

# **COD891: Minor Project**

Transliteration amongst Indian  
languages

Nikhil Chaturvedi  
2013CS50291

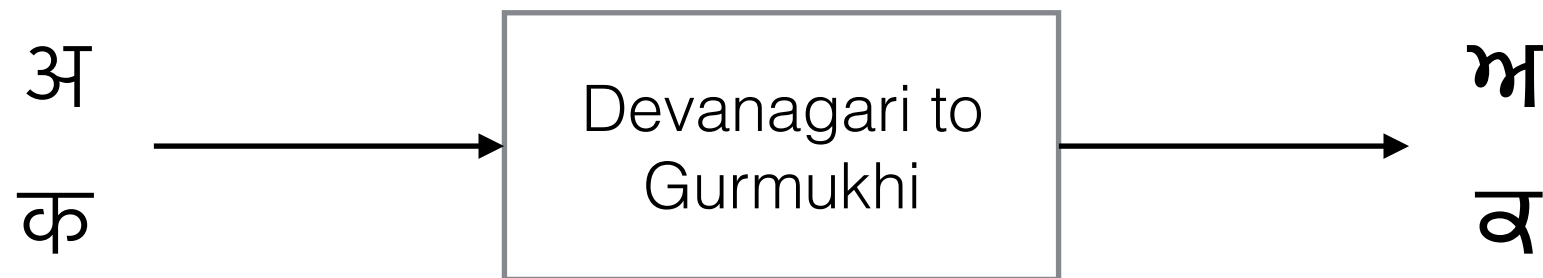
Supervisor: Prof. Rahul Garg

# Written Text

- Any computerised written text has 3 characteristics:
  - **Language:** Signifies the semantical meaning of words. Eg. Hindi, Punjabi, Sanskrit etc.
  - **Script:** Signifies the symbolic representation used for depicting the alphabet of the language. Eg. Devanagari, Gurmukhi, SLP1 etc.
  - **Encoding:** The binary encoding schema through which the symbols of the script are stored or presented. Eg. ASCII, UTF-8, UTF-16 etc.

# Transliteration

- Transliteration is a type of conversion of a text from one script to another that involves swapping letters (thus trans- + liter-) in predictable ways.
- It makes use of the alphabet correspondence across languages of similar origins.



# Motivation

- The Brahmi-derived writing systems of Indian languages are mostly rather similar in structure, but have different letter shapes.
- These scripts are based on similar phonetic values which allows for easy transliteration.
- The phonetic sound [ki] will be rendered as कि in Devanagari, as ਕਿ in Gurmukhi, and as கி in Tamil. Each having different code-points in Unicode and ISCII.
- This correspondence forms the basis of the motivation behind deriving a uniform encoding schema that is based on the underlying phonetic value rather than the symbolic representation.

# Existing Work: Unicode

- Unicode has designated code blocks for almost all major Indian scripts.
- The supported scripts are: Assamese, Bengali (Bangla), Devanagari, Gujarati, Gurmukhi, Kannada, Malayalam, Oriya, Tamil, and Telugu among others.
- Across scripts, Unicode respects alphabet correspondence and letters with similar phonetic values are assigned the same codepoints.
- As a result, transliteration can be done easily with a mere offsetting of the alphabet code. For example, अ is U+0905 while ऋ is U+0A05. क is U+0915 while ख is U+0A15.

# Devanagari and Gurmukhi

## Codeblocks

	090	091	092	093	094	095	096	097
0	ॐ 0900	ऐ 0910	ठ 0920	र 0930	ी 0940	ॐ 0950	ॠ 0960	ॐ 0970
1	ँ 0901	ऑ 0911	ड 0921	ॠ 0931	ु 0941	ं 0951	ॡ 0961	ं 0971
2	ं 0902	ओ 0912	ढ 0922	ल 0932	ॡ 0942	ॢ 0952	ॣ 0962	अँ 0972
3	ः 0903	ओ 0913	ण 0923	ळ 0933	ॢ 0943	े 0953	ॣ 0963	अं 0973
4	ऐ 0904	औ 0914	त 0924	ळ 0934	ॣ 0944	े 0954	। 0964	आ 0974
5	अ 0905	क 0915	थ 0925	व 0935	ँ 0945	ँ 0955	॥ 0965	औ 0975
6	आ 0906	ख 0916	द 0926	श 0936	े 0946	ॢ 0956	ॣ 0966	अ 0976
7	इ 0907	ग 0917	ध 0927	ष 0937	े 0947	ॢ 0957	ॣ 0967	अ 0977
8	ई 0908	घ 0918	न 0928	स 0938	ै 0948	ॢ 0958	ॣ 0968	अ 0978
9	उ 0909	ड 0919	न 0929	ह 0939	ँ 0949	ख 0959	३ 0969	ज़ 0979
A	ऊ 090A	च 091A	प 092A	ं 093A	ो 094A	ग 095A	४ 096A	ष 097A
B	ॠ 090B	छ 091B	फ 092B	ी 093B	ो 094B	ज 095B	५ 096B	ग 097B
C	ॡ 090C	ज 091C	ब 092C	ॢ 093C	ौ 094C	ड 095C	६ 096C	ज 097C
D	ँ 090D	झ 091D	भ 092D	ऽ 093D	ॢ 094D	ढ 095D	७ 096D	१ 097D
E	ऐ 090E	ज 091E	स 092E	ा 093E	ि 094E	फ 095E	८ 096E	ड 097E
F	ए 090F	ट 091F	य 092F	ि 093F	ौ 094F	य 095F	९ 096F	ब 097F

	0A0	0A1	0A2	0A3	0A4	0A5	0A6	0A7
0	ॐ 0A00	ॐ 0A10	ठ 0A20	र 0A30	ी 0A40	ॐ 0A50	ॐ 0A60	ॐ 0A70
1	ँ 0A01	ॐ 0A11	ड 0A21	ॐ 0A31	ॢ 0A41	ॣ 0A51	ॐ 0A61	ॐ 0A71
2	ं 0A02	ॐ 0A12	ढ 0A22	ल 0A32	ॢ 0A42	ॣ 0A52	ॐ 0A62	ॐ 0A72
3	ः 0A03	ॐ 0A13	ॢ 0A23	ल 0A33	ॣ 0A43	ॣ 0A53	ॐ 0A63	ॐ 0A73
4	ॐ 0A04	ॐ 0A14	ॢ 0A24	ॣ 0A34	ॣ 0A44	ॣ 0A54	ॐ 0A64	ॐ 0A74
5	अ 0A05	क 0A15	थ 0A25	व 0A35	ॐ 0A45	ॐ 0A55	ॐ 0A65	ॐ 0A75
6	आ 0A06	ख 0A16	द 0A26	स 0A36	ॐ 0A46	ॐ 0A56	ॐ 0A66	ॐ 0A76
7	इ 0A07	ग 0A17	ध 0A27	ॐ 0A37	ॐ 0A47	ॐ 0A57	ॐ 0A67	ॐ 0A77
8	ई 0A08	घ 0A18	न 0A28	स 0A38	ॐ 0A48	ॐ 0A58	ॐ 0A68	ॐ 0A78
9	उ 0A09	ड 0A19	ॐ 0A29	ह 0A39	ॐ 0A49	ॐ 0A59	ॐ 0A69	ॐ 0A79
A	ऊ 0A0A	च 0A1A	प 0A2A	ॐ 0A3A	ॐ 0A4A	ॐ 0A5A	ॐ 0A6A	ॐ 0A7A
B	ॐ 0A0B	ॐ 0A1B	ॐ 0A2B	ॐ 0A3B	ॐ 0A4B	ॐ 0A5B	ॐ 0A6B	ॐ 0A7B
C	ॐ 0A0C	ॐ 0A1C	ॐ 0A2C	ॐ 0A3C	ॐ 0A4C	ॐ 0A5C	ॐ 0A6C	ॐ 0A7C
D	ॐ 0A0D	ॐ 0A1D	ॐ 0A2D	ॐ 0A3D	ॐ 0A4D	ॐ 0A5D	ॐ 0A6D	ॐ 0A7D
E	ॐ 0A0E	ॐ 0A1E	ॐ 0A2E	ॐ 0A3E	ॐ 0A4E	ॐ 0A5E	ॐ 0A6E	ॐ 0A7E
F	ॐ 0A0F	ॐ 0A1F	ॐ 0A2F	ॐ 0A3F	ॐ 0A4F	ॐ 0A5F	ॐ 0A6F	ॐ 0A7F

# Problems with Unicode

- The encoding doesn't represent the language in its true essence.
- Hindi, Sanskrit and most other Indian languages are centred around phonetic values. Hence the encoded token should ideally represent the entire sound rather than it being split into different symbols for vyanjana and maatra.
- We cannot figure out anything about the letter from its corresponding encoding. Which section of vyanjana it belongs to, whether it has a sweet or a bitter sound etc.
- The vyanjana symbols have a pre-added 'अ' (क् + अ = क). It is this conjoined sound which gets representation in Unicode rather than the plain क्.

# Existing Work: Romanization

- There are several methods of transliteration from Devanagari to the Roman script (a process known as romanization) which share similarities, although no single system of transliteration has emerged as the standard. Eg. SLP1, Velthius, Harvard-Kyoto etc.
- These can represent not only the basic Devanagari letters, but also phonetic segments, phonetic features and punctuation. SLP1 also describes how to encode classical and Vedic Sanskrit.
- *This motivates the second part of the project which involves auto-detection of the script through Machine Learning. As figuring out the difference between script schemas which use latin letters is a non-trivial problem.*



# Universal Encoding

- In this project, we set out to design a universal encoding for all Indian languages.
- The encoding is 16-bit unlike Unicode which requires 24-bits to represent most Indian characters.
- Initial 5 bits are for specifying the script (hence we can support 32). Next 6 bits are for the vyanjana (consonant). Last 5 bits are for the swar/maatra (vowel).
- In most cases, the vyanjana as well as the mantra both require 24-bits each in Unicode. We fit both in 16 bits.
- The encoding respects the structure of the language as described by Panini and we can figure out important characteristics about the letter just by seeing the encoding.
- The first 3 bits of vyanjana represent the source of origin (tongue, throat, teeth etc). The last bit of vyanjana represents prayatna. The second last bit of vyanjana represents sweet or pungent sound. The last bit of swar represents deergh (long) or laghu (short).

# Encoding for Devanagari

## Vyanjana

	000	001	010	011	100	101	110	111
000	क	ख	ग	घ		ह	ङ	
001	च	छ	ज	झ		श	ञ	य
010	ट	ठ	ड	ढ		ष	ण	र
011	त	थ	द	ध		स	न	ल
100	प	फ	ब	भ			म	व
101	क्	ख़	ग़	ज़	ड़	ढ़	फ़	
110	Uddat	Anuddat	ः	ं	ँ	ऽ		
111	Latin	Space	Punc	Numbers	Vaidika			Null

# Encoding for Devanagari

Swara				
	00	01	10	11
000	अ	आ	ऐ	
001	इ	ई	ए	ए
010	ऋ	ॠ		ऐ
011	ॡ	ॢ	ओ	ओ
100	उ	ऊ	औ	औ
101	अं	आं	अु	अु
110	अँ	ऐँ	औँ	ॐ
111				Null

# Transliteration Pipeline

- **Fragmentation:** Splitting the given text into smaller fragments (words, sentences, paragraphs etc). The assumption shall be that the script and encoding remain same through these fragments if not through the entire text.
- **Script Detection:** Figuring out the scripts and encodings for the various fragments through a Naive Bayes model.
- **Tokenisation:** Splitting the fragment further into tokens, each of which represent a single sound. Similar to the concept of English syllables. कि will be seen as one single token under this model.
- **Universalisation:** Conversion of the token to the universal 16-bit encoding designed by us. Done through pre-populated hash maps.
- **Specification:** Conversion of the universal encoding to the specified script using pre-populated hash maps.

# Work Till Now

- Built a lexer for Devanagari and SLP1 using Flex.
- Populated code maps for the above two languages.
- Built an encoder and decoder using the above maps on C++.
- Can perform transliteration between the above two languages.
- Script addition to the above framework is easy and any new script can be added to the list in 5-10 mins.

# Results

- **Devanagari:** पक्षि इस पृष्ठ पर इन्टरनेट पर उपलब्ध विभिन्न हिन्दी एवं देवनागरी सम्बंधित साधनों की कड़ियों की सूची है। इसमें उपकरण टूल्स शामिल हैं।
- **Encoded:** 0400 001f 02a4 073f 07e4 03a0 073f 0408 02bf 0220 073f 0400 02e0 073f 07e4 03df 0200 02e0 03c7 0200 073f 0400 02e0 073f 07f0 0400 03e0 045f 0360 073f 04e4 0464 03df 03c0 073f 00a4 03df 0345 073f 07e7 04e0 067f 073f 0347 04e0 03c1 0040 02e5 073f 03a0 04df 0440 067f 0364 0300 073f 03a1 0360 03cf 067f 073f 0005 073f 0000 0584 01ef 067f 073f 0005 073f 03b1 0105 073f 00ab 075f 073f 07e4 03a0 04c7 067f 073f 07f0 0400 0000 02e0 02c0 073f 0211 03ff 03a0 073f 01a1 04c4 03e0 073f 00ab 067f 075f
- **SLP1:** pakzi isa pfzWa para inwaranewa para upalabDa viBinna hindl evaM devanAgarl sambaMDita sADanoM kl kaN\*iyom kl sUcl hE. isameM upakaraRa wUlsa SAmila hEM.

# Work Ahead

- Script addition to the current framework.
- Planned scripts are: Assamese, Bengali (Bangla), Devanagari, Gujarati, Gurmukhi, Kannada, Malayalam, Oriya, Tamil, Telugu, Velthius, ITRANS, Harvard-Kyoto.
- Building corpus for English, SLP1, Velthius, ITRANS, French, Italian and German.
- Building a naive bayes model on substring frequencies to differentiate the above scripts.

# References

- <http://unicode.org/charts/>
- [https://en.wikipedia.org/wiki/  
Indian\\_Script\\_Code\\_for\\_Information\\_Interchange](https://en.wikipedia.org/wiki/Indian_Script_Code_for_Information_Interchange)
- [https://en.wikipedia.org/wiki/  
Devanagari\\_transliteration](https://en.wikipedia.org/wiki/Devanagari_transliteration)
- <https://betterexplained.com/articles/unicode/>
- <https://en.wikipedia.org/wiki/SLP1>