

Tamil

Juneki Hong

February 17, 2015

தமிழ்



According to Wikipedia

¹2007 edition of Nationalencyklopedin.

²With about 8 million as a second language.

According to Wikipedia

- ▶ Somewhere around 70-100 million native speakers¹².

¹2007 edition of Nationalencyklopedin.

²With about 8 million as a second language.

According to Wikipedia

- ▶ Somewhere around 70-100 million native speakers¹².
 - ▶ India: 60.8 million (2001 census)
 - ▶ Sri Lanka: 4.7 million

¹2007 edition of Nationalencyklopedin.

²With about 8 million as a second language.

According to Wikipedia

- ▶ Somewhere around 70-100 million native speakers¹².
 - ▶ India: 60.8 million (2001 census)
 - ▶ Sri Lanka: 4.7 million
- ▶ 1% of all speakers worldwide.
- ▶ 5th most spoken language in India

¹2007 edition of Nationalencyklopedin.

²With about 8 million as a second language.



According to Wikipedia

- ▶ Somewhere around 70-100 million native speakers¹².
 - ▶ India: 60.8 million (2001 census)
 - ▶ Sri Lanka: 4.7 million
- ▶ 1% of all speakers worldwide.
- ▶ 5th most spoken language in India
- ▶ A *classical* language that is over 2000 years old.

¹2007 edition of Nationalencyklopedin.

²With about 8 million as a second language.

Map!

Here is a Map

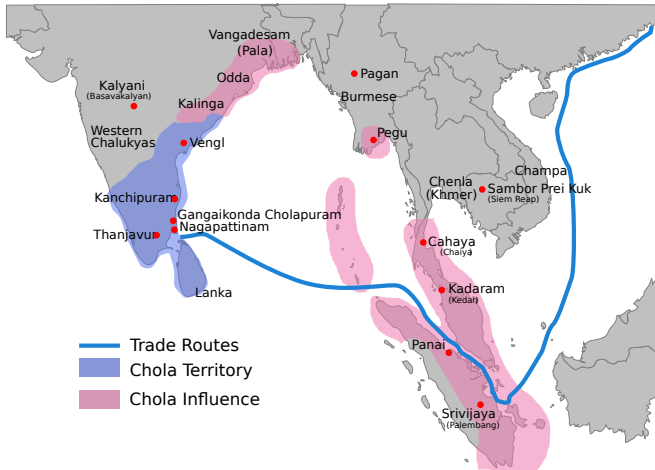


Figure : The Chola Dynasty: ~300BC - 1279CE

Sample Text³

மனிதப் பிறவியினர் சகலரும் சுதந்திரமாகவே பிறக்கின்றனர்; அவர்கள் மதிப்பிலும் உரிமைகளிலும் சமமானவர்கள். அவர்கள் நியாயத்தையும் மனசாட்சியையும் இயற்பண்பாகப் பெற்றவர்கள். அவர்கள் ஒருவருடனொருவர் சகோதர உணர்வுப் பாங்கில் நடந்துகொள்ளல் வேண்டும்.

Tamil Text

Maṇitap piṛaviyiṇar čakalarum čutantiramākavē piṛakkiṇṇar; avarkaḷ matippilum urimaikaḷilum čamamaṇavarkaḷ. Avarkaḷ niyāyattaiyum maṇačāṭčiyaiyum iyarpaṇpākap peṇṇavarkaḷ. Avarkaḷ oruvaruṇoruvar čakōtara uṇarvup pāṅkiḷ naṭantukoḷḷal vēṇṭum.

Transliterated Text

All human beings are born free and equal in dignity and rights. They are endowed with reason and conscience and should act towards one another in a spirit of brotherhood.

Translated Text

³www.omniglot.com/writing/tamil.htm



Tamil Script

- ▶ 12 vowels, 18 consonants, 1 special character.

Tamil Script

- ▶ 12 vowels, 18 consonants, 1 special character.
 - ▶ Script is syllabic, with 216 combinations of characters.
 - ▶ Total of 247 characters ($12 + 18 + 1 + (12 \cdot 18)$)

Tamil Script

- ▶ 12 vowels, 18 consonants, 1 special character.
 - ▶ Script is syllabic, with 216 combinations of characters.
 - ▶ Total of 247 characters (12 + 18 + 1 + (12*18))

Formation	Compound form	ISO 15919	IPA
க் + அ	க	ka	[kʌ]
க் + ஆ	கா	kā	[ka:]
க் + இ	கி	ki	[ki]

Syntax

- ▶ SOV
 - ▶ But with flexible word orders denoting different meanings.
- ▶ Null-Subject language.
 - ▶ A verb by itself can be a valid sentence.
 - ▶ Or the subject and object without a verb.

Parallel Corpora

- ▶ Corpus taken from `www.tamilnet.com` manually translated into English
 - ▶ 1,300 sentences
 - ▶ USC, 2001
- ▶ Joshua Decoder offers Indian parallel corpora including English-Tamil.
 - ▶ 36,252 sentences
 - ▶ JHU and University of Edinburgh, 2012
- ▶ EnTam: An English-Tamil Parallel Corpus taken from bible, cinema, and news domains.
 - ▶ 169,871 sentences
 - ▶ Charles University in Prague, 2014

MT Systems

	BLEU
Google (2011)	13.51
Post et al. 2012	9.81

JHU, University of Edinburgh ⁴

	BLEU _{suffix-sep} ⁵
Ramasamy et al. 2012	15.12

Charles University in Prague

⁴joshua-decoder.org/data/indian-parallel-corpora/

⁵Reference and Hypothesis translations suffix-separated before evaluation