pashto) Naskh

Scripts of South Asia: Examples of scripts used to write **Dravidian** languages

- **கோன்ஸ்டான்ஸ்** - Tamil
- **కొన్నంజ్** - Telugu
- **ಕೊನ್ನಂಜ್** - Kannada
- **കോൺസ്റ്റാൻസ്** - Malayalam

Writing systems

- **Brahmi scripts: Tamil, Sanskrit, Punjabi and others.**
 - **Left to right** writing system
 - **Abugida/alphasyllabary** writing system: vowels, consonants, and composites.
 - Composite = Consonant+Vowel modifier = One Unit/Grapheme
 - க் (k) + அ (a) = கா (ka)
 - க் (k) + உ (u) = கூ (ku)
 - க் (k) + ஆ (aa) = கா (kaa)
 - க் (k) + ஐ (ai) = கை (kai)
 - க் (k) + ஓ (oo) = கொ (koo)

Writing systems

- **Arabic scripts: Urdu, Punjabi, Sindhi, Pastho**
 - **Abjad** writing system
 - **Right to left**
 - Consonants and long vowels are written – short vowels exist, but **optional**.
 - Different shapes for the character

1 st character	2 nd character	3 rd character	Result
س	پ	ب	سپب

Language vs Script

- **One script -> Many languages**

- आप कैसे हैं? (*Āp kaise hain?*) - Hindi
- आप केम बा? (*Āp kem bā?*) - Bhojpuri
- तुम्ही कसे आहात? (*Tumhī kase āhāt?*) - Marathi
- कथं असि? (*Katham asi?*) - Sanskrit
- तपाइंलाई कस्तो छ? (*Tapā'inlā'ī kasto cha?*) - Nepali
- ...

Language vs Script

අ	ආ	ඈ	ඉ	ඊ	උ
a	ā	i	ī	u	ū
ඍ	ඎ	ඏ	ඐ		
e	ē	o	ō		
ඌ	ඍ	ඎ	ඏ		
l	ll	o	oo		
ඒ	උ	ඌ	ඍ		
e	ai	o	au		
ආ	භ				
ā	ḥ				

- **Many scripts -> One language**

- Sanskrit: Devanagari and Grantha
- Punjabi and Sindhi are written using different writing systems
 - Punjabi: Shahmukhi and Gurmukhi
 - Shahmukhi - **ਤਸੀ ਕੀਵੇਂ ਹੋ** (tusī kiveN ho)
 - Gurmukhi - **ਤੁਸੀਂ ਕਿਵੇਂ ਹੋ** (tusī kivēṁ hō)
 - Sindhi: Perso-Arabic, Devanagari, Khudabadi, and Khojki
 - Naskh - **تون ڪيئن آهين**
 - Devanagari - **तुं कीअं आहीनि** (tūṁ kī'am āhīni)
- All the languages are also written in the Roman script.

अ	a	
इ	i	
उ	u	
ऋ	ṛ	ṝ
ॠ	ṝ	ṝ̄
ए	e	ē
ओ	o	ō
अं / ं ^{1,2}	ṁ	ṁ̄
आ	ā	
ई	ī	
ऊ	ū	
ॠ	ṝ	ṝ̄
ॡ	ṝ̄	ṝ̄̄
ऐ	ai	
औ	au	
अः / ः ¹	ḥ	

Scripts have evolved: Example - Tamil

HISTORY OF TAMIL SCRIPT									
நூற்றாண்டு	a ā i ī u ū e ē ai o ō								
Century	அ ஆ இ ஈ உ ஊ எ ஏ ஐ ஒ ஓ								
BC 3 rd C	𑌀	𑌁	𑌂	𑌃	𑌄	𑌅	𑌆	𑌇	𑌈
AD 2 nd C	𑌀	𑌁	𑌂	𑌃	𑌄	𑌅	𑌆	𑌇	𑌈
AD 3 rd C	𑌀	𑌁	𑌂	𑌃	𑌄	𑌅	𑌆	𑌇	𑌈
AD 4 th C	𑌀	𑌁	𑌂	𑌃	𑌄	𑌅	𑌆	𑌇	𑌈
AD 5 th C	𑌀	𑌁	𑌂	𑌃	𑌄	𑌅	𑌆	𑌇	𑌈
AD 6 th C	𑌀	𑌁	𑌂	𑌃	𑌄	𑌅	𑌆	𑌇	𑌈
AD 7 th C	𑌀	𑌁	𑌂	𑌃	𑌄	𑌅	𑌆	𑌇	𑌈
AD 8 th C	𑌀	𑌁	𑌂	𑌃	𑌄	𑌅	𑌆	𑌇	𑌈
AD 9 th C	𑌀	𑌁	𑌂	𑌃	𑌄	𑌅	𑌆	𑌇	𑌈
AD 10 th C	𑌀	𑌁	𑌂	𑌃	𑌄	𑌅	𑌆	𑌇	𑌈
AD 11 th C	𑌀	𑌁	𑌂	𑌃	𑌄	𑌅	𑌆	𑌇	𑌈
AD 12 th C	𑌀	𑌁	𑌂	𑌃	𑌄	𑌅	𑌆	𑌇	𑌈
AD 13 th C	𑌀	𑌁	𑌂	𑌃	𑌄	𑌅	𑌆	𑌇	𑌈
AD 14 th C	𑌀	𑌁	𑌂	𑌃	𑌄	𑌅	𑌆	𑌇	𑌈
AD 15 th C	𑌀	𑌁	𑌂	𑌃	𑌄	𑌅	𑌆	𑌇	𑌈
AD 16 th C	𑌀	𑌁	𑌂	𑌃	𑌄	𑌅	𑌆	𑌇	𑌈
AD 17 th C	𑌀	𑌁	𑌂	𑌃	𑌄	𑌅	𑌆	𑌇	𑌈
AD 18 th C	𑌀	𑌁	𑌂	𑌃	𑌄	𑌅	𑌆	𑌇	𑌈
AD 19 th C	𑌀	𑌁	𑌂	𑌃	𑌄	𑌅	𑌆	𑌇	𑌈

நூற்றாண்டு	K ṅ c ṇ ṭ ṇ t n p m y r l v ḷ ḷ ṛ ṇ												
Century	க்ங்சஞ்ஞட்ணத்ந்ப்ம்யர்ல்வழ்ளற்ன்												
BC 3 rd C	𑌕	𑌖	𑌗	𑌘	𑌙	𑌚	𑌛	𑌜	𑌝	𑌞	𑌟	𑌠	𑌡
AD 2 nd C	𑌕	𑌖	𑌗	𑌘	𑌙	𑌚	𑌛	𑌜	𑌝	𑌞	𑌟	𑌠	𑌡
AD 3 rd C	𑌕	𑌖	𑌗	𑌘	𑌙	𑌚	𑌛	𑌜	𑌝	𑌞	𑌟	𑌠	𑌡
AD 4 th C	𑌕	𑌖	𑌗	𑌘	𑌙	𑌚	𑌛	𑌜	𑌝	𑌞	𑌟	𑌠	𑌡
AD 5 th C	𑌕	𑌖	𑌗	𑌘	𑌙	𑌚	𑌛	𑌜	𑌝	𑌞	𑌟	𑌠	𑌡
AD 6 th C	𑌕	𑌖	𑌗	𑌘	𑌙	𑌚	𑌛	𑌜	𑌝	𑌞	𑌟	𑌠	𑌡
AD 7 th C	𑌕	𑌖	𑌗	𑌘	𑌙	𑌚	𑌛	𑌜	𑌝	𑌞	𑌟	𑌠	𑌡
AD 8 th C	𑌕	𑌖	𑌗	𑌘	𑌙	𑌚	𑌛	𑌜	𑌝	𑌞	𑌟	𑌠	𑌡
AD 9 th C	𑌕	𑌖	𑌗	𑌘	𑌙	𑌚	𑌛	𑌜	𑌝	𑌞	𑌟	𑌠	𑌡
AD 10 th C	𑌕	𑌖	𑌗	𑌘	𑌙	𑌚	𑌛	𑌜	𑌝	𑌞	𑌟	𑌠	𑌡
AD 11 th C	𑌕	𑌖	𑌗	𑌘	𑌙	𑌚	𑌛	𑌜	𑌝	𑌞	𑌟	𑌠	𑌡
AD 12 th C	𑌕	𑌖	𑌗	𑌘	𑌙	𑌚	𑌛	𑌜	𑌝	𑌞	𑌟	𑌠	𑌡
AD 13 th C	𑌕	𑌖	𑌗	𑌘	𑌙	𑌚	𑌛	𑌜	𑌝	𑌞	𑌟	𑌠	𑌡
AD 15 th C	𑌕	𑌖	𑌗	𑌘	𑌙	𑌚	𑌛	𑌜	𑌝	𑌞	𑌟	𑌠	𑌡
AD 16 th C	𑌕	𑌖	𑌗	𑌘	𑌙	𑌚	𑌛	𑌜	𑌝	𑌞	𑌟	𑌠	𑌡
AD 17 th C	𑌕	𑌖	𑌗	𑌘	𑌙	𑌚	𑌛	𑌜	𑌝	𑌞	𑌟	𑌠	𑌡
AD 18 th C	𑌕	𑌖	𑌗	𑌘	𑌙	𑌚	𑌛	𑌜	𑌝	𑌞	𑌟	𑌠	𑌡
AD 19 th C	𑌕	𑌖	𑌗	𑌘	𑌙	𑌚	𑌛	𑌜	𑌝	𑌞	𑌟	𑌠	𑌡

Numerals and other symbols - Arabic scripts

- Different shapes of some digits

Each numeral in the Persian variant has a different [Unicode](#) point even if it looks identical to the Eastern Arabic numeral counterpart. However, the variants used with [Urdu](#), [Sindhi](#), and other [Languages of South Asia](#) are not encoded separately from the Persian variants.

Western Arabic	0	1	2	3	4	5	6	7	8	9	10
Eastern Arabic ^[a]					٤	٥	٦	٧			
Persian ^[b]	۰	۱	۲	۳	۴	۵	۶		۸	۹	۱۰
Urdu ^[c]					۴	۵	۶				

Numerals and other symbols - the Tamil script

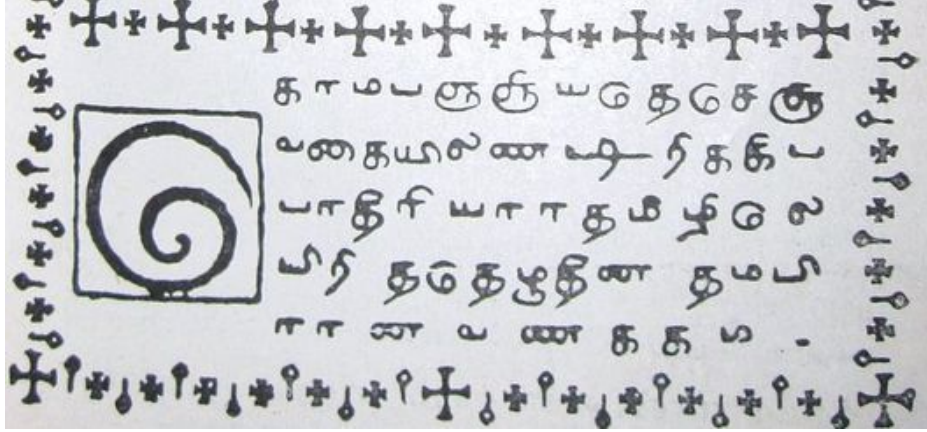
- Examples

Comparing Ninal Decimal and Tamil Decimal Counting

0	1	2	3	4	5	6	7	8	9	10	100	1000
	க	உ	ந	ச	ஐ	கூ	எ	அ	கை	ஓ	௩	௧௦௦
11	12	13	14	15	16	17	18	19	20			
கக	கஉ	கந	கச	கஐ	ககூ	கஎ	கஅ	ககை	கஓ			
21	22	23	24	25	26	27	28	29	30			
உக	உஉ	உந	உச	உஐ	உகூ	உஎ	உஅ	உகை	உஓ			
31	32	33	34	35	36	37	38	39	40			
நக	நஉ	நந	நச	நஐ	நகூ	நஎ	நஅ	நகை	நஓ			
								49	50			
								சகை	ஐஓ			
								59	60			
								நகை	கஓ			
								69	70			
								ககை	எஓ			

கூ	ந	$\frac{3}{8,400}$
கூ	ஹ	$\frac{1}{2,800}$
கூ	உ	$\frac{1}{8,200}$
கூ	பற	$\frac{1}{5,120}$
கூ	ப	$\frac{1}{6,400}$
கூ	சு	$\frac{3}{25,600}$
கூ	சு	$\frac{1}{12,800}$
கூ	பி	$\frac{1}{25,600}$
கூ	ந	$\frac{1}{5,120,000}$
கூ	ஹ	$\frac{1}{1,02,400}$

Writing system vary based on medium/technology



தமிழ், உலகில் உள்ள முதன்மையான மொழிகளில் ஒன்றும் செம்மொழியும் ஆகும். இந்தியா, இலங்கை, மலேசியா, சிங்கப்பூர் ஆகிய நாடுகளில் அதிக அளவிலும், ஐக்கிய அரபு அமீரகம், தென்னாப்பிரிக்கா, மொரிசியசு, பிசி, இரீயூனியன், திரினிடாடு போன்ற நாடுகளில் சிறிய அளவிலும் தமிழ் பேசப்படுகிறது.

- Written together **without spaces**, different shapes for the same character in different medium
- Consonants are written without *pulli* (a diacritic)

Character Encoding, Fonts, Input methods of South Asian Scripts

ASCII & Unicode

- **ASCII**

- American Standard Code for Information Interchange - 128 code points.
- 1 byte per character.

- **Unicode**

- Version 15.1 of the standard defines 149,813 code points and 161 scripts.
- <https://unicode.org/charts/>
- 3 (UTF-8) or 2 (UTF-16) bytes per Unicode (by default).

Fonts

- Understanding fonts is important for South Asian languages.
- More than just providing different shapes, some of the styles and ligatures are managed at the font rendering level.

Shakespeare	Naskh شیکسپیئر	Nastaliq شیکسپیئر
-------------	----------------	-------------------

Unicode - Ligatures

- **Multiple characters -> single ligature**

- By rendering - using fonts
- By using special Unicode characters

- **Ligatures by rendering characters**

- श्री (shri) - श (sh) + री (ri)

1st character	2nd character	3rd character	Result
س	ب	ب	سبب

- **Ligatures by using special Unicode characters**

ଠ + ଠ + ZWJ + ଠ --> ଠ

ଠ + ଠ + ଠ --> ଠଠ

<https://r12a.github.io/app-conversion/>

Unicode Character, Glyphs, and Ligatures

- Unicode Character

- U+0627 ا , U+0628 ب , U+062A ت , U+062B ث

beh-arab	beh-arab.isol	beh-arab.init	beh-arab.medi	beh-arab.fina
ب	ب	ب	ب	ب
teh-arab	teh-arab.isol	teh-arab.init	teh-arab.medi	teh-arab.fina
ت	ت	ت	ت	ت

- Glyphs

- Ligatures

Correct

لا ← ا + ل
U+0627 U+0644

Wrong

ا ← ا + ل ← ا + ل
Final shape Initial shape U+0627 U+0644

Challenges

Identification

- **How many characters in South Asian Languages?**
- **Language identification:**
 - Several languages are being written using a single script.
 - Single language is being written using multiple scripts.
 - Dialectal variations.

Confusing shapes

- Visually confusing graphemes/numbers/symbols:
 - கண் (*kan* - Eye) - கண் (*1n(?)* - ?)
 - அரசு - அரசு (*arasu* - Government) -
- Optical Character Recognition:
 - Character, Word, and Sentence segmentations.
 - Evolution of scripts.
 - Similar shapes.

	0B8	0B9	0BA	0BB	0BC	0BD	0BE	0BF
0		ஐ 0B90		ர 0BB0	ீ 0BC0	ூ 0BD0		ய 0BF0
1				ற 0BB1	ு 0BC1			ா 0BF1
2	ஃ 0B82	ஒ 0B92		ல 0BB2	ூ 0BC2			து 0BF2
3	ஃ 0B83	ஒ 0B93	ண 0BA3	ள 0BB3				உ 0BF3
4		ஒள 0B94	த 0BA4	ழ 0BB4				ம் 0BF4
5	அ 0B85	க 0B95		வ 0BB5				ஸ் 0BF5
6	ஆ 0B86			ய 0BB6	ெ 0BC6		ஃ 0BE6	பு 0BF6
7	இ 0B87			ஷ 0BB7	ே 0BC7	ள 0BD7	க 0BE7	ஸ் 0BF7
8	ஈ 0B88		ந 0BA8	ஸ 0BB8	ை 0BC8		உ 0BE8	ஷ 0BF8
9	உ 0B89	ங 0B99	ன 0BA9	ஹ 0BB9			ங 0BE9	நு 0BF9
A	ஊ 0B8A	ச 0B9A	ப 0BAA		ொ 0BCA		சு 0BEA	நீ 0BFA
B					ோ 0BCB		ரு 0BEB	
C		ஐ 0B9C			ெள 0BCC		கூ 0BEC	
D					ஃ 0BCD		எ 0BED	
E	எ 0B8E	ஞ 0B9E	ம 0BAE	ா 0BBE			அ 0BEF	
F	ஏ 0B8F	ட 0B9F	ய 0BAF	ரி 0BBF			கூ 0BEF	

Modifiers

- Short Vowels as diacritics - usually not transcribed resulting in ambiguity

- Non-joiner Characters - ambiguity for word boundaries

برابر	برابر	ابر	بر
برابر	brAbar	abar	bar
برابر	'equal/next'	'cloud'	'land'

Sorting

- Cannot sort just by using Unicode points

1.	க	(U+0B95) - KA
2.	ங	(U+0B99) - NGA
3.	ச	(U+0B9A) - CA
4.	ஞ	(U+0B9E) - NYA
5.	ட	(U+0B9F) - ṬA
6.	ண	(U+0BA3) - ṆA
7.	த	(U+0BA4) - TA
8.	ந	(U+0BA8) - NA
9.	ப	(U+0BAA) - PA
10.	ம	(U+0BAE) - MA
11.	ய	(U+0BAF) - YA
12.	ர	(U+0BB0) - RA
13.	ல	(U+0BB2) - LA
14.	வ	(U+0BB5) - VA
15.	ழ	(U+0BB4) - ḷA
16.	ள	(U+0BB3) - ḷA
17.	ற	(U+0BB1) - ṟA
18.	ள்	(U+0BA9) - ṇA

1.	(U+0627) - A
2.	ب (U+0628) - B
3.	پ (U+067E) - P
4.	ت (U+062A) - T
5.	ٹ (U+0679) - Ṭ
6.	ث (U+062B) - S
7.	ج (U+062C) - J
8.	چ (U+0686) - Ch
9.	ح (U+062D) - Ḥ
10.	خ (U+062E) - Kh
11.	د (U+062F) - D
12.	ڈ (U+0688) - Ḍ
13.	ذ (U+0630) - Dh
14.	ر (U+0631) - R
15.	ڑ (U+0691) - Ṛ
16.	ز (U+0632) - Z
17.	ژ (U+0698) - Zh
18.	...

Rendering

- **Application support for rendering characters**

- Special glyphs are not supported - in some applications these are considered as symbols during the processing.
- Rendering issues - Glyphs are not stored in the order that we see.
 - தே - த ே
 - <https://r12a.github.io/app-conversion/>
- Input methods - keyboard software.

- **Not all the letters in South Asian Languages are encoded**

- Cannot process old text

Encodings

- **Different Encodings:**
 - Standard Unicode
 - Other regional encoding (Tamil - TACE)
 - ASCII for South Asian Languages (Tamil - TAB/TAM)
- **Conversation challenges**
- **Consume more memory**
 - <https://onlinetools.com/unicode/count-unicode-characters>

Unicode normalisation

- Some characters can be typed in multiple ways
 - Example:
 - கொக்கு (*kokku* - Egret) vs கொக்கு
 - க ொ க் கு vs க ெ ா க் கு
- <https://r12a.github.io/app-conversion/>

Normalisation - more confusions

Display	Unicode Character Sequence					Transformation
	C ₁	C ₂	C ₃	C ₄	C ₅	
رئيس	<i>reh</i> U+0631	<i>yeh with hamza above</i> U+0626	<i>yeh</i> U+064A	<i>seen</i> U+0633		
رئيس	<i>reh</i> U+0631	<i>yeh</i> U+064A	<i>hamza above</i> U+0654	<i>yeh</i> U+064A	<i>seen</i> U+0633	Unicode NFC
رئيس	<i>reh</i> U+0631	<i>alef maksura</i> U+0649	<i>hamza above</i> U+0654	<i>yeh</i> U+064A	<i>seen</i> U+0633	Visual Normalization
رئيس	<i>reh</i> U+0631	<i>yeh with hamza above</i> U+0626	<i>farsi yeh</i> U+06CC	<i>seen</i> U+0633		Visual Normalization
رئيس	<i>reh</i> U+0631	<i>farsi yeh</i> U+06CC	<i>hamza above</i> U+0654	<i>yeh</i> U+064A	<i>seen</i> U+0633	Reading Normalization
رئيس	<i>reh</i> U+0631	<i>farsi yeh</i> U+06CC	<i>hamza above</i> U+0654	<i>farsi yeh</i> U+06CC	<i>seen</i> U+0633	Reading Normalization

Table 1: Six different spellings of the Arabic word for “leader” (MSA: /ra.ʔi:s/) rendered in Naskh.

<https://arxiv.org/ftp/arxiv/papers/2210/2210.12273.pdf>

Normalisation - more confusions

- Different unicode points for Arabic (U+0660 to U+0669) and other languages (U+06F0 to +06F9)

Each numeral in the Persian variant has a different [Unicode](#) point even if it looks identical to the Eastern Arabic numeral counterpart. However, the variants used with [Urdu](#), [Sindhi](#), and other [Languages of South Asia](#) are not encoded separately from the Persian variants.

Western Arabic	0	1	2	3	4	5	6	7	8	9	10
Eastern Arabic ^[a]					٤	٥	٦	٧			
Persian ^[b]	•	۱	۲	۳	۴	۵	۶		۸	۹	۱۰
Urdu ^[c]					۴	۵	۶				

Processing related challenges

- **Bidirectional Algorithm**

- Left to Right (LTR)
- Right to Left (RTL)
- Processing needs to be done in two level: Character and String levels

bahrain مصر kuwait

1 2 3

Fig. 3. The same directional runs in a LTR context. [See live demo](#)

kuwait مصر bahrain

3 2 1

Fig. 4. The same directional runs in a RTL context. [See live demo](#)

egypt

→ LTR

مصر

← RTL

Thank you

Kengatharaiyer Sarveswaran

University of Jaffna, Sri Lanka.

University of Konstanz, Germany

sarves@univ.jfn.ac.lk

sarves.github.io

Tafseer Ahmed

Senior NLP Scientist, Alexa Translations

tafseer@gmail.com