

A UNICODE PRIMER

Neo-Brāhmī Generation Panel
16 July 2015

Agenda

1. From Stone to Chip
2. Fundamentals
3. Why Unicode: Need for Standards
4. What is Unicode ?
5. Basic Features of Unicode

FROM STONE TO CHIP



FROM STONE TO CHIP

Data Preservation has a fascinating history:

Text inscribed on Stones, Palm Leaves, Copper-plates, Clay Tablets gave way to Wooden Block printing by the 9th Century and by the 15th Century with Gutenberg's moveable metal type printing press created the Gutenberg Galaxy which has survived for over 5 centuries and is slowly being replaced by the Digital Galaxy.

Fundamentals:

Unlike the earlier means for data preservation, storage on chips has an important difference. A book can be replicated and shared and will be legible all those who share it. Sharing of electronic data which can ensure that anybody can access the data posed a major issue.

गाजियाबाद। मुख्य न्यायिक मजिस्ट्रेट की अदालत ने कथित रूप से तिरंगे के 'अपमान' के लिए अभिनेता अमिताभ बच्चन और अभिषेक बच्चन के खिलाफ मामला दर्ज कराने वाले व्यक्ति का सोमवार को यहां बयान दर्ज किया।

Fundamentals: Issues with Digital data

Before Unicode was invented, there were hundreds of different encoding systems for assigning these numbers. No single encoding could contain enough characters: for example, the European Union alone requires several different encodings to cover all its languages. Even for a single language like English no single encoding was adequate for all the letters, punctuation, and technical symbols in common use.

These encoding systems also conflict with one another. That is, two encodings can use the same number for two different characters, or use different numbers for the same character.

To resolve this confusion - in the late eighties, it was decided to have a single unified standard called Unicode. The aim of Unicode is to assign a unique number to each and every character of each and every "reasonable" writing system used on the earth, "no matter what the platform, no matter what the program, not matter what the language".

This ensure that data entered with a keyboard using Unicode compliant fonts is shared seamlessly with another machine which may or may not necessarily have the font.

Fundamentals: Issues with Digital data

This involves standards. Adherence to standards ensures compatibility, safeguarded data, avoids vendor locking, proper exchange of data between various systems, applications, databases, devices, etc.

Standards exist for Input mechanisms: Keyboards, Display and Storage. GIST has been at the forefront in development of standards for all three.

However this presentation will focus only on Unicode and its relevance to Internationalised Domain Names.

UNICODE

- Storage standard
- Unicode is for all languages of the world: *the world speaks Unicode.*
- Enables seamless exchange of data – desktops, printers, databases, browsers, devices.



UNICODE

- *The Unicode consortium defines Unicode as :*
- *“Unicode is the universal character encoding, maintained by the Unicode consortium. This encoding standard provides the basis for processing, storage and interchange of text data in any language in all modern software and information technology protocols.”*
- *It is the superset of all the languages in the world which also includes punctuation, special characters (shapes), currency symbols, mathematical symbols etc. Using Unicode, more than 65000 different characters can be represented. Unicode comprises of many code charts.*
- *The Unicode code charts can be referred at:
<http://www.unicode.org/charts>*

Unicode Features

In what follows using the Devanagari code-block, basic features of Unicode required for an understanding of IDN will be presented:

1. Unicode's main source is ISO 10646 which maintains the world's scripts. However Unicode is also Community Based and receives proposals from both members as well as individuals. Each proposal has to be substantiated with concrete evidence for introducing a given character. The proposals are put up for comments and eventually accepted
2. Unicode is script based. Each script is stored in "Blocks". Each code-block bears the name of the Script to which it refers. Unicode has a total of 262 code blocks of which 160 are in the active plane termed as BMP or **Basic Multilingual Plane** and the remainder refer to obsolete scripts which are no longer in use are stored in the **Plane 1**, the **Supplementary Multilingual Plane (SMP)**, which contains historic scripts such as Linear B, Egyptian hieroglyphs, and cuneiform scripts.

Unicode Features

3. Thus Devanagari script is stored in a code-block (often termed also as code-page or code-chart).

Devanagari

	090	091	092	093	094	095	096	097
0	ै	ऐ	ठ	र	ी	ॐ	ऋ	ॠ
1	ँ	ऑ	ड	र	ु	ं	ऌ	ॡ
2	ं	ओ	ढ	ल	ॡ	ॢ	ॣ	अँ
3	ः	ओ	ण	ळ	ॣ	।	॥	अँ
4	अे	औ	त	ळ	ॣ	।	॥	औ
5	अ	क	थ	व	ँ	ँ	॥	औ
6	आ	ख	द	श	े	ॢ	ॣ	अ
7	इ	ग	ध	प	े	ॢ	ॣ	अ
8	ई	घ	न	स	ै	ॣ	ॣ	र
9	उ	ङ	त	ह	ँ	ख	३	ज
A	ऊ	च	प	ं	ो	ग	४	य
B	ऋ	ॡ	फ	ी	ो	ज	५	ग
C	ऌ	ॢ	व	ॣ	ौ	ड	६	ज
D	ँ	झ	भ	।	॥	ढ	७	२
E	ऐ	अ	म	ा	ि	फ	८	ड
F	ए	ट	य	ि	ौ	य	९	ब

Unicode Features

4. Each block is made up of **code-points**. Each code-point refers to a shape i.e. a glyph. The number of code points in a Unicode block is a multiple of 16. Unicode blocks range in size from the minimum of 16 to a maximum of 65,536 code points. A block may contain unassigned code points, which are reserved. In case a code-block gets full, an extended code-block is created. Thus Devanagari also admits an extended Devanagari block.
5. The totality of code-points within a code-block is termed as the **repertoire** of that particular script. All the code-points pertaining to the character of Devanagari in the Devanagari code block constitute its repertoire.

Unicode Features

6. The code-points within a code-block are pertinent to that script and that script alone. No mixing of scripts is permitted. If a given code-point i.e. a character glyph is shared by many scripts, Unicode provides instructions to that effect. Thus the danda and double danda of Devanagari are used by other scripts as in the case of the code-block Bengali and Assamese which covers 0900 -097F

Reserved

For viram punctuation, use the generic Indic 0964 and 0965.

Note that these punctuation marks are referred to as dahri and double dahri in Bangla.

09E4  <reserved>

→ 0964 I devanagari danda

09E5  <reserved>

→ 0965 II devanagari double danda

Unicode Features

7. The code-points cover not only the characters but all and every element including diacritics and even weights and measures as in the case of Tamil.

0900 Devanagari 094F

Various signs

0900	ॐ	DEVANAGARI SIGN INVERTED CANDRABINDU
		= vaidika adhomukha candrabindu
0901	ं	DEVANAGARI SIGN CANDRABINDU
		= anusika
		→ 0310 ॐ combining candrabindu
0902	ँ	DEVANAGARI SIGN ANUSVARA
		= bindu
0903	ॐ	DEVANAGARI SIGN VISARGA

Independent vowels

0904	अ	DEVANAGARI LETTER SHORT A
0905	आ	DEVANAGARI LETTER A
0906	इ	DEVANAGARI LETTER AA
0907	ई	DEVANAGARI LETTER I
0908	उ	DEVANAGARI LETTER II
0909	ऊ	DEVANAGARI LETTER U
090A	ऋ	DEVANAGARI LETTER UU
090B	ॠ	DEVANAGARI LETTER VOCALIC R
090C	ॡ	DEVANAGARI LETTER VOCALIC L
090D	ए	DEVANAGARI LETTER CANDRA E
090E	ऐ	DEVANAGARI LETTER SHORT E

		• Kashmiri, Bihari languages
		• also used for transcribing Dravidian short e
090F	ऋ	DEVANAGARI LETTER E
0910	ॠ	DEVANAGARI LETTER AI
0911	औ	DEVANAGARI LETTER CANDRA O
0912	ॐ	DEVANAGARI LETTER SHORT O
		• Kashmiri, Bihari languages
		• also used for transcribing Dravidian short o

0913	ओ	DEVANAGARI LETTER O
0914	औ	DEVANAGARI LETTER AU

Consonants

0915	क	DEVANAGARI LETTER KA
0916	ख	DEVANAGARI LETTER KHA
0917	ग	DEVANAGARI LETTER GA
0918	घ	DEVANAGARI LETTER GHA
0919	ङ	DEVANAGARI LETTER NGA
091A	च	DEVANAGARI LETTER CA
091B	छ	DEVANAGARI LETTER CHA
091C	ज	DEVANAGARI LETTER JA
091D	झ	DEVANAGARI LETTER JHA
091E	ञ	DEVANAGARI LETTER NYA
091F	ट	DEVANAGARI LETTER TTA
0920	ठ	DEVANAGARI LETTER TTHA
0921	ड	DEVANAGARI LETTER DDA
0922	ढ	DEVANAGARI LETTER DDHA
0923	ण	DEVANAGARI LETTER NNA
0924	त	DEVANAGARI LETTER TA
0925	थ	DEVANAGARI LETTER THA
0926	द	DEVANAGARI LETTER DA
0927	ध	DEVANAGARI LETTER DHA
0928	न	DEVANAGARI LETTER NA
0929	न्	DEVANAGARI LETTER NNNA

		• for transcribing Dravidian alveolar n
		≡ 0928 ण 093C ॢ
092A	प	DEVANAGARI LETTER PA
092B	फ	DEVANAGARI LETTER PHA
092C	ब	DEVANAGARI LETTER BA
092D	भ	DEVANAGARI LETTER BHA
092E	म	DEVANAGARI LETTER MA
092F	य	DEVANAGARI LETTER YA

0930	र	DEVANAGARI LETTER RA
0931	ॠ	DEVANAGARI LETTER RRA

		• for transcribing Dravidian alveolar r
		• half form is represented as "Eyelash RA"
		≡ 0930 र 093C ॢ

0932	ल	DEVANAGARI LETTER LA
0933	ळ	DEVANAGARI LETTER LLA
0934	ॢ	DEVANAGARI LETTER LLLA

		• for transcribing Dravidian l
		≡ 0933 ळ 093C ॢ

0935	व	DEVANAGARI LETTER VA
0936	श	DEVANAGARI LETTER SHA
0937	ष	DEVANAGARI LETTER SSA
0938	स	DEVANAGARI LETTER SA
0939	ह	DEVANAGARI LETTER HA

Dependent vowel signs

These dependent vowel signs are used in Kashmiri and in the Bihari languages (Bhojpuri, Magadhi, and Maithili).

093A	ॣ	DEVANAGARI VOWEL SIGN OE
093B	।	DEVANAGARI VOWEL SIGN OOE

Various signs

093C	ॣ	DEVANAGARI SIGN NUKTA
		• for extending the alphabet to new letters
093D	।	DEVANAGARI SIGN AVAGRAHA

Dependent vowel signs

093E	॥	DEVANAGARI VOWEL SIGN AA
093F	॥	DEVANAGARI VOWEL SIGN I
		• stands to the left of the consonant

0940	॥	DEVANAGARI VOWEL SIGN II
0941	॥	DEVANAGARI VOWEL SIGN U
0942	॥	DEVANAGARI VOWEL SIGN UU
0943	॥	DEVANAGARI VOWEL SIGN VOCALIC R
0944	॥	DEVANAGARI VOWEL SIGN VOCALIC RR
0945	॥	DEVANAGARI VOWEL SIGN CANDRA E

0946	॥	DEVANAGARI VOWEL SIGN SHORT E
		• Kashmiri, Bihari languages
		• also used for transcribing Dravidian short e

0947	॥	DEVANAGARI VOWEL SIGN E
0948	॥	DEVANAGARI VOWEL SIGN AI
0949	॥	DEVANAGARI VOWEL SIGN CANDRA O
094A	॥	DEVANAGARI VOWEL SIGN SHORT O

		• Kashmiri, Bihari languages
		• also used for transcribing Dravidian short o
094B	॥	DEVANAGARI VOWEL SIGN O
094C	॥	DEVANAGARI VOWEL SIGN AU

Virama

094D	॥	DEVANAGARI SIGN VIRAMA
		= halant (the preferred Hindi name)
		• suppresses inherent vowel

Dependent vowel signs

094E	॥	DEVANAGARI VOWEL SIGN PRISHTHAMATRA E
		• character has historic use only
		• combines with E to form AI, with AA to form O, and with O to form AU

094F	॥	DEVANAGARI VOWEL SIGN AW
		• Kashmiri, Bihari languages

Tamil numerics

0BF0	௮	TAMIL NUMBER TEN
0BF1	௮	TAMIL NUMBER ONE HUNDRED
0BF2	௮	TAMIL NUMBER ONE THOUSAND

Tamil calendrical symbols

0BF3	௮	TAMIL DAY SIGN
		= naal, naali/padi
		• denotes a measure of grain that equals 2 uri or 4 ulakku
		= pillaiyaar suli
		• denotes auspiciousness

0BF4	௮	TAMIL MONTH SIGN
		= maatham

0BF5	௮	TAMIL YEAR SIGN
		= varudam

Tamil clerical symbols

0BF6	௮	TAMIL DEBIT SIGN
		= patru

0BF7	௮	TAMIL CREDIT SIGN
		= eduppu

		• denotes incoming cash which is set aside for unknown expenses
		• sometimes used as the credit sign; the traditional credit sign is distinct

0BF8	௮	TAMIL AS ABOVE SIGN
		= merpadi

Currency symbol

0BF9	௮	TAMIL RUPEE SIGN
		= rupai

Tamil symbol

0BFA	௮	TAMIL NUMBER SIGN
		= enn, niluvai
		• denotes balance
		• sometimes used as the number sign; the traditional number sign is distinct

Unicode Features

8. Each code-point or glyph has a name termed as **Unicode Name** which and also a **unique number** which refers to its position in the code-block. This is normally shown as U+0000 .

Thus क has the following value:

क U+0915

DEVANAGARI LETTER KA

9. In addition each code-point has **properties** which are not defined in the code-chart/code-block. These properties are crucial for Internationalised Domain Names since they define whether a given character will be accepted or not in resolving an IDN typed in a browser. This means that not all characters can be used to create an address. Properties such as PValid or Protocol valid are assigned to each code point [RFC 5987]

Typing डबांग०.भारत

will result in an error since the abbreviation marker ० is not permissible and will not be resolved.

Unicode Features: PRINCIPLES

Unicode is governed by 3 major principles. These are

1. The Principle of Economy : Only one and one glyph afferent to a particular script will be maintained and stored as a code-point. Normally this is the most frequent one. In other words allographs are not permitted. Thus the older forms of Devanagari /a/ /jh/ are not retained. However a font developer if he so desires can replace the /jh/ by the older shape the older shape in his font and Unicode will render it since both will refer to the same code-point.

अ = आ

झ = झ

Unicode Features: Principles

2. The principle of Simplicity: This principle states that only atomic glyphs will be retained. These are characters which are not decomposable. This is important in the case of complex scripts of the Neo-Brahmi family since these use the Halanta or the Vowel killer (with the possible exception of Tamil) to render composite characters known as ligatures. Unicode only permits akhanda forms and composing complex conjuncts is rendered possible by the adjunction of akhanda code-points

क्ष = क+्+ष

कात्स्न्य = क+ा+र+्+त+्+स+न्+य

Unicode Features

3. The principle of Comprehensiveness:

This principle simply states that a given repertoire of a script represented in Unicode shall be as comprehensive as possible, and this without compromising the other two principles.

Thus, as per the principle of Economy, allographs of a given character shall not be included.

Thus ल ल variants of a single shape shall not be admitted. Only one form shall be permitted as per the first principle.

Unicode Features: Normalization

An exception to the above principle arises when a given character within a language can be rendered in two different ways. Unicode defines this as removing alternate representations of equivalent sequences from textual data, to convert the data into a unique form. In the Unicode Standard, normalization refers specifically to processing to ensure that strings have unique representations.

रिज़र्व

ज + ्र = 091C+093C

ज़ = 095B

रिज़र्व = र + ि + ज + . + र + + व

Also = र + ि + ज़ + र + + व

This ensures in the case of IDN's that phishing and spoofing are reduced as far as possible, since the browser itself "normalizes" the two forms.

Unicode Features: Normalization

Groundwork for IDN. Points covered:

1. Unicode
2. ISO 10646
3. BMP: Basic Multilingual Plane/SMP: Supplementary Multilingual Plane
4. Code Charts/Code Blocks
5. Code point
6. Character Set/Character Repertoire/Allographs
7. Reading a Unicode Code-Block
8. Character Properties
9. Unicode Normalization
10. Governing Principles: Simplicity; Comprehensiveness, Economy

Further reading:

Pvalid code-points: <https://tools.ietf.org/html/rfc5892>

•Glossary of Unicode terms: <http://www.unicode.org/glossary/>

•Properties of South Asian Scripts:

<http://www.unicode.org/versions/Unicode7.0.0/ch12.pdf>

धन्यवाद !