Hindi-Kannada Named Entity Transliteration: Issues and Possible Solutions

Annarao Kulkarni^{*}, B. R. Srivatsa and Chetan Baji

Centre for Development of Advanced Computing (C-DAC), Bengaluru - 560100, Karnataka, India; annaraogk@yahoo.com

Abstract

Indian languages belong to four language families, namely, the Indo-Aryan, Dravidian, Tibeto-Burman and the Austro-Asiatic. Hindi and Kannada belong to Indo-Aryan and Dravidian family respectively and are evolved from the ancient Brahmi script and have a common phonetic structure. But the Named Entity writing convention is different due to dialectic influence, language specific rules, and other factors. Due to this, the Named Entity Transliteration from Hindi to Kannada and vice versa is not one to one character mapping. This introduces many problems in Machine Translation (MT), Cross Lingual Information Retrieval (CLIR) and Parallel corpus creation between Hindi and Kannada. The paper discusses the Named Entity Transliteration issues encountered between Hindi and Kannada during the parallel corpora creation from Hindi to Kannada for the Indian Language Corpus Initiative (ILCI) project. In this paper, we discuss cases of no exact equivalence character between Hindi and Kannada, multiple mappings, diacritic marks, loan words and language specific transliteration issues in detail and propose the possible solution to resolve the problem. At implementation level, one may make use of either Finite-State Transducers (FST) or Regular Expressions

Keywords: Hindi, Kannada, Named Entity, Regular Expressions, Transliteration

1. Introduction

Named Entities (NE) are union of subsets of person name, place name, organization name, monetary expressions, dates, numerical expressions. Parallel pair of Named Entities between two languages of homogeneous and heterogeneous family and their classification plays an important role in Natural Language Processing (NLP) applications. Extracting NE is a challenging task. Various methods such as Hidden Marko Model (HMM), Conditional Random Fields (CRF), rule-based and hybrid approaches are in practice for Named Entity Extraction. Typically, NE of person name and place name are transliterated and organization name is translated for CLIR and MT.

Transliteration is not translation; the text remains the same, only the script in which it is rendered is changed³. Transliteration is usually done so as to preserve the

accuracy of pronunciation as much as practically feasible without distorting the target language rules. In reality, it is a bit difficult in case of Named Entity Transliteration between Hindi and Kannada. Practically, we came across many issues while translating tourism domain ILCI corpora from Hindi to Kannada. In this paper, we address various issues we encountered during transliteration and propose possible solutions to resolve the issues. Rule based solution is better compared to other techniques. Techniques such as use of NE annotated parallel corpora and use of comparable corpora for NE transliteration may suffer due to lack of coverage and constrained resources. Methods based on phonetic information and statistical methods are also not feasible because they are computing incentive and language family dependent. Hence, designing a rule to address Named Entity transliteration and then implement those using either FST or regular expressions is more feasible.

^{*} Author for correspondence

2. Related Works

Suggests finite state Hindi Urdu Machine Transliteration using Finite-State Transducers based on language specific characteristics, FST and UIT¹. Problems have been faced at various levels like no equivalence of character set between two languages, diacritic marks and language specific issues for Urdu-Hindi Machine translation. Have come up with rule base for Hindi-Punjabi transliteration⁸. Suggest transliteration among Indian Languages using WX notation, by making use of common representation to transliterate among Indian languages⁶. In Hindi, syllabic vowels occur in the middle and end e.g., आइहोले (AihOlE), नकिारागुआ (nikArAguA). In Kannada, sequence of syllabic vowel is not allowed and pure vowel can occur only as a first syllable. This makes WX notation methodology difficult to implement.5 have proposed statistical method for transliteration system across Indian languages using parallel corpora and algorithm similar to Soundx. However, Indian languages are resource constrained and due to lack of resources and NE aligned parallel corpora, rule based approach is suitable for Named Entity transliteration. ILCI parallel corpus is one such effort to create parallel corpora for 17 Indian languages with Hindi as source language. We have used ILCI parallel corpora for Hindi-Kannada to arrive at rule base for transliteration of NE from Hindi to Kannada.

3. Hindi and Kannada

Hindi is one of the official languages of India and is written in Devanagari Script. Devanagari is a script shared by ten other official languages of India and more than 70% Indians can understand and speak Hindi to a certain level. Hindi language has 49 dialects as per census of India 2001. The most spoken of them are Bhojpuri, Rajasthani, Chhattisgarhi, Magahi, Pahari, Bundeli, Bagheli, Awadhi, Marwari, Mewari etc6. Various phonological transformations take place and deviation is observed from Standard Hindi in written script due to dialectic influence thus resulting in different writing conventions. This is one of the issues while transliterating Named Entities from Hindi to other Indian languages for NLP applications. Guidelines of Central Hindi Directorate are considered as authentic for writing in Hindi. Both Hindi and Kannada languages evolved from the ancient Brahmi script and have a common phonetic structure with some differences in vowels, consonants, diacritic marks. This makes transliteration difficult. Kannada is one of the four major Dravidian languages. Kannada is rich in morphology and is agglutinative in nature. Kannada has many dialects and the major dialects are the old Mysore dialect, the coastal or Mangalore dialect and the northern or Dharwad dialect⁷. Kannada is a diglossic language. Literary variety differs in several respects from the spoken variety in phonology, morphology, lexicon and syntax. But, literary variety is the same across dialects.

4. Hindi-Kannada NE Transliteration

In the following section, we discuss issues with respect to no equivalent characters, multiple mapping, diacritic marks, language specific issues, loan words and suggest suitable solutions for NE transliteration from Hindi to Kannada.

4.1 Vowels

 Table 1.
 Hindi-Kannada Vowels

Hindi	Kannada	ITRANS
अ	అ	А
आ	ಆ	А
इ	Q	Ι
ई	ਚੱ	Ι
3	സ	U
ক	സ	U
뀪	ಋ	R
-	ఎ	E
ए	ධ	Е
ऐ	ສ	ai
-	ఒ	0
ओ	ఓ	0
औ	꿊	au

Hindi contains 11 vowels and 10 dependent vowel symbols. Examples for them are Γ , Υ . These dependent vowel symbols are also called matras. Apart from these 11 vowels, Devanagari $\breve{\Psi}$ [U+090D] and $3\breve{\Pi}$ [U+09011] are also considered as vowels for transliterating loan words borrowed from Urdu and English with their corresponding matras [U+0945] and $\breve{\Gamma}$ [U+0940]. This

issue is discussed in section 4.6. Kannada contains 13 vowels and 12 dependent vowel symbols. If we compare Hindi-Kannada vowel set, Hindi doesn't have short [e] and [o]. Kannada contains both long and short [e] and [o]. Kannada language doesn't have \breve{V} and \Im equivalents.

Case-1 Diphthongs: Sequence of syllabic vowels does not occur in Kannada language (S.N. Sridhar, 2007). But, there is a tendency to write Kannada diphthongs \mathfrak{D} [ai] and \mathfrak{Z} [au] as sequence of syllabic vowels in Hindi like \mathfrak{HF} [a|A+i|I] and \mathfrak{HF} [a|A+u|U]. In such cases, replace Hindi sequence of syllabic vowels \mathfrak{HF} [a|A+i|I] with Kannada dipthong \mathfrak{D} [ai] and replace sequence of Hindi syllabic vowels \mathfrak{HF} [a|A+u|U] with Kannada dipthong \mathfrak{Z} [au]. For example

ఐ [ai] – आइहोले - ఐహొళి (aihoLe) आइजोल - ఐజూల (aijOla) ఔ [au] – आउटडोर – ఔటడూలరో (auTDOr) आउटलुक – ఔటోలుకో (auTluk)

Diphthong matras: Replace Γ + ड़|ई with Kannada ' [ai] matra. Replace Γ + 3]ऊ with Kannada ^ভ [au] matra. For example,

कराईकुडी – ಕರ**ೈಕುಡೆ (karaikuDi)** ताइवान – ತೈವಾನ (taivAna) राउडी – ರೌಡೆ (rauDi)

Case-2 Long Vowels: Long vowels (and matras) tend not to occur word finally, except as grammatical morphemes⁷. In Kannada, long vowel matra occurs only in initial and middle positions for Named Entities. Replace long vowel matra for the word finally syllable with short [a],[i],[e],[u] and [o] matra respectively. For example

```
भविानी – ಭೆವಾನೆ (bhivAni)
मनाली – ಮೆನಾಲೆ (manAli)
मसूरी – ಮಸೂರೆ (masUri)
पणजी – ಪಣಜೆ (paNaji)
रायबरेली – ರಾಯಬರೇಲೆ (rAyabarEli)
गनिी – ಗಿನೆ (gini)
पेरू – ಪೆರು (peru)
जम्मू – ಜಮ್ಮ (jammu)
बसािऊ – ಬಿಸಾಮ (bisAvu)
कन्याकुमारी – ಕನ್ಯಾಕುಮಾರೆ (kanyAkumAri)
पुरी – ಮರೆ (puri)
बेल्लारी – ಬಳ್ಳಾರೆ (baLLAri)
```

Case-3 Glide Insertion: When two vowels come together, they are replaced with glides [y] and [v] in

Kannada. If the first vowel is front vowel [i, I] then [y] is inserted and [v] is inserted if the first vowel is back or low⁷. This rule is not applicable to diphthongs [ai] and [au]. Kannada words do not end with pure vowel and sequence of syllabic vowels does not occur in Kannada. If the final syllable of the NE ends with [A], [i, I] and [u, U] then replace with [vA], [yi] and [vu] respectively. For example,

बहराईच – थळठन्र्यार्थ (baharAyica) तरिअनन्तपुरम् – डैपेवर्त्तरुं (tiruvanaMtapuraM) नोइडा – ನೋಯೆढन (nOyiDA) हसिुआ – ಹैस्रावेन (hisuvA) गोआ – तौश्विन (gOvA) नकािरागुआ – तौरुन्ठनेरोत्वेन (nikArAguvA) हरदोई - ळेटवौश्ट्यी (haradOyi) राऊ – ठन्वे (rAvu) बसािऊ – थीरुन्वे (bisAvu)

Case-4 Multiple Vowel Mapping: Hindi does not contain short [e] and [o]. For the syllable final [E] and [O] word, case 3 will be used, otherwise replace Hindi [E] and [O] with Kannada short vowel [e] and [o] respectively. This is a case of multiple mapping. However, both these possibilities are more suitable for data mining and Kannada data searching task. For example

राऊ – ರಾಮ (rAvu) हरदोई – ಹರದೋಯೆ (haradOyi) एर्नाकुलम – ಎರ್ನಾಕುಲಂ (ernAkulaM) क्रोएशयाि – ಕ್ರೊಯೇಶೆಯ (kroyEshiya) वेनेजुएला – ವೆನೆಜುವೆಲ (venejuvela)

4.2 Vowel Modifiers

Table 2. Hindi-Kannada Vowel Modifiers

Vowel modifiers	Hindi	Kannada	ITRANS
Chandrabindu	U	-	.N
Anuswara		0	М
Visarga	:	8	Н

After a consonant, vowel or matra, a character can be used which modifies the vowel sound and is called vowel modifier⁹. This can be Chandrabindu, Anuswar or Visarga. Kannada doesn't have Chandrabindu. Following table shows the list of vowel modifiers. Case-5 Chandrabindu [.N]: Chandrabindu denotes nasalization of the preceding vowel. It is substituted with anuswara in Kannada. For example

राँची - ರಾಂಚಿ (rAMci) हाँगकाँग – ಹಾಂಗ ಕಾಂಗ (hAMga kAMga)

Case-6 म [ma] as word final syllable: Word final syllable [ma] is replaced by anuswara [M] in Kannada. For example

मज़ोिरम – ಮಿಜೊರಂ (mijoraM) सकि्कमि – ಸೆಕ್ಕಿಂ (sikkiM) भद्राचलम – ಭದ್ರಾಚಲಂ (bhadrAcalaM) चदिंबरम – ಚಿದಂಬರಂ (cidaMbaraM) श्रीकाकुलम – ಶ್ೀಕಾಕುಲಂ (shrIkAkulaM)

4.3 Diacritic Marks

A diacritic mark called Nukta is used in Devanagari for some languages to represent sounds from other languages. It takes the form of a dot placed below a character.

Case-7 Nukta: In Hindi, Nukta is used to modify consonants क, ख, ग, ज, इ, ढ, फ as क, ख़, ग, ज़, इ, ढ़, फ़ respectively. However, there is no equivalent Nukta character in Kannada. It is necessary to drop Nukta while transliterating from Hindi to Kannada. Replace क with क ख़ with ख, ज़ with ज, ढ़ with ढ, फ़ with फ respectively. For example

धारवाड़ – ಧಾರವಾಡ (dhAravADa) चंडीगढ़ – ಚಂಡಿಗಢ (caMDigaDha) फ़्रान्स – ಫ್ರಾನ್ಸ (phrAnsa) ज़मि्बाब्वे – ಜಿಂಬಾಬ್ವೆ (jiMbAbve) अज़रबैजान – ಅಜರಬ್ಶಿಜಾನ (ajarabaijAna)

4.4 Consonants

Hindi and Kannada share the same consonant set except for the non-varga consonant [La – \forall] i.e., Devanagari consonant $\overline{\omega}$. Devanagari $\overline{\omega}$ character is not used in Hindi. Consonants are categorized according to their phonetic properties into 5 Vargas and the remaining ones fall under the Non-Varga category. The last consonant in each Varga is called nasal consonant.

Case-8 Nasal consonants: If the nasal consonant is followed by any consonant belonging to that Varga, then, that nasal consonant is replaced by anuswar [.m] in Kannada. For example

```
अङ्कोला – ७०ँँ०९७ (aMkOla)
केन्द्र – हैं१०८ँ४ (kEMdra)
```

 ङ followed by क | ख | ग | घ is replaced by [M]

 ञ followed by च | छ | ज | झ is replaced by [M]

 ण followed by ट | ठ | ड | ढ is replaced by [M]

 न followed by त | थ | द | ध is replaced by [M]

 म followed by प | फ | ब | भ is replaced by [M]

Table 3.	Hindi-Kannada	Varga	Consonants
----------	---------------	-------	------------

क / रु [ka]	ख / ಖ [kha]	ग / त [ga]	घ / द्ध [gha]	ङ / छ [~Na]
च / ಚ [ca]	छ / ಛ [cha]	ज / ಜ [ja]	झ / ಝ [jha]	স / জ [~na]
て / ಟ [Ta]	ठ / ರ [Tha]	ਤ / ਕ [Da]	ढ / द्व [Dha]	ण / ಣ [Na]
त / छ [ta]	थ / ऴ [tha]	द / ದ [da]	೮ / ಧ [dha]	೯ / ನ [na]
प / ळ [pa]	फ / ಫ [pha]	ब / ಬ [ba]	भ / ಭ [bha]	म / ಮ [ma]

Table 4.	Hindi-Kannada	Non-Varga	Consonants
----------	---------------	-----------	------------

य / ಯ [ya]	र / ರ [ra]	ਕ / ಲ [la]	ಡ / ವ [va]	श / শ্ত [sha]
ष / ಷ [Sha]	स / रु [sa]	ह / ळ [ha]	당 [La]	

Case-9 [\mathfrak{G} - La] consonant in Kannada: In most of the cases (not always), Hindi $\overline{\mathfrak{R}}$ [la] is substituted with \mathfrak{G} [La] in Kannada. This is a case of multiple mapping for transliterating from Hindi to Kannada. For data mining, this rule is more useful. In such cases, search with both [la] and [La] replacements in Kannada. For example

केरल – हैंश्टेंध (kEraLa) तमलिनाडु – डंಮೆ'धेर्जन्खे (tamiLunADu) बेल्लारी – थर्ध'धन्ट (baLLAri) बेंगल्रू – थ'oriधन्छ (beMgaLUru) भुसावल – थ्रेरान्चध (bhusAvaLa)

4.5 Language Specific

Case-10 Conjunct formation for \overline{er} [la]: \overline{er} followed by consonant is often replaced by \overline{er} + halant to form a conjunct cluster. However, this is not a general rule. For example

कोलकाता – ಕೊಲ್ಕತ್ d (kOlkatta) सुलतानगंज – ಸುಲ್ತಾನಗಂಜ (sultAnagaMja) जैसलमेर – ಜೈಸಲ್ಮೇರ (jaisalmEra) तरिुनलवेली – ತಿರುನಲ್ವೇಲೆ (tirunalvEli)

Case-11 $\mathbf{\zeta}$ [ra] followed by Consonant: If $\mathbf{\zeta}$ [ra] is followed by a consonant, then, replace $\mathbf{\zeta}$ with $\mathbf{\zeta}$ + halant.

In Kannada, repha is inserted to form a conjunct cluster. This is not a general rule, though. For example

```
कारगलि – ಕəರ್rেೆಲ್ (kArgil)
इटारसी – ಇಟಾರ್ಸೆ (iTArsi)
करनूल – ಕರ್ನೂಲ (karnUla)
बरमूडा – ಬರ್ಮುಡ (barmuDa)
```

```
Case-12 Replace Hindi लै and बै with ಲ್ಯಾಂ [lyAM]
and ಬ್ಯಾಂ [byAM] in Kannada. For example
फनिलैंड - ಫೆನ್ಲ್ಯಾಂಡ್ (phinlyAMD)
आइसलैण्ड – ಐಸ್ಲ್ಯ್ಯಾಂಡ್ (aislyAMD)
न्यूज़ीलैंड – ත්ಯೂಜಿಲ್ಯಾಂಡ್ (nyUjilyAMD)
बैंकाक – ಬ್ಯಾಂಕಾಕ್ (byAMkAk)
बैंड – ಬ್ಯಾಂಡ್ (byAMD)
बैंक – ಬ್ಯಾಂಕ್ (byAMk)
```

4.6 Loan Words

Case-14 While transliterating loan words from Hindi or English to Kannada, replace 3Π with [A] and $\overline{\Psi}$ with [a]e].

5. Conclusion

The rules discussed in this paper are applied to translation activities for creating parallel corpora from Hindi to Kannada for the ILCI project and these rules will help to get better quality of transliteration and translation output. Creating parallel dictionary of NE between two languages is a difficult task. Techniques, such as, use of parallel corpora for transliteration is not feasible due to constrained resources and lack of coverage. However, parallel corpora help in creating rule base for effective Named Entity transliteration and translation for CLIR and MT applications. This work can be extended to other Dravidian languages. However, there is no one-to-one character set mapping among Dravidian languages but the rules are the same. This work can be used as preprocessor and can be used with other existing transliteration techniques like phonetics based and parallel corpora based transliteration techniques.

6. Acknowledgments

While writing this paper, we referred Hindi-Kannada Tourism Domain ILCI parallel corpora with Hindi as source language. Examples of place name for this work have been taken from ILCI parallel corpora of Tourism domain. We would like to thank ILCI project team for this.

7. References

- 1. Abbas Malik MG, Boitet C, Bhattacharyya P. Hindi Urdu Machine Transliteration using Finite-State Transducers. COLING; 2008. p. 537-44.
- Jha GN. The TDIL program and the Indian Languages Corpora Initiative (ILCI). Proceedings of the Seventh Conference on International Language Resources and Evaluation (LREC); 2010. p. 982-5.
- Murthy KN, Badugu S. Roman Transliteration of Indic Scripts. Proceedings of the Tenth International Conference on Computer Applications (ICCA); 2012. p. 19-26.
- 4. Zhang M, Duan X, Pervouchine V, Li H. Machine Transliteration: Leveraging on Third Language. Proceedings of the 23rd International Conference on Computational Linguistics; 2010. p. 1444–52.
- Srivastava R, Bhat RA. Transliteration System across Indian Languages Using Parallel Corpora. Proceedings of the 27th Pacific Asia Conference on Language, Information and Computation; 2013. p. 390-8.
- 6. Gupta R, Goyal P, Diwakar S. Transliteration among Indian Languages Using WX Notation. KONVENS; 2010. p. 147-50.
- Sridhar SN. Modern Kannada Grammar. New Delhi: Manohar Publishers; 2007.
- 8. Goyal V, Lehal GS. Hindi to Punjabi Machine Transliteration System. Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies; 2011.
- 9. Indian Script Code for Information Interchange (ISCII). Bureau of Indian Standards; 1991.