

English to Indian Languages Machine Transliteration System at NEWS 2010

**Amitava Das¹, Tanik Saikh², Tapabrata
Mondal³, Asif Ekbal⁴, Sivaji
Bandyopadhyay⁵**

**Department of Computer Science and
Engineering^{1,2,3,5}
Jadavpur University**

**Department of Computational
Linguistics⁴
University of Heidelberg**

Transliteration Unit (TU)

**Bengali, Hindi and Kannada and Tamil
and English words divided into
Transliteration Units (TUs):**

English word TU: C*V*

**where C represents a consonant and V
represents a vowel.**

Indian Languages Words TU: C+M?

**where C represents a vowel or a
consonant or a conjunct and M
represents the vowel modifier or
matra**

**Contextual information in the form of
collocated TUs considered**

Firdausi → Fi | r | dau | si
फिं | रौं | दौं | सी

Banphool → Ba | n | phoo | l
बाौ | नौ | फूौ | लौ

Babooji → Ba | boo | ji
बाौ | बूौ | जीौ

English to Indian Languages Machine Transliteration

- ❖ Plausibility of transliteration from each English TU to various Indian Languages (ILs) candidate TUs calculated
 - ILs candidate TU chosen with **maximum transliteration probability**
 - ↳ Equivalent to choosing the most appropriate sense of a word in the source language to identify its representation in the target language
- ❖ System learns mappings automatically from the **bilingual training corpus**
 - **Learning guided by linguistic features**
 - **Output of mapping is a decision list classifier**
- ❖ Machine transliteration obtained by direct orthographic mapping or phonetic mapping
 - Equivalent ILs TU for each English TU in the input identified and placed in ranked order.

Transliteration Models

Orthographic Transliteration Models

- ❖ **Trigram Model**
- ❖ **Joint Source-Channel Model (JSC)**
- ❖ **Modified Joint Source-Channel Model (MJSC)**
- ❖ **Improved Modified Joint Source-Channel Model (IMJSC)**

International Phonetic Alphabet (IPA) Model

Trigram Model

**Previous and next source
TUs as the context**

$$P(S | T) = \prod_{k=1}^K P(< s, t >_k | s_{k-1}, s_{k+1})$$

$$S \rightarrow T(S) = \arg \max_T \{P(T) \times P(S | T)\}$$

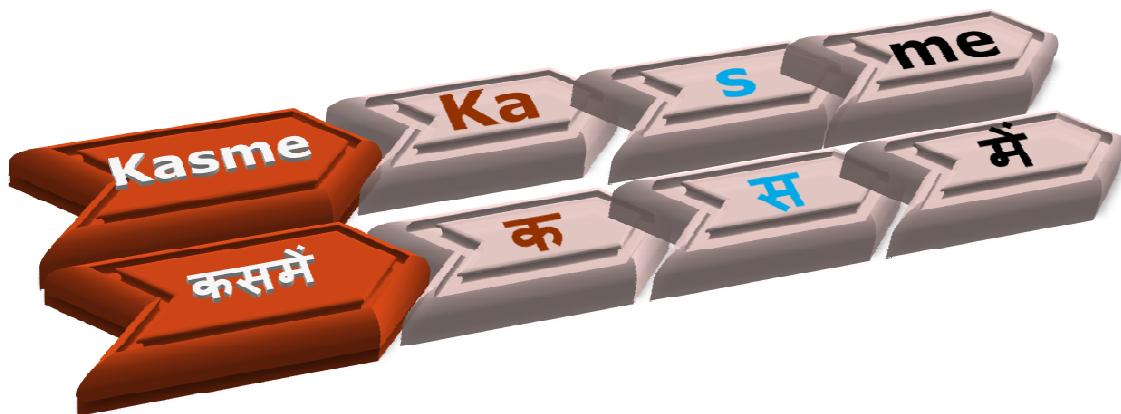


Joint Source-Channel Model (JSC)

Previous TUs with reference to the current TUs in both the source (s) and the target sides (t) are considered as the context.

$$P(S|T) = \prod_{k=1}^K P(< s, t >_k | < s, t >_{k-1})$$

$$S \rightarrow T(S) = \arg \max_T \{P(T) \times P(S|T)\}$$

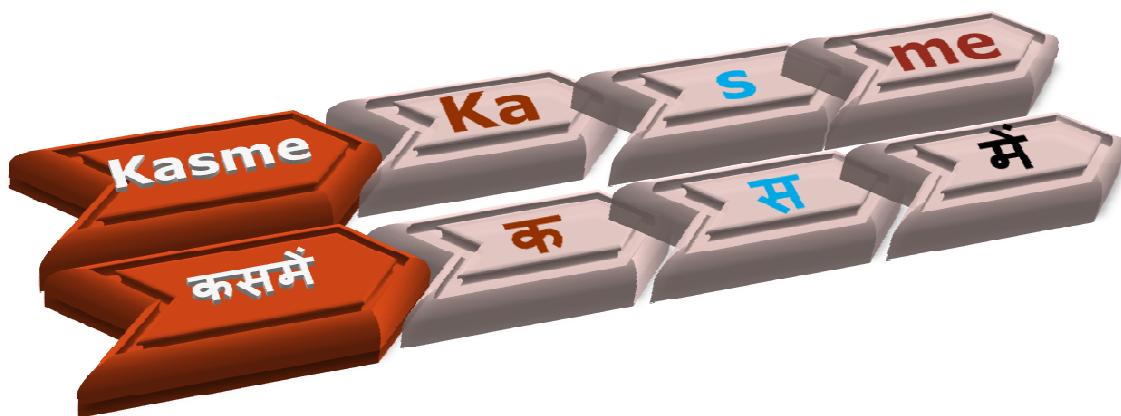


Modified Joint Source-Channel Model (MJSC)

the previous and the next TUs in the source and the previous target TU are considered as the context.

$$P(S | T) = \prod_{k=1}^K P(< s, t >_k | < s, t >_{k-1}, s_{k+1})$$

$$S \rightarrow T(S) = \arg \max_T \{P(T) \times P(S | T)\}$$



Improved Modified Joint Source-Channel Model (IMJSC)

The previous two and the next TUs in the source and the previous target TU are considered as the context.

$$P(S | T) = \prod_{k=1}^K P(< s, t >_k | s_{k+1}, < s, t >_{k-1}, s_{k+1})$$

$$S \rightarrow T(S) = \arg \max_T \{P(T) \times P(S | T)\}$$



International Phonetic Alphabet (IPA) Model

Large number of dictionary words in NEWS 2010 data set

	Training	Development	Test
Bengali	7.77%	5.14%	6.46%
Hindi	27.82%	15.80%	3.7%
Kannada	27.60%	14.63%	4.4%
Tamil	27.87%	17.31%	3.0%

IPA Symbol Set

Phoneme	Example	Translation
AA	<i>odd</i>	AA-D
AE	<i>at</i>	AE-T
AH	<i>hut</i>	HH-AH-T
AO	<i>ought</i>	AO-T
AW	<i>cow</i>	K-AW
AY	<i>hide</i>	HH-AY-D
B	<i>be</i>	B-IY
CH	<i>cheese</i>	CH-IY-Z
D	<i>dee</i>	D-IY
DH	<i>thee</i>	DH-IY
EH	<i>Ed</i>	EH-D
ER	<i>hurt</i>	HH-ER-T
EY	<i>ate</i>	EY-T
F	<i>fee</i>	F-IY
G	<i>green</i>	G R-IY-N
HH	<i>he</i>	HH-IY
IH	<i>it</i>	IH-T
IY	<i>eat</i>	IY-T
JH	<i>gee</i>	JH-IY
K	<i>key</i>	K-IY
L	<i>lee</i>	L-IY
M	<i>me</i>	M-IY
N	<i>knee</i>	N-IY
NG	<i>ping</i>	P-IH-NG
OW	<i>oat</i>	OW-T
OY	<i>toy</i>	T-OY
P	<i>pee</i>	P-IY
R	<i>read</i>	R-IY-D
S	<i>sea</i>	S-IY
SH	<i>she</i>	SH-IY
T	<i>tea</i>	T-IY
TH	<i>theta</i>	TH-EY-T-AH
UH	<i>hood</i>	HH-UH-D
UW	<i>two</i>	T-UW
V	<i>vee</i>	V-IY
W	<i>we</i>	W-IY
Y	<i>yield</i>	Y-IY-L-D
Z	<i>zee</i>	Z-IY
ZH	<i>seizure</i>	S-IY-ZH-ER

CRF Based IPA Model

- A **pre-processing** module checks whether a targeted source English word is a valid **dictionary word** or not.
- In the target side we use our TU segregation logic to get **phoneme wise transliteration pattern**.
- Modeled as a **sequence labeling problem** as transliteration pattern changes depending upon the contextual phonemes in source side and TUs in the target side.

Output Ranking

	Ranking Order				
Word Type	1	2	3	4	5
Dictionary	IPA	IMJSC	MJSC	JSC	TRI
Non-Dictionary	IMJSC	MJSC	JSC	TRI	-

Bengali Experimental Result

Parameters	Accuracy		
	BSR	BNSR1	BNSR2
Accuracy in top-1	0.232	0.369	0.430
Mean F-score	0.818	0.845	0.875
Mean Reciprocal Rank (MRR)	0.325	0.451	0.526
Mean Average Precision (MAP) _{ref}	0.232	0.369	0.430

- ✓ **BNSR2 has achieved the highest score among all the submitted Runs.**

Hindi Experimental Result

Parameters	Accuracy		
	HSR	HNSR1	HNSR2
Accuracy in top-1	0.150	0.254	0.170
Mean F-score	0.714	0.752	0.739
Mean Reciprocal Rank (MRR)	0.308	0.369	0.314
Mean Average Precision (MAP) _{ref}	0.150	0.254	0.170

- ✓ **HNSR1 and HNSR2 runs are ranked as the 5th and 6th among all submitted Runs.**

Kannada Experimental Result

Parameters	Accuracy	
	KSR	KNSR1
Accuracy in top-1	0.056	0.055
Mean F-score	0.663	0.662
Mean Reciprocal Rank (MRR)	0.112	0.169
Mean Average Precision (MAP) _{ref}	0.056	0.055

Tamil Experimental Result

Parameters	Accuracy	
	TSR	TNSR1
Accuracy in top-1	0.013	0.082
Mean F-score	0.563	0.760
Mean Reciprocal Rank (MRR)	0.121	0.142
Mean Average Precision (MAP) _{ref}	0.013	0.082