

A Model for Identifying and Predicting deteriorated text in Manuscripts using *Vrutta* Pattern Matching

¹Keertimalini A, ²T. Venkata Padmaja, ³Dr. B. Chandrasekaram

^{1,2}P. G Students, ³Associate professor,
Dept. of Computer Science,
National Sanskrit University, Tirupati 517 507.

Abstract: Manuscripts are important cultural artifacts which provide significant insights into the history and culture of ancient India. Manuscripts form an invaluable part of India's documentary heritage. But due to years of time lapse, physical & chemical factors caused naturally or artificially by humans (foreign invaders) led to loss of valuable information. Manuscripts that may have undergone just a little amount of damage, which when identified through the remaining available text, can be used to extract the complete knowledge it contains. This paper presents one such model to predict and retrieve the deteriorated text from the Sanskrit manuscripts using metrical analysis and laghu-guru pattern matching techniques. This paper also extends the scope of prediction of deteriorated texts to metered and unmetered text and deals with the more efficient corpus searching algorithms using Regular expressions.

Keywords: Manuscript, Metrical analysis (Vruttas), Pattern matching, Chandas Shastra, Deterioration, Regular expressions, Sanskrit corpus.

I. Introduction

1.1 Manuscripts - Memory of the world

National Mission for Manuscripts (NAMAMI), a project undertaken by the Ministry of Tourism and Culture defines a Manuscript as — “Handwritten composition on paper, bark, cloth, metal, palm leaf or any other material dating back at least seventy-five years that has significant scientific, historical or aesthetic value.” Manuscripts are a treasure house of knowledge. It not just covers the traditional domains such as religion, culture & philosophy, but also contains innumerable scientific information on Medicine (like *Sushruta and Charaka Samhita*), Astronomy and Astrology (*Jyothisa*), Vedic Mathematics (*Sulbasutras*), Yoga (*Patanjali Yogasutras*), Literature (*Sahitya*), History (*Purana-Itihasa*), Language Grammar (*Vyakaran*) and Lexicons (*Amarakosha*): India possesses one of the largest collections of manuscripts comprising more than five millions of manuscripts on different materials. They capture our thoughts, achievements, experience and lessons learnt from history; in other words, they constitute our ‘memory.’ UNESCO provided this recognition to the most valuable documentary heritage of the world - INDIA. Thus manuscripts form an important part of the Indian Knowledge system (IKS).

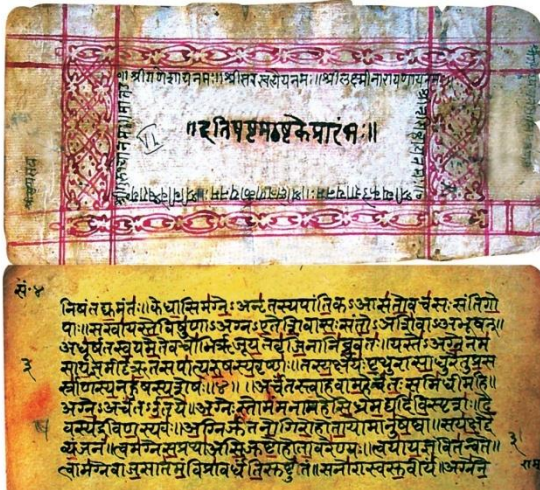


Fig. 1 Rigveda manuscript

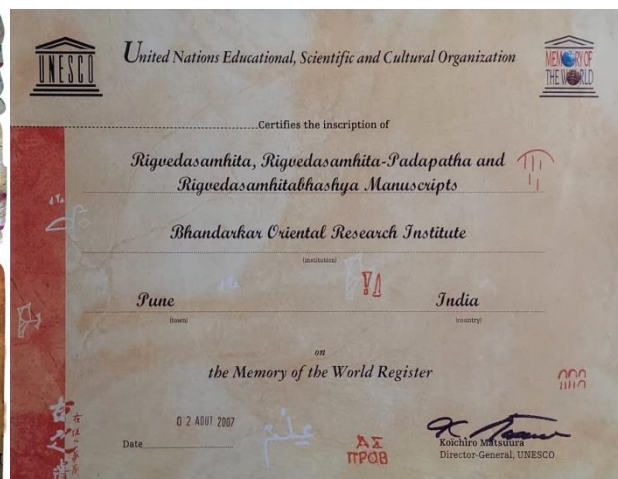


Fig. 2 UNESCO Certificate

1.2 Deterioration of Manuscripts

Deterioration is a loss of structural capacity with time by the action of the external agents or material leaching (Saiz and Laiz, 2000). These aesthetically valuable resources, which are easily more than seventy five years old, are usually written on birch bark, palm leaf, handmade paper, cloth, leather etc. Found in various shapes and sizes, the deterioration of these collections can be in the form of wear and tear, shrinkage, cracks, warping, bio-infestation, discoloration, abrasion, holes, etc. The ravages of time, extreme

climatic conditions and biological agents often destroy priceless cultural property and records^[1]. Here in this paper we deal with manuscripts which have lesser degree of deterioration (Point sized deterioration). The remaining text, when recognized, can help us retrieve the complete knowledge provided in the manuscript. This paper serves useful in preserving the archaic document and thus preserve the valuable knowledge it possesses.



Fig. 3 Deteriorated manuscript sample

1.3 Sanskrit Prosody - Chhandas Sastra & Vrutta ratnakara

Prosody or chhandas is the science of metres in Sanskrit literature. Metre is an important tool for poetical compositions. Sanskrit prosody was started with pingala's chhandas sastra, followed by Vrutta rathnakara of kedara and chandomanjari of Gangadasa. Chhandas sastra is one of the six vedangas or limbs of vedas. ("Chandah padau tu vedasya"). These texts deal with metres which specify the name of the metre and its particular gana through verses i.e. via 1) Sutra form or 2) Karika form. Dandin, a proficient sanskrit poet remarks "Prosody is like a boat/ship for the reader who wants to cross the ocean of poetry."

Characteristics of Classical Metres:

The Sanskrit prosody is divided into two namely vedic and classical. Metres used in vedas are known as vedic, and metres used in sanskrit sahitya (*Purana Itihasa, Kavyas*) are known as classical.

- Pada (Quarter)** - Sanskrit shlokas or verses generally contain 4 lines, each of which is referred to as Pada. Depending upon the Vrutta chosen by the poet, each quarter is either characterized by matras or aksharas (Syllables).
- Akshara** - A sanskrit syllable which can be pronounced distinctly. Eg. रा, मः in रामः
- Maatra** - each syllabic instant or morae of an Akshara. Eg. र्+आ+म्+अः in the word रामः
- Laghu** - A syllable is short when the vowel is short. The vowels अ, इ, उ, ऋ, ॠ (a,i,u,r,lr) are short. It is denoted by the symbol ' I '
- Guru** - i) A syllable is long when the vowels are long. आ, ई, ऊ, ऋ, ॠ.
ii) A short vowel can become long when it is followed by anusvaaram अं (am) or visargam अः (aH).
iii) If a conjunct consonant - called samyuktaakshara in Sanskrit) follows a short vowel, then the short vowel becomes long. Except प्र, ब्र, क
- Any short vowel at the end of पाद. A guru syllable is denoted by the symbol 'U'
- Vrutta** - Sanskrit metres that are classified according to the number of syllables occurring in each Paadam (quarter) of the verse.

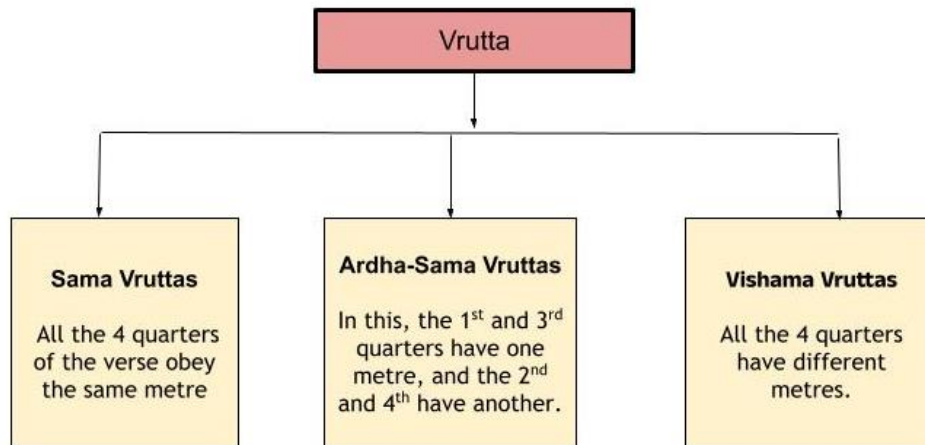


Fig. 4 Classification of Vruttas

- Gana** - A metrical feet is of combination of three syllables, which can either be guru or laghu. Thus $2^3 = 8$ combinations of ganas are possible.

आदिमध्यावसानेषु चरतायान्ति लाघवम् ।

भजसागौरवं यान्ति मनौतु गुरुलाघवम् ॥

This means, there is a short syllable at the start (aadi), middle (madhya) and end (avasaana) positions for the ganas य, र

¹ Deterioration of Manuscripts and Indigenous Methods used for Preserving These Documentary Cultural Heritage Ravindra Goswami, Department of Botany, R.B.S. College, Agra.

and त respectively (the other 2 syllables being long). Similarly, for भ, ज, and स there is a long syllable in the same positions (the other 2 syllables being short). Ma and na are fully long and short respectively.

“ यमाताराजभानसलगं ” - is usually the acronym used to refer to these ganas.

| यगण (ya) | रगण (ra) | तगण (ta) | भगण (bha) | जगण (ja) | सगण (sa) | मगण (ma) | नगण (na) |
|----------|----------|----------|-----------|----------|----------|----------|----------|
| IUU | UIU | UUI | UII | IUI | IUU | UUU | III |

Table 1 The 8 Ganas and their patterns

1.4 Regular Expressions

A regular expression (also known as regex) is a sequence of characters that defines a search pattern used for matching or manipulating strings of text. It can be used to search for specific patterns or sequences of characters within a larger text, and can be used in a variety of programming languages and applications. Regular expressions can include special characters and symbols, which can be used to specify certain types of characters or patterns. Here we have used a basic regular expression ‘([IU]*)’ to find and return the missing laghu-guru pattern of the deteriorated zone. Using regular expressions, the efficiency of searching and pattern matching of laghu-guru mappings are increased and so is the accuracy of output as well.

1.5 Pattern matching

According to DeepAI Pattern matching is defined as “An algorithmic task that finds pre-determined patterns among sequences of raw data or processed tokens.” Using regular expressions, the given data is checked and then matched with a process-of-elimination approach, like backtracking. This task makes exact matches from the existing *Vrutta* database and returns the pattern of the missing text.

1.6 Sanskrit Corpus

Sanskrit corpus data is a collection of Sanskrit text which contains text extracted from the epics like Ramayana, Mahabharata, Bhagavad Gita and Vedic texts like Vedas. In this paper the corpus data is used to extract the words which are predictable in the deteriorated zones. Some of the corpus data from SanskNET was collected and saved as files. When we find the pattern for the deteriorated zone, say P(A), the corresponding word is returned from this corpus via enhanced corpus searching methods.

II. Literature Survey

The restoration of damaged manuscripts has been a topic of interest for many researchers. Different works have been carried out related to the restoration and analysis of damaged manuscripts, as well as a system developed for Sanskrit verse metrical analysis. This paper has been mainly inspired from the research paper of Dr.B.Chandrasekharam, Associate Professor, National Sanskrit University titled “ Automatic digital rebuilding of text from deteriorated zones of manuscripts containing metered text. The author proposes a model via which it may be possible to minimize the loss of text, for the deteriorated zone of a manuscript which contains metered text. This is done by taking advantage of Chhanda sastra. It proposes a model/tool for rebuilding text from deteriorated/corrupted zones of manuscripts containing metered text.

Apart from this, the Jñānasaṅgrahaḥ system platform developed at IIT Kanpur, hosts an application for Sanskrit verse metrical analysis. This application allows users to obtain the metrical analysis of Sanskrit verses. Techniques for restoring damaged manuscripts via image processing proposed by Sai Siddharth Kota et al and the digital restoration of erased and damaged manuscripts, published by Keith T et al, are some prominent works in this arena. Another such work, authored by Diwakar Mishra, is titled "Strategies for Metrical Analysis of Sanskrit Text" and provides details about the metrical analysis of Sanskrit texts.

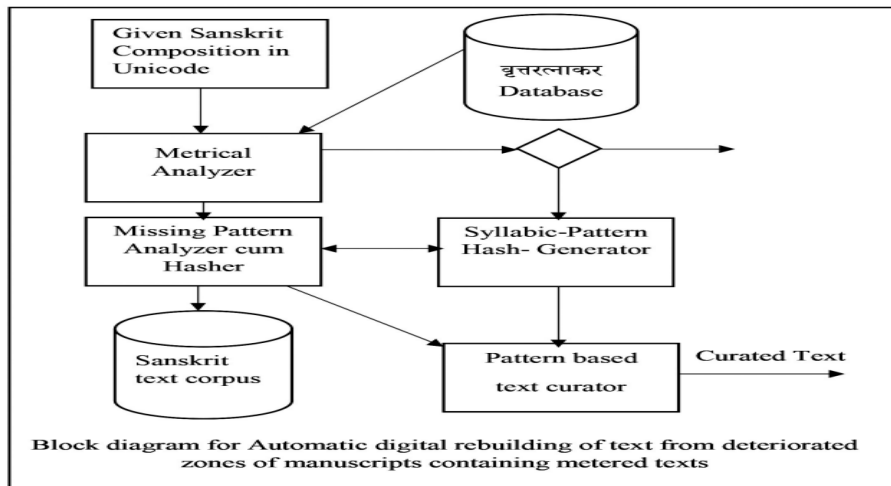


Fig. 5 Block diagram of the previous Model

III. Methodology

Manuscripts which have been deteriorated need to be pre-processed to extract the information from it. For identifying the available content and marking the deteriorated zones, we first digitize the manuscript and upload it in the web interface that has been created.

Step 1 - Digitization of manuscript and image enhancement

In this web application there is an image input box where the user will input a corrupted manuscript image. The noise is removed from the raw image using Guassian bandpass filtering as suggested in the paper “Digital Enhancement of Indian Manuscripts.”² The enhanced image is now of the readable form to the users.

Step 2 - Keying in of Unicode Sanskrit Text and marking missing zones

Below the image input box there is a text input box. There is a technique called Optical Character Recognition (OCR) Technology which will convert the image into editable text. But here the user is going to input the text manually in the text box and mark the missing or deteriorated parts in capital letters (A,B,C etc.) The text will be in Devanagari Unicode text format.

Process involved in the above 2 steps

1: Input - Sanskrit shloka which contains some deterioration.

मातामत्वितारामच

स्वामीमत्सखारामचन्द्रः।

2: Combine both lines into a single line.

मातामत्वितारामचस्वामीमत्सखारामचन्द्रः।

3: Label the deteriorated zones using alphabets.

मातामत्वितारामचस्वामीमत्सखारामचन्द्रः।

Step 3 - Mapping of laghu guru:

Using the laghu-guru assigning process that has been explained above, the laghu-guru patterns of the given input text is found. The source code has been designed in such a way that the words will be split into Aksharas which will then be patternised into binary laghu-guru pattern. The output of the inserted shloka will be represented as:

Input :

मातामत्वितारामचस्वामीमत्सखारामचन्द्रः।

Output:

UU A UIUUIU BUU CUUIUUIUU

```

93 def display_lagu_guru(txt):
94     char_split=list(split_clusters(txt))#UUU
95     lagu_list=अउरकवएँऑओओ कणखगजघडडणणरुपयभमयरलळवशषफभस्ररु क ख
96     guru_list=अउरकवएँऑओओ कणखगजघडडणणरुपयभमयरलळवशषफभस्ररु क ख
97     lagu_str=अउरकवएँऑओओकणखगजघडडणणरुपयभमयरलळवशषफभस्ररु क ख
98     guru_str=अउरकवएँऑओओ कणखगजघडडणणरुपयभमयरलळवशषफभस्ररु क ख
99     deergamtras=अउरकवएँऑओओ
100     laghumatras=अउरकवएँऑओओ
101
102     final_list=''
103     pattern = re.compile('[A-Za-z]')
104     for i in char_split:
105         #print(i)
106         if pattern.match(i): ## added 3 line to check whether given char is alphet or add to deteriorated string
107             final_list=final_list+i
108             continue #end of the code
109         if len(i)==1:
110             if i in lagu_str:
111                 final_list=final_list+'U'
112             elif i in guru_str:
113                 final_list=final_list+'I'
114
115         elif len(i)==2 :
116             j=i[-1]
117             if j in lagu_str:
118                 final_list=final_list+'1'
119             elif j in guru_str:
120                 final_list=final_list+'0'
121             elif j in laghumatras:
122                 final_list=final_list+'I'
123             elif j in deergamtras:
124                 final_list=final_list+'0'
125         elif len(i)>=3 and char_split[-1]!=i:
126             final_list=final_list[:-1]
127             final_list=final_list+'UU'
128
129     #print(final_list)
130     final_str=final_list
131
    
```

²Digital Enhancement of indian Manuscript, Yashodhar Charitra , Sai siddarth Kota, Raja Massand, Abhinaya agrawal and Preety Singh

Fig. 6 Source code for laghu-guru mapping

Step 4 : Detection of P(A) via pattern matching

To identify the exact I,U denoted as P(A), for deteriorated text A, firstly a database containing a set of Vrutta patterns (*Sama, Ardha-sama, Vishama*) is created. Through this database, the P(A) is found by comparing the P(I) with the entries of the data set. The chandas of the shloka can also be found. The following steps will assist in understanding the process in detail:-

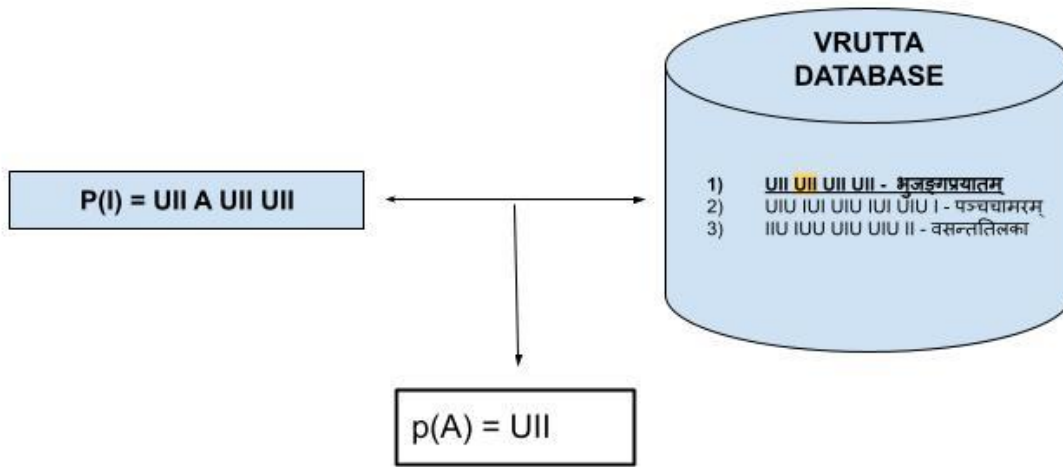


Fig. 7 Pattern matching process

1) Creating a database of containing Vruttas and its corresponding pattern.

| | S.NO | NAME OF THE VRUTTA | GANNA | NO. OF SYLLABLES IN EACH PADA | NO. OF PADAS | PATTERN | PATTERN |
|------------|------|--------------------|-------|-------------------------------|--------------|---|--|
| सम वृत्ता: | 1 | शालिनी | मलतमग | 11 | 4 | UUU UUU UUU UUU | UUU UUU UUU UUU UUU UUU UUU UUU UUU UUU UUU |
| | 2 | भुजङ्गप्रयातम् | यययय | 12 | 4 | UUU UU UU UU UU UU UU UU UU UU UU UU | UUU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU |
| | 3 | पञ्चचामरम् | जजजज | 13 | 4 | UUUU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU | UUUU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU |
| | 4 | इन्द्रकक्ष | तततत | 11 | 4 | UUUU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU | UUUU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU |
| | 5 | उपेन्द्रकक्ष | जजजज | 11 | 4 | UUUU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU | UUUU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU |

Fig. 8 Sama Vrutta Database

| S.N O | NAME OF THE VRUTTA | NO. OF SYLLABLES IN 1st and 3rd PADA | NO. OF SYLLABLES IN 2nd and 4th PADA | GANNA OF 1ST AND 3RD PADA | GANNA OF 2ND AND 4TH PADA | NO. OF PADAS | PATTERN | PATTERN |
|-------|--------------------------|--------------------------------------|--------------------------------------|---------------------------|---------------------------|--------------|---|---|
| 1 | हरिणीवृत्ता | 11 | 12 | ससस+ IU | नसस | 4 | UUU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU | UUU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU |
| 2 | औपच्यन्दसिकः पुष्पिलम्बा | 12 | 13 | ननरय | नजजर+U | 4 | UUU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU | UUU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU |
| 3 | त्रिवारिणी | 10 | 11 | ससजग | ससस+U | 4 | UUU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU | UUU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU |
| 4 | अपरकक्ष | 11 | 12 | ननरजग | नजजर | 4 | UUU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU | UUU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU |
| 5 | उपचित्र | 8 | 10 | सससजग | सससजग | 4 | UUU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU | UUU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU |
| 6 | वेगवती | 7 | 11 | सससजग | सससजग | 4 | UUU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU | UUU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU |
| 7 | वसन्ततिलका | 9 | 11 | ससजग | सससजग | 4 | UUU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU | UUU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU UU |

Fig. 9 Ardha-Sama Vrutta Database

2) By having the above data set we will compare P(I) with the predicted match pattern. If any of the patterns is similar to the existing patterns then it will identify the missing pattern.

Input pattern:

UU A UIUUIU BUU CU UIUUIU

Existing pattern which is similar:

UU UU UIUUI UU UU UU UIUUIU

Output: Identified pattern for deteriorated zone: A= UU, B= U, C= UU

- [4]. Natarajan Meghanathan et al.:WiMONE, NCS, SPM CSEIT-2014 pp.199-207, 2014 doi:10.5121/csit.2014.412216
- [5]. Internet resource:<https://www.researchgate.net/publications/265644163>
- [6]. Roger L., Easton.Jr, Keith.T.Knox, "Digital Restoration of Erased and Damaged Manuscripts", Proceedings of the 39th Annual Convention of the Association of Jewish Libraries (Brooklyn, NY-June 20-23,2004)
- [7]. <https://sanskrit.iitk.ac.in/jnanasangraha/chanda/verse>
- [8]. Pingala, Chandas-Sastra, Kavyamala series no.91 (3rd Edition) Bombay, 1938.
- [9]. RAMA N. AND Meenaksh Lakshmanan, A Computational Algorithm for Classification of Verse, IJCSI International Journal Computer Science Issues, Vol.7, Issue 2, No.1 March 2010, ISSN:1694-078