

# The recent work on Internationalized Domain Names (IDNs) in South and East Asian Scripts and especially in Sinhala

Harsha Wijayawardhana B.Sc.,FBCS,CITP  
COO/CTO Theekshana R & D,  
Director, BoC and LK Domain Registry



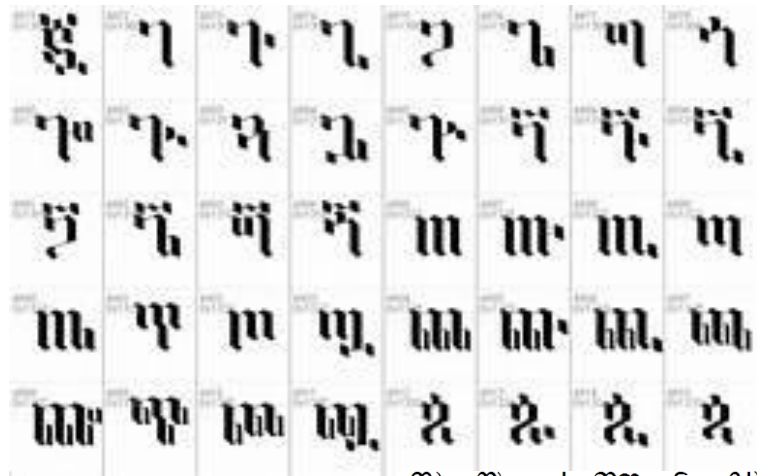
# About Sinhala

- Sinhala belongs to Indo-European Language family and sub family of Indo- Arya.
- Closest Relatives of Sinhala Language are from North Indian family such as Hindi, Urdu, Bengali, Orya, etc.
- Europe has several of Indo-European Sub-families such as German, Romanic, and East Slavic.
- English, German and Swedish belong to German Sub Family whereas French, Portuguese, Spanish and Italian belong to Romanic Languages.

සිංහල

# Written Script

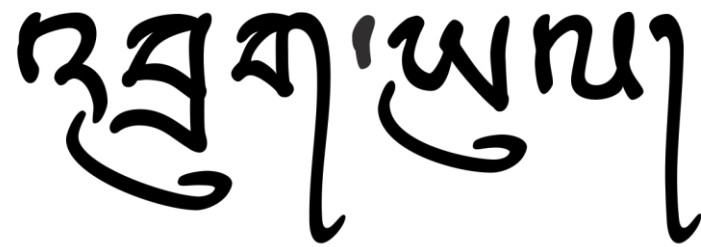
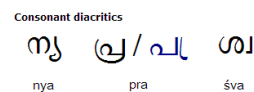
- Modern Indian and East Asian scripts derive from Brahmi Script and they belong to Abugida Scripts



## Ethiopic Script



## Malayalam



# DZonka

ก ข ค ฅ ง จ ฉ ช  
ช ฌ ญ ฎ ฏ ฐ ฑ ฒ  
ณ ด ต ถ ท ธ น บ ป  
ผ ฝ พ ฟ ภ ม ย ร ล  
ว ศ ษ ส ห ฬ อ ฮ

## Thai

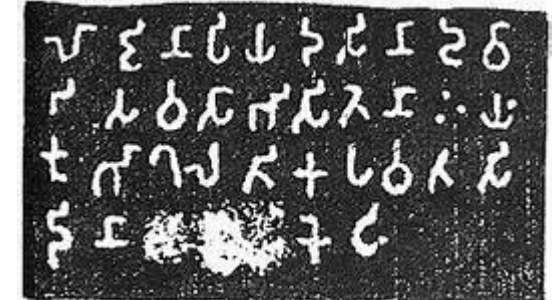
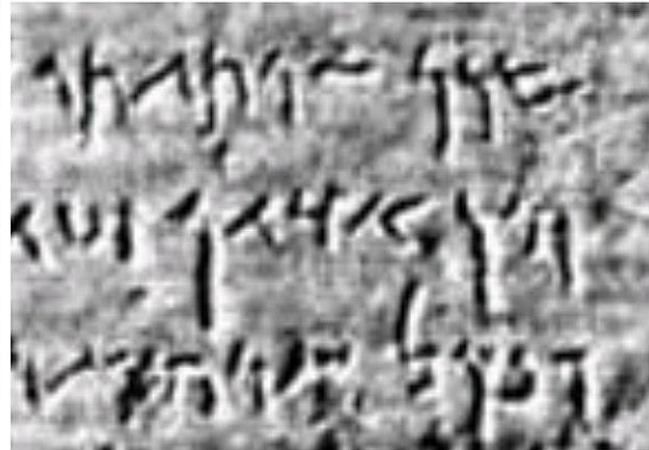
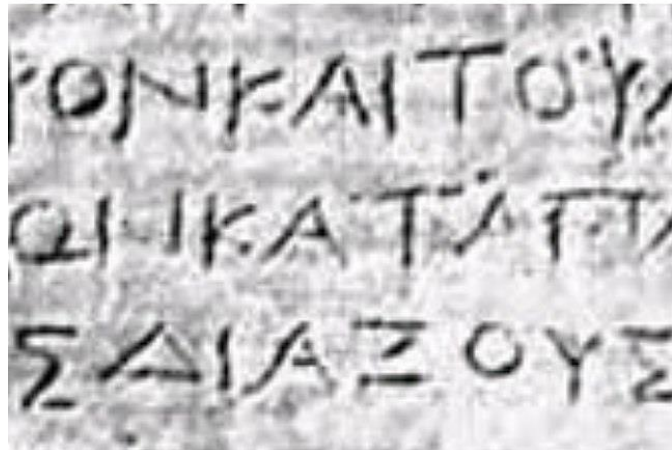
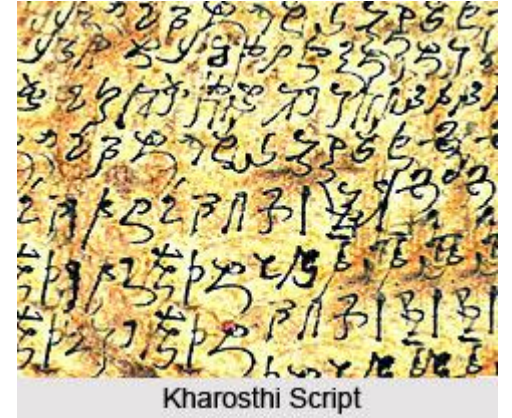
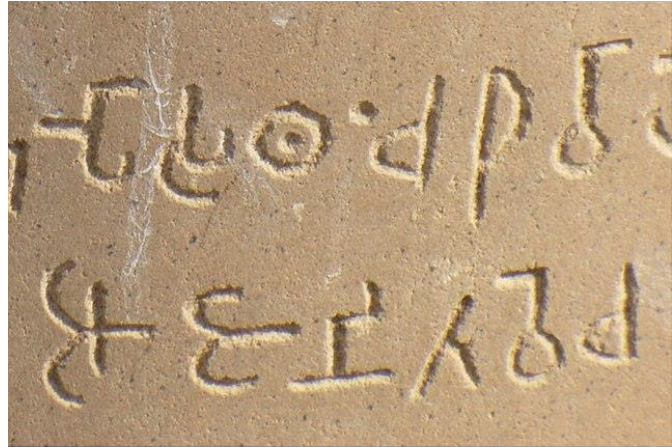
[illegible]

# Brahmi

- All these Indian, East Asian Scripts derived from Brahmi.
- Although it was once believed that Emperor Asoka used Brahmi as his regal script, it has discovered that oldest Brahmi is found in Sri Lanka and Southern India.
- Brahmi is usually written from right to left. In Sri Lanka, eighteen (18) Rock inscriptions have discovered with Brahmi written from right to left dating to 3<sup>rd</sup> century BC.
- In North Eastern part of Asoka's Empire, his edicts were written in Prakrit using Kharoshthi Script. Kharoshthi was written from right to left.

# Brahmi cont.

Brahmi  
and  
Kharoshthi





# Brahmi Script cont.

- It was earlier postulated that Brahmi derived from Aramaic.
- Another hypothesis is Brahmi derived from Indus Valley script.



# Sinhala and other Indic Scripts

- Sinhala Script derived from Southern Brahmi Script.
- Subsequently, Sinhala Script was influenced by Pallava Grantha.
- Eastern Scripts such as Khmer, Myanmar, etc. influenced by Pallava Script.
- Sinhala has two Vowels අ and ඞ.
- These two vowels are found in

East Asian Scripts.

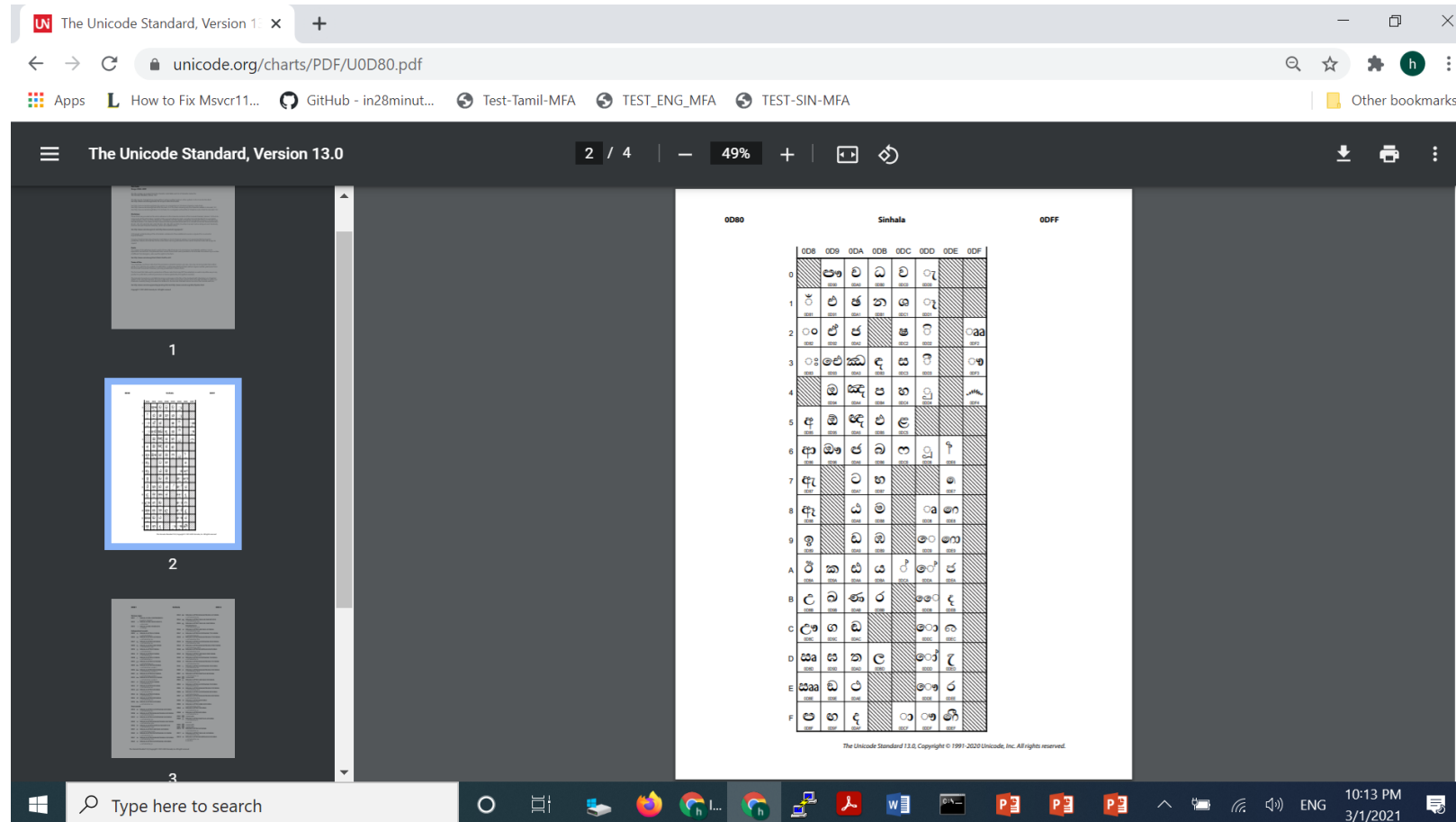


# Sinhala Computerization

- Various attempts on Sinhala Computerization begun in early eighties.
- Sinhala was not encoded until 1996.
- Sinhala Unicode Implementation began in early two thousand.,
- Sri Lanka came up with input standard SLS 1134 and in 2010, Sinhala Unicode was updated for the last time with the third revision by including Sinhala Numerals.
- Sri Lanka also has its own standard for Tamil input since Sri Lanka Tamil community prefers to use input method deviating from Indian Implementation: SLS 1326 of 2008.



# Unicode Version 13.1



# Sinhala Archaic Numbers in Supplementary Multilingual Plane (SMP)

ISO/IEC 10646:2011 FDIS

Not secure | unicode.org/L2/L2011/11220-n3968-pdam10-all.pdf

Apps | How to Fix Msvcrt11... | GitHub - in28minut... | Test-Tamil-MFA | TEST\_ENG\_MFA | TEST-SIN-MFA | Other bookmarks

ISO/IEC 10646:2011 FDIS 19 / 33 50%

18

19

20

Proposed Draft Amendment (PDAM) 1  
111E0

ISO/IEC 10646:2012/Amd.1: 2012 (E)  
Sinhala Archaic Numbers  
111F0

This number system is also known as Sinhala Illakam. This number system does not have a zero place holder concept unlike the Sinhala categorical numbers, Sinhala Lith Illakam, encoded in the U+0040-00FF range.

**Historical digits**

These digits are not used with a zero

111E1 𑆀 SINHALA ARCHAIC DIGIT ONE  
111E2 𑆁 SINHALA ARCHAIC DIGIT TWO  
111E3 𑆂 SINHALA ARCHAIC DIGIT THREE  
111E4 𑆃 SINHALA ARCHAIC DIGIT FOUR  
111E5 𑆄 SINHALA ARCHAIC DIGIT FIVE  
111E6 𑆅 SINHALA ARCHAIC DIGIT SIX  
111E7 𑆆 SINHALA ARCHAIC DIGIT SEVEN  
111E8 𑆇 SINHALA ARCHAIC DIGIT EIGHT  
111E9 𑆈 SINHALA ARCHAIC DIGIT NINE

**Historical numbers**

111EA 𑆉 SINHALA ARCHAIC NUMBER TEN  
111EB 𑆊 SINHALA ARCHAIC NUMBER TWENTY  
111EC 𑆋 SINHALA ARCHAIC NUMBER THIRTY  
111ED 𑆌 SINHALA ARCHAIC NUMBER FORTY  
111EE 𑆍 SINHALA ARCHAIC NUMBER FIFTY  
111EF 𑆎 SINHALA ARCHAIC NUMBER SIXTY  
111F0 𑆏 SINHALA ARCHAIC NUMBER SEVENTY  
111F1 𑆐 SINHALA ARCHAIC NUMBER EIGHTY  
111F2 𑆑 SINHALA ARCHAIC NUMBER NINETY  
111F3 𑆒 SINHALA ARCHAIC NUMBER ONE HUNDRED  
111F4 𑆓 SINHALA ARCHAIC NUMBER ONE THOUSAND

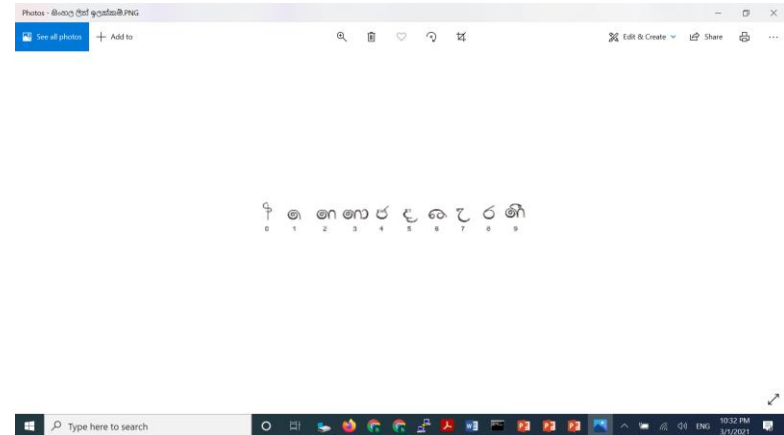
© ISO/IEC 2011 - All rights reserved

Type here to search

10:17 PM  
3/1/2021

# Sinhala Numerals

- Although five ways of writing numbers found, only two have been encoded.
- The Sinhala Illakkam and Lith Illakkam have evolved from Brahimi numerals.
- Lith Illakkam has a zero but not the Sinhala Illakkam.



# IDNs

- Sri Lanka began working in Internationalized Domain Names early Two Thousand becoming very first country to come up with two Non Latin Script ccTLDs: *.ලංකා* and *.இலங்கை*
- Sri Lanka further worked on Second Level for ccTLDs.
- Language Technology Research Lab (LTRL) of University of Colombo School of Computing (UCSC) under Dr. Ruvan Weerasinghe and University Moratuwa under the guidance of Prof. Gihan Dias carried out research into IDNs.
- One of the earliest research papers on IDNs by LTRL:
- [https://www.researchgate.net/profile/Chamila-Liyanage-2/publication/228715570\\_Implementation\\_of\\_Internet\\_Domain\\_Names\\_in\\_Sinhala/links/0deec52148a10198fc000000/Implementation-of-Internet-Domain-Names-in-Sinhala.pdf](https://www.researchgate.net/profile/Chamila-Liyanage-2/publication/228715570_Implementation_of_Internet_Domain_Names_in_Sinhala/links/0deec52148a10198fc000000/Implementation-of-Internet-Domain-Names-in-Sinhala.pdf)

# Sinhala Generation Panel of ICANN

- ICANN formed the Sinhala Generation Panel in December 2017. Indian Generation Panels met in Sri Lanka while Sri Lanka formed its panel and began in earnest. the public comment. At present, we have released rules for second level.
- By mid 2018, Sri Lanka was ready to release first- draft on rules and by August, we managed to finish most of the work and released for



# Sinhala GP – Languages Using Sinhala Script

- Languages covered by the script:
  - Sinhala
  - Pali
  - Sanskrit
- Starting from (Maximal Starting Repertoire)MSR-3, the repertoire includes:
  - 72 code points
  - 4 sequences

# Code Point Repertoire

- ☉ Starting from MSR-3, the repertoire includes:
  - 72 code points
  - 4 sequences
- ☉ The repertoire excludes:

#	Unicode Code Point	Glyph	Character Name	Reason for exclusion
1	0D8E	සෲ	SINHALA LETTER IRUUYANNA	Usage unknown
2	0D8F	ඌ	SINHALA LETTER ILUYANNA	Usage unknown
3	0D90	ඌඹ	SINHALA LETTER ILUUYANNA	Usage unknown
4	0D9E	ඳ	SINHALA LETTER KANTAJA	Not in modern usage
5	0DA6	ඳ	SINHALA LETTER SANYAKA	Only used in the word ‘ඳුල්ල’ (this word is used to call dogs)
6	0DDF	ඹ	SINHALA VOWEL SIGN GAYANUKITTA	Usage unknown
7	0DF3	ඹ	SINHALA VOWEL SIGN DIGA GAYANUKITTA	Usage unknown

# In-Script Variant Analysis

- ◎ Sinhala GP decided the following are in-script variant code points due to similar shapes and characters which could be used interchangeably
  - ස (U+0DC3) and ස (U+0D9D)
  - බ (U+0DB6) and බ (U+0D9B)
  - හ (U+0DC4) and හ (U+0DB7)
  - ච (U+0DA0) and ච (U+0DC0)
  - ඔ (U+0D94) and ඔ (U+0DB9)
  - එ (U+0D91) and එ (U+0DB5)
  - සෘ (U+0D8D) and සෘ (U+0DC3 U+0DD8)
  - ඔඑ (U+0D93) and ඔඑ (U+0DB5 U+0DD9)
  - ඒ (U+0D92) and ඒ (U+0DB5 U+0DCA)
  - ඕ (U+0D95) and ඕ (U+0DB9 U+0DCA)

# Cross-Script Variant Analysis

- ⦿ Sinhala GP concluded there is no cross-script variant rules
- ⦿ Following are confusable cases
  - U+0D82 (SINHALA SIGN ANUSVARAYA, ඌ)

Sinhala	Telugu	Kannada	Malayalam
ඌ (U+0D82)	ౠౠ (U+0C02)	ౠೠ (U+0C82)	ౠౠ (U+0D02)

- U+0D83 (SINHALA SIGN VISARGAYA, ඌಃ)

Sinhala	Devanagari	Gujarati	Telugu	Kannada	Malayalam
ඌಃ (U+0D83)	ౠഃ (U+0903)	ౠಃ (U+0A83)	ౠః (U+0C03)	ౠః (U+0C83)	ౠഃ (U+0D03)

# Cross-Script Variant Analysis

◎ Following are confusable cases (cont.)

○ Sinhala and Malayalam

Sinhala	Malayalam
ඹ (U+0D9C)	൩ (U+0D17)
ඹ (U+0DC1)	൪ (U+0D36)
ඹ (U+0DCF)	൫ (U+0D3E)

○ Sinhala and Myanmar

Sinhala	Myanmar
ඹ (U+0D9C)	၂ (U+1010)
ඹ (U+0DC1)	၃ (U+107B)



# Whole Label Evaluation Rules

- ⊙ Code point category
  - C → Consonant
  - V → Vowel
  - M → Matras / Vowel Signs
  - B → Anusvara (Bindu)
  - X → Visarga
  - H → Halanta / Virama
  - J → Sannjakas
- ⊙ Whole Label Evaluation Rules
  - H: must be preceded by C
  - M: must be preceded by C or J
  - X: must be preceded by either V, C, or M
  - B: must be preceded by either V, C, J or M

# Some of the Issues -

- ⦿ IDNA 2003 did not support hidden Characters such ZWJ, ZWNJ, etc.
- ⦿ Without those characters, in Sinhala and some other Indic scripts had issue of rendering some of commonly used characters such as ශ්‍රී without country name cannot be written. ZWJ is needed for Rakar, and Yansa forms to render.
- ⦿ IDNA 2008 now support the above and puny code allows these characters. And they allowed now in the second level.

Thank you

Questions ?