



Chandojñānam

A Sanskrit Meter Identification and Utilization System

Hrishikesh Terdalkar, Arnab Bhattacharya

18th World Sanskrit Conference, 2023

Department of Computer Science and Engineering,
Indian Institute of Technology Kanpur

Introduction

Recite

योऽन्तः प्रविश्य मम वाचमिमां प्रसुप्तां yo'ntaḥ praviśya mama vācamimāṃ prasuptāṃ
सञ्जीवयत्यखिलासक्तिधरः स्वधाम्ना। sañjīvayatyakhilāśaktidharaḥ svadhāmnā |
अन्यांश्च हस्तचरणश्रवणत्वगादीन् anyāṃśca hastacaraṇaśravaṇatvagādīn
प्राणान्नमो भगवते पुरुषाय तुभ्यम्॥ prāṇānnamo bhagavate puruṣāya tubhyam | |

What is the **chanda** in this verse?

Ans: वसन्ततिलका (Vasantatilakā)

Recite

योऽन्तः प्रविश्य मम वाचं प्रसुप्तां yo'ntaḥ praviśya mama vācaṃ prasuptāṃ

Why does it feel odd?

- Deviation from a known pattern
- How do we know these patterns?
 - Sanskrit Prosody!

Background

- Classification of syllables
 - Pronunciation dependent
 - **Laghu** (*short*)
 - Letters with short vowels
 - **Guru** (*long*)
 - Letters with long vowels
 - **Laghu** letters followed by a joint letter (**saṃyogaḥ**)
 - Last letter of a **pāda** (conditional)
- **Mātrā**: Laghu 1, Guru 2
- **Gaṇa**: Sequence of three letters ($2^3 = 8$)

Chanda

- Types
 - **Akṣaracchanda**: Sequences of laghu-guru
 - **Samavṛtta**, **Ardhasamavṛtta** and **Viṣamavṛtta**
 - **Mātrācchanda**: Counts of mātrā
- Literature: **Vṛttaratnākaraḥ**, **Chandovicitih**, **Chandomañjarī** etc.

Motivation



Enthusiast



Poet



Teacher

Line
...
Scansion
...



Researcher



Why another meter identification tool?

Aim

- Add more user-friendly features.
 - Catch errors in the text and suggest corrections!
-
- Web-based application
 - Python library
 - Three input modes: (1) plain text, (2) images (3) text files
 - Two OCR Engines: (1) Google Drive OCR (2) Tesseract OCR
 - Transliteration support (powered by `indic-transliteration`)
 - Two meter identification modes: (1) line mode (2) verse mode
 - Fuzzy matching support using edit distance comparison
 - Informative scansion display
 - Downloadable results

Feature Comparison

Features		[Mis07]	[MGS13]	[Raj20]	[Nei22]	Chandojñanam
Availability	Web Interface	✓ ¹	✓ ²	✓	✓	✓
	Software Library			✓	✓	✓
Input	Text	✓	✓	✓	✓	✓
	Arbitrary Lines					✓
	Multiple Verses					✓
	Textfile Upload				✓	✓
	Image Upload					✓
Functionality	Meter Identification	✓	✓	✓	✓	✓
	Error Tolerance			✓	✓	✓
	Fuzzy Matching			✓		✓

Table 1: Feature comparison of extant meter identification systems

¹<http://sanskrit.sai.uni-heidelberg.de/Chanda/HTML/> no longer functional.

²<https://sanskritlibrary.org:8080/MeterIdentification/> no longer functional.

System

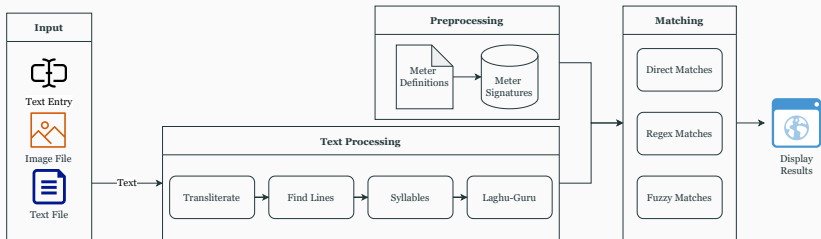


Figure 1: Workflow of Chandojñanam

वृत्त	पाद	गण	लक्षण	अक्षरसङ्ख्या	मात्रा	यति
शार्दूलविक्रीडित		मसजसततग	गगगललगलगललललगगगलगगलग	19	30	12,7
शालिनी		मततगग	गगगगगलगगलगग	11	20	4,7
अपरवक्त्र	1	ननरलग	लललललललललगलगलग	11	14	
अपरवक्त्र	2	नजजर	लललललगलललगलगलग	12	16	
सौरभ	1	सजसल	ललगलगललललललगल	10	13	
सौरभ	2	नसजग	लललललललललगलगलग	10	13	
सौरभ	3	रनभग	गलगललललललगललग	10	14	
सौरभ	4	सजसजग	ललगलगललललललगलगलग	13	18	
अनुष्टुभ्	1	----लग--	----लग--	8		
अनुष्टुभ्	2	----लगल-	----लगल-	8		

Figure 2: Generic Chanda definition format

- Column 'Pāda': index of pāda in the meter
- Uniform treatment of samavṛtta, ardhasamavṛtta and viṣamavṛtta
- Regex pattern (*regular expression*) definition

Metrical Database

- Two types of dictionaries
 - Signature of individual **pādas** (CHANDA_SINGLE)
 - Signature of consecutive **pādas** (CHANDA_MULTIPLE).

```
CHANDA_SINGLE = {  
    'LGGLGGLGGLGG': ['Bhujāṅgaprayāta'],  
    '[LG] [LG] [LG] [LG] LG [LG] [LG]': ['Anuṣṭubh (Pāda 1)'],  
    '[LG] [LG] [LG] [LG] LGL [LG]': ['Anuṣṭubh (Pāda 2)']  
}
```

```
CHANDA_MULTIPLE = {  
    'LGGLGGLGGLGGLGGLGGLGGLGG': ['Bhujāṅgaprayāta (Pāda 1-2)'],  
    '[LG] [LG] [LG] [LG] LG [LG] [LG] [LG] [LG] [LG] [LG] LGL [LG]': [  
        'Anuṣṭubh (Pāda 1-2)'  
    ]  
}
```

Input

- Text input
- Textfile input
- Image file input

ध्यायेदाजानुबाहुं धृतशरधनुषं बद्धपद्मासनस्थं
पीतं वासो वसानं नवकमलदलस्पर्धिनेत्रं प्रसन्नम् ।
वामाङ्कारूढ सीतामुखकमलमिलल्लोचनं नीरदाभं
नानालङ्कारदीप्तं दधतमुरुजटामण्डनं रामचन्द्रम् ॥

UploadChoose image fileBrowse

Google OCR Tesseract OCROCR + Identify

Figure 3: Upload a screenshot of a verse for meter identification

Text Processing

- Common text processing pipeline for all input modes
- Transliteration (*powered by indic-transliteration*³)
 - Detect scheme
 - Convert to internal scheme (Devanagari)
- Line identification
 - Standard line-end markers: '\n', '।', '॥', ''
- Syllabification (*powered by sanskrit-text*⁴)
 - भारत = भा + र + त
- **Laghu-Guru** markers
 - Standard rules
 - \exists **Chanda** where last letter is **laghu** (**Padānta Laghu**)
 - Last letter not forced to be **guru**

³<https://pypi.org/project/indic-transliteration/>

⁴<https://pypi.org/project/sanskrit-text>

Meter Identification Algorithm

Algorithm 1: Meter Identification

Data: Metrical Database (MD)

Input: lg -signatures of every 'line' in the input ($T = \{lg_1, lg_2, \dots, lg_n\}$)

Output: Result set containing exact or fuzzy matches

```
1 forall  $lg \in T$  do
2    $SM_1 = \text{FindDirectMatch}(lg, \text{'CHANDA\_SINGLE'})$ 
3    $SM_2 = \text{FindDirectMatch}(lg, \text{'CHANDA\_MULTIPLE'})$ 
4    $RM = \text{FindRegexMatch}(lg, \text{'CHANDA\_SINGLE'} + \text{'CHANDA\_MULTIPLE'})$ 
5    $DM = SM_1 + SM_2 + RM$ 
6    $FM = \phi$ 
7   if  $DM = \phi$  then
8     |  $FM = \text{FindFuzzyMatch}(lg)$ 
9   end
10  return  $DM + FM$ 
11 end
```

Direct Matching

Algorithm 2: Direct Matching

Input: *lg-signature*

Output: Result set containing exact matches

```
1 Function FindDirectMatch(lg, 'MD')
2    $M_1 = \text{Query}(\textit{lg}, \textit{'MD'})$            // dictionary lookup
3    $M_2 = \phi$ 
4   if  $M_1 = \phi$  then                       // if no match found
5     | if the last letter of lg is laghu then
6     | |    $lg_1 = \textit{replace last letter of lg with guru}$ 
7     | |    $M_2 = \text{Query}(lg_1, \textit{'MD'})$ 
8     | end
9   end
10  return  $M_1 + M_2$ 
```

What?

Finding approximate and close matches if exact match not found

Why?

- Digitally available Sanskrit text can be erroneous
 - Manual data entry
 - Post-scanning OCR followed by manual correction
- Types of Errors
 - Characters may be misspelt, e.g., रु (ru) as रू (rū)
 - Characters may be missing, e.g., वर्गे (vargai) as वगै (vagai)
 - Characters may be misidentified, e.g., ऋ (ṛ) as क्र (kra)
 - Characters may get split, e.g., ख (kha) as रव (rava)
- Several such errors can affect the metrical pattern of the text

How?

- **Problem:** Finding the *nearest matching string* for the *lg-signature* of the text
- Compute Levenshtein edit distance of the observed pattern (powered by `python-Levenshtein`⁵)
- Normalize the edit distance by the length of target pattern

$$\text{Similarity} = 1 - \frac{\text{Levenshtein distance}}{\text{length of target match}}$$

- Topmost k matches as the possible fuzzy matches ($k = 10$)
- Suggestions: changes required to transform input into target
 - insert, delete, replace




⁵<https://pypi.org/project/python-Levenshtein/>

Fuzzy Matching Example – Matching

Input text Output Scheme: Match Input

नमस्ते सदा वत्सले मातृभुमे
त्वया हिन्दुभूमे सुखं वर्धितोऽहम्।
महामङ्गले पुण्यभूमे त्वदर्थे
पतत्वेष कायो नमस्ते नमस्ते॥

Verse Mode Line Mode Identify

Results   

Akṣarāṇi	न	म	स्ते	स	दा	व	त्स	ले	मा	तृ	भु	मे
Laghu-Guru	ल	ग	ग	ल	ग	ग	ल	ग	ग	ल	ल	ग
Gaṇa	य		य		य			स				
Counts	12 अक्षराणि, 19 मात्राः											
Jāti	जगती											
Chanda	भुजङ्गप्रयात (1 edit)											+ Fuzzy

Figure 4: Meter identification with fuzzy matching

Fuzzy Matching Example – Suggestions

Chanda		भुजङ्गप्रयात (1 edit)	- Fuzzy	
Fuzzy Matches				
#	Chanda	Gaṇa	Cost	Similarity
1	भुजङ्गप्रयात	यययय	1	91.7%
	[[['न', 'म', 'स्ते'], ['स', 'दा'], ['व', 'त्स', 'ले'], ['मा', 'वृ', 'r(भु)[G]{भ्रू}', 'मे']]]			
2	स्रग्विणी	रररर	2	83.3%
	[[['i(G)', 'd(न)', 'म', 'स्ते'], ['स', 'दा'], ['व', 'त्स', 'ले'], ['मा', 'वृ', 'भ्रु', 'मे']]]			
3	विध्वङ्कमाला	तततगग	2	81.8%
	[[['d(न)', 'म', 'स्ते'], ['स', 'दा'], ['व', 'त्स', 'ले'], ['मा', 'वृ', 'r(भु)[G]{भ्रू}', 'मे']]]			
4	हंसमाला (पाद 1-2)	सरभतगग	3	78.6%
	[[['न', 'i(L)', 'i(L)', 'स्ते'], ['स', 'दा'], ['व', 'त्स', 'ले'], ['मा', 'वृ', 'r(भु)[G]{भ्रू}', 'मे']]]			
5	इन्द्रवंशा	ततजर	3	75.0%
	[[['d(न)', 'म', 'स्ते'], ['स', 'दा'], ['व', 'त्स', 'r(ले)[L]'], ['मा', 'वृ', 'i(G)', 'i(G)', 'मे']]]			

Figure 5: Fuzzy matching suggestions

Identification Modes

- *Line mode*: Treat the input as a set of arbitrary lines
 - Useful for checking meter of a single line or a set of unrelated lines
- *Verse Mode*: Treat the input as a collection of verses
 - Useful for identifying meter of a single or multiple verses
 - Utilizes information from other lines of the verse
 - Re-order results if required

Input text
Output Scheme: Match Input

Input text: माता रामो मम पिता रामचन्द्रः। स्वामी रामो मल्लखा रामचन्द्रः। सर्वस्वं मे रामचन्द्रो ददातुर् नान्यं जाने नैव जाने न जाने॥

Output text: माता रामो मम पिता रामचन्द्रः। स्वामी रामो मल्लखा रामचन्द्रः। सर्वस्वं मे रामचन्द्रो ददातुर् नान्यं जाने नैव जाने न जाने॥

Verse Mode Line Mode

Results

1. शालिनी

Akṣarāṇi	मा	ता	रा	मो	म	म	पि	ता	रा	म	च	न्द्रः
Laghu-Guru	ग	ग	ग	ग	ल	ल	ग	ग	ल	ग	ग	ग
Gaṇa	म			भ			य			य		
Counts	12 अक्षराणि, 20 मात्राः											
Jāti	जगती											
Chanda	वातोर्मि (1 edit)											

Figure 6: Meter identification in (a) *Line mode* and (b) *Verse mode*

Chanda					Chanda				
वातोर्मि (1 edit)					शालिनी (2 edits)				
Fuzzy Matches					Fuzzy Matches				
#	Chanda	Gaṇa	Cost	Similarity	#	Chanda	Gaṇa	Cost	Similarity
1	वातोर्मि	मभतगग	1	90.9%	1	शालिनी	मतगग	2	81.8%
	[[[मा, 'ता], [रा, 'मो], [म, 'म], [द(म), 'ता], [रा, 'म, 'च, 'न्द्रः]]]					[[[मा, 'ता], [रा, 'मो], [द(म), 'र(म)[G], [पि, 'ता], [रा, 'म, 'च, 'न्द्रः]]]			
2	प्रहृषिणी	मनजरग	2	84.6%	2	वातोर्मि	मभतगग	1	90.9%
	[[[मा, 'ता], [रा, 'र(मो)[L], [म, 'म], [पि, 'ता], [L(L), 'र(L), 'म, 'च, 'न्द्रः]]]					[[[मा, 'ता], [रा, 'मो], [म, 'म], [द(पि), 'ता], [रा, 'म, 'च, 'न्द्रः]]]			

Figure 7: Fuzzy matches in (a) *Line mode* and (b) *Verse mode*

- Transliteration-based primitive⁶ multilingual support

The screenshot displays two side-by-side panels for identifying Sanskrit prosody from other Indian languages. Each panel has an 'Input text' field, an 'Output Scheme' dropdown, and an 'Identify' button. Below the input is a 'Results' section with a 'Hide Scansion' button and a table of prosodic data.

(a) Marathi:

Input text: सुगंगि सदा पठो सुजनवाक्य कानी पठो कलंक मलिचा छटो विषय सर्वचा नावटो संदंभिकमयी दडो मुरविशा हटने अडो विद्योग घडला रडो मन भवच्छरिणी अडो

Output Scheme: Devanagari

Results Table:

1. पृथ्वी	Hide Scansion																
Akṣarāṇi	सु	सं	ग	ति	स	दा	ष	डो	सु	ज	न	वा	क्य	का	नी	ष	डो
Laghu-Guru	ल	ग	ल	ल	ल	ग	ल	ग	ल	ल	ग	ल	ग	ल	ग	ल	ग
Gaṇa	ज	स	ज	स	य	ल	ग										
Counts	17 अक्षराणि, 24 मात्राः																
Jāti	अत्यष्टिः																
Chanda	पृथ्वी																

(b) Telugu:

Input text: పిలికిం జిన్నడు: శంఖ చక్ర యుగముంటేయి పంపించడే పరిసంఘము శీఠ దర్శనశం బిన్నెల రాళ్ళింఱాం ఈ ధున్దము జక్క రాళ్ళకు పొన్నెల్లకొత్త శ్రీ పంచ పరిసంఘంపైను పడడు గణ ప్రాణాచార్యునిచ్చు

Output Scheme: Telugu

Results Table:

1. మల్లేశ్వరీశీఠ	Hide Scansion																						
Akṣarāṇi	పి	రి	కి	ం	జి	న్న	డు	శం	ఖ	చ	క్ర	యు	గ	ము	ం	జే	దో	యి	పం	ధి	ం	వ	డే
Laghu-Guru	ల	ల	క	క	ల	ల	ల	ల	ల	ల	ల	ల	ల	ల	ల	ల	ల	ల	ల	ల	ల	ల	ల
Gaṇa	ప	భ	ర	వ	మ	య	ల	క															
Counts	20 syllables, 30 morae																						
Jāti	కృతిః																						
Chanda	మల్లేశ్వరీశీఠ																						

Figure 8: Identification from other Indian languages, (a) Marathi (b) Telugu

⁶Uses the rules and metrical database of Sanskrit prosody

Evaluation for Error Correction

- Single text from different sources can differ in several places
- Three versions of **Meghadūta**⁷ composed by **Kālidāsa**
 - Wikisource, sanskritdocuments.org and GRETIL
- Texts with more metrical variety
 - **Śāntavilāsa** (36 verses) (12 distinct meters)
 - **Śrīrāmarakṣāstotra** (39 verses) (9 distinct meters)
 - **Rājendrakarṇapūra** (72 verses) (4 distinct meters)
- Manually tagged meters for each verse from these texts
- Realistic evaluation: Simulate digitization pipeline for all four texts
 - Generate PDF from Wikisource text
 - Run both the OCR systems
 - Obtain the OCR-ed versions of the text
- Total 14 text versions, 1038 verses, exhibiting 17 distinct meters

⁷Also used by [Raj20] for evaluation

Results

		Meghadūta					Śāntavilāsa			Rāmarakṣā			Rājendrakarṇapūra			Total
		SD	GR	WS	GO	TO	WS	GO	TO	WS	GO	TO	WS	GO	TO	
Number of Verses		117	111	123	123	123	36	36	36	39	39	39	72	72	72	1038
Unique Chanda		1	1	1	1	1	12	12	12	9	9	9	4	4	4	17
Erroneous Verses		20	79	2	31	77	13	16	31	1	4	13	12	26	71	396
Correct	[Nei22]	20	79	2	30	66	11	13	14	0	2	9	12	24	36	318 (80.3%)
Meters	[Raj20]	19	79	2	30	75	12	15	24	1	2	9	12	26	58	364 (91.9%)
Identified	Chandojñānam	20	79	2	31	77	13	16	29	1	3	9	12	26	71	389 (98.2%)

Table 2: Error tolerance of meter identification systems. (Versions are WS: Wikisource, GO: Google OCR, TO: Tesseract OCR, SD: sanskritdocuments.org, GR: GRETIL.) **Chandojñānam** is able to detect correct **chanda** from erroneous verses 98.2% of the times.

- Successfully detected and corrected two errors from Wikisource version of **Meghadūta**

Error #1

- Line: कालक्षेपं ककुभसुरभौ पर्वते पर्वते ते (Pāda 3, Śloka 1.23)
- Incorrect word पर्वते (should be पर्वते)
- Likely due to OCR error and an oversight by the curator
- Suggestion: [[['का', 'ल', 'क्षे', 'पं'], ['क', 'कु', 'भ', 'सु', 'र', 'भौ'], ['प', 'र्व', 'ते'], ['प', 'र्व(र्वे)[L]', 'ते'], ['ते']]]
- System correctly points to the location where a change is required

Error #2

- Line: साभिज्ञानप्रहितकुशलैस्ततद्वचोभिर्ममापि (Pāda 3, Śloka 2.53)
- Extra letter (त) present in the **sandhi** of words कुशलैः and तद्वचोभिः
- Suggestion: [[['सा', 'भि', 'ज्ञा', 'न', 'प्र', 'हि', 'त', 'कु', 'श', 'लै', 'd(स्त)', 'त', 'द्व', 'चो', 'भि', 'र्म', 'मा', 'पि']]]
- Points out correctly that a syllable needs to be deleted
- However, points to an incorrect syllable स्त to be deleted
 - Both स्त and त are **laghu** letters
 - Deletion of either letter \implies the correct metrical signature
 - Impossible for a meter identification based system

Conclusions

Conclusions and Future Work

- User-friendly meter identification system
- Several input options
- Focus on error detection and correction

Future Work

- Inclusion of **Mātrācchandās**
- Improvements to meter correction algorithm
- Possible consideration of semantic aspects
- More extensive support for Indian languages that use similar rules of prosody

Links

System: <https://sanskrit.iitk.ac.in/jnanasangraha/chanda/>

Source: <https://github.com/hrishikeshrt/chanda/>

References

References

- [MGS13] Keshav S Melnad, Pawan Goyal, and Peter Scharf.
Meter identification of sanskrit verse.
The Sanskrit Library, USA, 2013.
- [Mis07] Anand Mishra.
Sanskrit metre recognizer.
2007.
- [Nei22] Tyler Neill.
Skrutable: Another step toward effective sanskrit meter identification.
2022.
- [Raj20] Shreevatsa Rajagopalan.
A user-friendly tool for metrical analysis of sanskrit verse.
Computational Sanskrit & Digital Humanities, page 113, 2020.

Thank you!

Questions?