

# SANSKRIT SANDHI SPLITTING AND MERGING TOOLS

Shubham Bhardwaj  
Neelamadhav Gantayat  
(Mini Project)  
under the guidance of  
Prof. Rahul Garg  
Prof. Sumeet Agarwal

February 24, 2016

# Outline

- 1 Sandhi
  - What governs this interference?
  - Sandhi Corpus
- 2 Sanskrit Transliteration
  - Need for a transliteration tool
  - Sanskrit Transliteration
- 3 Sandhi Corpus Format
  - Processing corpus
- 4 Dictionary Creation
- 5 Online Sandhi Tools
  - Evaluation

# Sandhi

## Sandhi



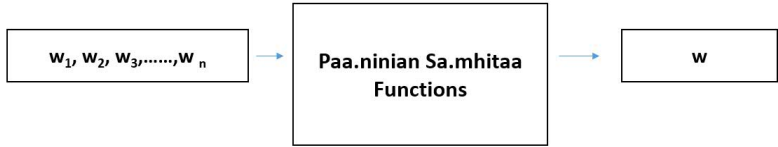
For every  $i$ ,  $w_i$  is a word  $\rightarrow$  External Sandhi

e.g., तस्मै + एतत्  $\rightarrow$  तस्मायेतत्

For some  $i$ ,  $w_i$  is a prefix, verb root or a suffix  $\rightarrow$  Internal Sandhi

e.g. वि + छेद  $\rightarrow$  विच्छेद

# What governs this interference?



- A.s.taadhyayii of Paa.nini
- Sa.mhita governs sutras 73-157 of Chapter 1 of Book 6 and all sutras of Chapter 3 and 4 of Book 8
- Total number of Paa.ninian Sandhi Sutras - 271
- Paa.ninian Sutras codified in the form of Sets and Functions

# Automatic Evaluation of Sandhi Splitters and Generators

- Sanskrit Sandhi Analyzer and Generator<sup>1</sup>  
(Dr. G.N. Jha, Special Centre for Sanskrit Studies, JNU)
- Sandhi-Splitter and Generator<sup>2</sup>  
(Dr. Amba Kulkarni, Department of Sanskrit Studies, University of Hyderabad)
- The Sanskrit Reader Companion and the Sandhi Engine<sup>3</sup>  
(Dr. Gerard Huet ,Computational Linguistics, INRIA, France)

---

<sup>1</sup><http://sanskrit.jnu.ac.in/sandhi/viccheda.jsp?itext=>

<sup>2</sup><http://52.25.246.194/sc1/>

<sup>3</sup><http://sanskrit.inria.fr/DICO/reader.fr.html> ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻

# Sandhi Corpus

- 40 different sandhi-split corpora available on the University of Hyderabad website<sup>4</sup>
- Problems - Insufficient splits, no splits, wrong splits, typos

<b>Text (first 100 cases)</b>	<b>Errors</b>
Vinodini	42%
Manjusa	44%
Vyutpattivada	58%
Dutvakyam	8%
Short-stories	9%

- Creating a Srimad Bhagvad Gita Corpus

---

<sup>4</sup><http://sanskrit.uohyd.ac.in/Corpus/>

# Need for a transliteration tool

- Various data sources of Sanskrit literature were encoded in different scripts.
- Different tools use different transliteration schemes for input and output
- From our survey identified some web tools and some standalone tools
  - No bulk conversion option
- Need a tool which can convert from one script to another in bulk (Files and folders)

# Sanskrit Transliteration

Devanāgarī	IAST	Harvard-Kyoto	ITRANS	Velthuis	SLP1
अ	a	a	a	a	a
आ	ā	A	A/aa	aa	A
इ	i	i	i	i	i
ई	ī	I	I/ii	ii	I
उ	u	u	u	u	u
ऊ	ū	U	U/uu	uu	U
ए	e	e	e	e	e
ऐ	ai	ai	ai	ai	E
ओ	o	o	o	o	o
औ	au	au	au	au	O
फ	f	R	RR/Rʰi	.r	f
भ	b	RR	RR/Rʰi	.rr	F
क्ष	ṣ	IR	LL/Lʰi	.l	x
ख	ṭ	IRR	LL/Lʰi	.ll	X
म	m	M	M/ n/ m	.m	M
ह	h	H	H	.h	H
ञ			.N		~

Figure: Example: One step conversation<sup>5</sup>

- Among these SLP1 script has complete encoding scheme for Devanagari

<sup>5</sup>[https://en.wikipedia.org/wiki/Devanagari\\_transliteration](https://en.wikipedia.org/wiki/Devanagari_transliteration)



# Sandhi Corpus Format

- Manually curated corpus from University of Hyderabad <sup>6</sup>

## Formats

- Single words split into multiple words

चार्थे => च+अर्थे

- Multiple words split into multiple words

सङ्केतो लक्षणा => सङ्केतः+ लक्षणा

- Multiple words span across multiple lines

व्याख्यां विजहितां => व्याख्याम्+ विजहिताम्+  
विजहितां कुर्वे => विजहिताम्+ कुर्वे

<sup>6</sup><http://sanskrit.uohyd.ac.in/Corpus/>

# Processing corpus

- Filtered out single word, multi split instances
- Combined multi lines into single line and converted them to multi word instance.
- One way of ensuring correctness is checking for the split words in the dictionary

# Dictionary Creation

Various types of dictionary, example: <sup>7</sup>

- Sanskrit to Sanskrit dictionary
- Sanskrit to English dictionary
- English to Sanskrit dictionary

Dictionary words creation

- Parsed “Sanskrit to Sanskrit” and “Sanskrit to English”
- Created a comprehensive list of words from all the dictionary.

---

<sup>7</sup><http://www.sanskrit-lexicon.uni-koeln.de/>

# Filtering Corpus

- Identified all the sandhi splits from the corpus where all the split words are available in the dictionary
- out of 97674 samples we identified 15325 samples where all the splits are present in the dictionary

# Evaluation

- Sanskrit Sandhi Analyzer and Generator - 1523/7217 (21.10%)
- Sandhi-Splitter and Generator - 4248/7217 (58.86%)
- The Sanskrit Reader Companion and the Sandhi Engine - Not yet done

# Future Work

- Validate dictionary output
- Check for Sandhi splitting and Sandhi Merging tools apart from the above three
- Create corpus using Sandhi Merging tool by merging words or dictionary.
- Manually create Sandhi corpus from “Sri Mad Bhagwat Gita”
- Create our own Merging and splitting tool if required.