



Word Sense Disambiguation for Telugu based on Telugu WordNet

Srinivas Mulkalapalli¹, Dr. B. Padmaja Rani², Dr. G. Nagaraju³

¹Associate Professor, CSE, Geethanjali College of Engineering and Technology, Telangana, India

²Professor, CSE, College of Engineering, JNTUH, Telangana, India

³Assistant Professor, CSE(AIML&IOT), VNRVIJET, Telangana, India

Abstract: Telugu (తెలుగు) is a Dravidian language which is an official language in Telangana and Andhra Pradesh the states of India. It is also spoken by many people of other states for communication in India. Like all other languages, there are many ambiguous words in Telugu. Determining the implied meaning of an ambiguous word is essential for various Natural Language Processing (NLP) tasks namely Information Retrieval, Machine Translation, and Text Summarization. In this paper, we have proposed Word Sense Disambiguation (WSD) algorithms for Telugu language that uses Telugu WordNet using Adapted Lesk's algorithm, and an extension of Adapted Lesk's algorithm based on one measure of semantic relatedness- Leacock – Chodorow. Adapted Lesk's algorithm determines the most appropriate sense of an ambiguous word based on the highest overlapping of the sense's definition in the gloss with the context words definitions including glosses and example sentences. In the second algorithm, we are using Leacock – Chodorow semantic relatedness measure that is computed by utilizing the distance between noun concepts in an IS-A hierarchy. Evaluation is performed on the dataset prepared by us that consists of 30 polysemous Telugu noun words and compared the results. We obtained an accuracy of 72% and 78% respectively for these algorithms.

Keywords: Natural Language Processing, WordNet, Word Sense Disambiguation, ambiguity, Highest overlap, semantic relatedness measure, Telugu.

Introduction:

In natural language processing, resolving ambiguity is one important research problem. Ambiguous words, having multiple senses or meaning, is prevalent in almost every language used by human race to communicate their opinion or ideas. Telugu language, spoken in South India, also has many ambiguous words. Human beings are very skillful in resolving ambiguity if arises during their communication using their world knowledge. But, for computers the task of finding the most appropriate sense of an ambiguous word in a sentence is very difficult. It is even considered as one of AI Complete problem. The task of determining the most appropriate sense of an ambiguous word implied through the context in a sentence or domain is known as Word sense Disambiguation (WSD). WSD is one of the important research areas where development of good systems helps to improve the performance of several NLP applications like Machine Translation, Information Retrieval, and similarity analysis etc. Substantial research was



completed and in continuation for English and other European languages because of the availability of required lexical resources such as WordNet [1], Corpus [2] and Golden Dataset for Evaluation [3] and very less work has been done for Telugu language [4][5]. The research on WSD for Indian languages such as Hindi, Tamil, Bengali, Telugu, and Malayalam... etc. is hampered due to unavailability of well-established lexical resources such as WordNet, Corpus etc. But, Indo-WordNet [6] is a lexical database developed for the Indian languages by Pushpak Bhattacharya et al. and providing access through <http://www.cfilt.iitb.ac.in/indowordnet/index.jsp>.

Telugu language also has many ambiguous words. Recently, with the development of ICT many NLP Tasks such as Sentiment Analysis, Text Classification, Question- Answering Systems, Information Retrieval, Machine Translation etc. for Telugu are also in need. WSD for Telugu is not matured enough to be useful in the above mentioned tasks. But, WSD for other languages like English, German, French, Italy, Chinese is well developed because of sophisticated lexical resources and corpus. So, we are interested in developing good WSD systems for Telugu in hope of being useful in several NLP tasks.

1. RELATED WORK

Lesk's algorithm [7] is one in determining the most appropriate sense of an ambiguous word implied from its context. It proceeds as follows: A separate bag of words is created for every possible sense of the target word from its definition found in Machine Readable Dictionaries. They have used Oxford Advanced Learner's Dictionary of English. Similarly, a context bag is created from all other words surrounding the target word taking their definitions. The algorithm determines the number of common words between each sense bag and the context bag. Finally, the sense whose sense bag has highest common words is selected as the appropriate meaning. Later, glosses from WordNet are used for Word Sense Disambiguation in [8]. They have achieved an accuracy of 32% by evaluating their algorithm through participation in SENSEVAL-2, which involves an evaluation exercise on English lexical sample data.

Knowledge based WSD for Hindi is proposed in [9]. Prity Bala has evaluated the proposed system on a dataset of 100 ambiguous words and got 50% accuracy. Alok Ranjan Pal et al in [10] proposed knowledge based WSD using Bengali WordNet. They shared the results obtained by conducting experiments on dataset containing 9 frequently used Bengali ambiguous words.



Dependency parsing [11] and word embeddings [12] are also plays crucial role for resolving disambiguate for words. A knowledge-based approach to WSD exploiting Structural Semantic Interconnection is presented in [11] by Roberto Navigli.

WSD based on conceptual density is proposed in [13] by Agirre and Rigau. Conceptual density is computed from is-a hierarchy from the WordNet. C. Leacock and M. Chodorow proposed a WSD based on combining local context and WordNet similarity in [15]. D. Lin explored the possibility of using syntactic dependency to improve WSD in [16]. Semantic distance between topics in WordNet is used for Word Sense Disambiguation by Sussan in [17]. Sussan proposed a weighting scheme based on WordNet relations. The synonymy relation is assigned a weight of zero, whereas the weights in the range [1, 2] are assigned to other relations hypernymy, hyponymy, holonymy and meronymy. WSD for Hindi based on measure of Semantic Relatedness is proposed by Satyendr Singh et al in [18]. They have achieved 60.65% as an overall average accuracy by conducting an experiment on a sense tagged dataset prepared by them which consists of several instances for 20 polysemous Hindi nouns. Knowledge based WSD for Telugu is proposed in [19] by Suneetha Eluri and Vishala Siddu. They have achieved an accuracy of 65.4 by evaluating on Telugu sentences formed for 150 ambiguous words both nouns and verbs.

3. WordNet

Traditional dictionaries arrange the words in alphabetical order making it difficult to access similar words during WSD. WordNet arranges the words semantically. It is an electronic lexical database that consists of nouns, verbs, adjectives, and adverbs. It is created in 1990 at Princeton University. It groups synonymous words together to form synsets or synonym sets. A word having multiple senses known as polysemous naturally appears in multiple synsets. For example, line occurs in 4 noun synsets {argumentation, logical argument, argument, line of reasoning}, {telephone line, phone line, telephone circuit, subscriber line}, {occupation, business, job, line of work}, and {note, short letter, billet} and the verb synset {trace, draw, describe, delineate}. WordNet3.0 has 117798 nouns, 11529 verbs, 21479 adjectives, and 4481 with a total of 155287 words.

WordNet has a definition gloss for every synset. This contains a short description of the meaning of the corresponding synset. The gloss of the synset {argumentation, logical argument, argument, line of reasoning} is “a course of reasoning aimed at demonstrating a truth or falsehood; the methodical process of logical reasoning”. WordNet provides a sense tag for each synset that is a unique identifier.

Semantic relations exist between the synsets to connect them. Synsets of one part of speech are mostly connected to synsets of the same part of speech. Hyponymy and hypernymy are two most important relations for nouns. Suppose that synset S_1 is a kind of synset S_2 , then we say S_1 is the



hyponymy of S_2 , and S_2 is the hypernymy of S_1 . Similarly, hypernymy and troponymy is one pair of related relations for verbs. If S_2 is one way to S_1 , then S_1 is the hypernymy of S_2 , and S_2 is the troponym of S_1 . Attribute is one relation for adjectives. It relates adjective to noun.

WordNet employs inheritance relation, one of the principles in OOP. So, WordNet allows us to start at any node and permits to ascend up or descent down to determine the broader or narrower meanings to use them in variety of ways.

The Telugu WordNet is a component of IndoWordNet developed by Indian Institute of Technology (IIT), Bombay. It is developed using expansion approach in the lines of Princeton English WordNet. This lexical resource consists of Gloss, Example sentence, Synset, POS tag, various lexical and semantic relations of the word as information. Latest Telugu WordNet consists of 12078 nouns, 2795 verbs, 5776 adjectives, and 442 adverbs, a total of 21091 words. But, still it needs to cover more words.

We are sharing some part of the information retrieved for the Telugu word: అడుగు from Telugu WordNet

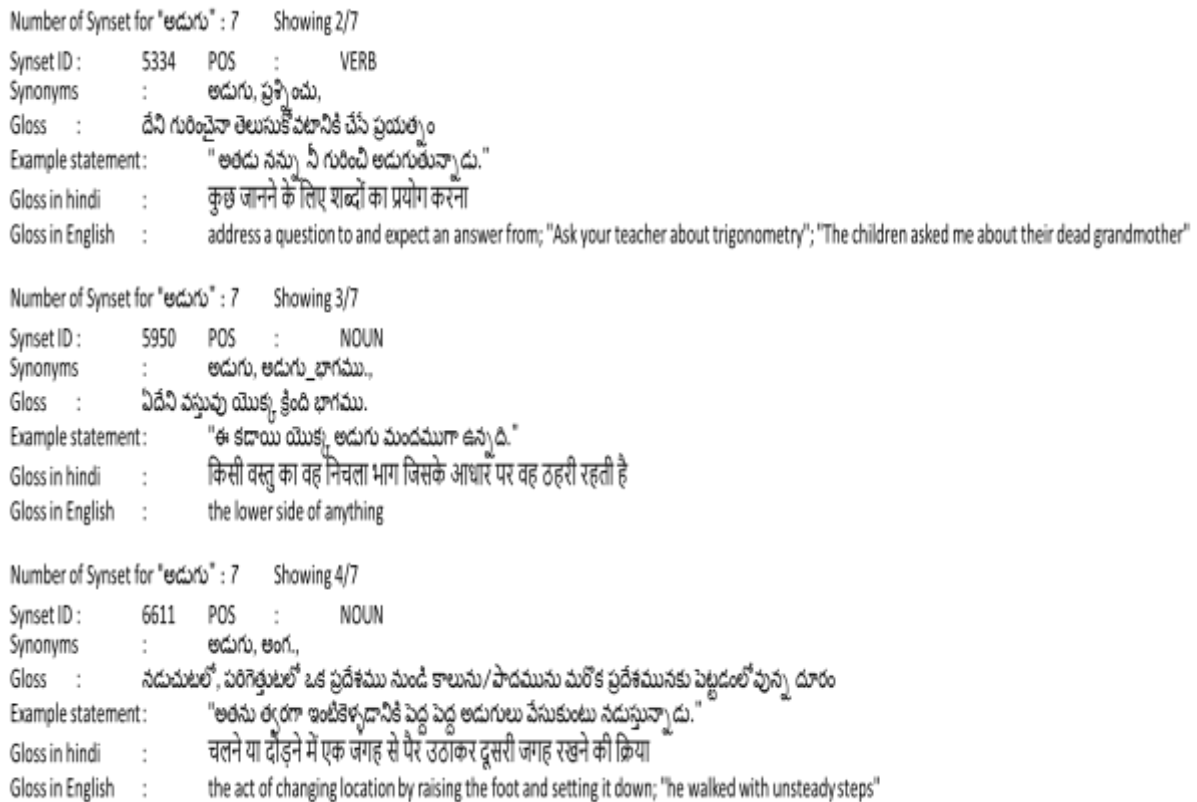


Fig 1: Part of Synset Id's and Information of అడుగు in Telugu WordNet



Showing Hypernymy			
Hypernymy ID	Synonyms	Gloss	Example
13376	మాట్లాడు	ఒక విషయాన్ని గురించి చర్చించుకోవడం	"పిల్లలు రామ్-రామ్ అని అందున్నారు"
246	వెళ్లడంవేయి, ప్రకటించు	మనస్సులోని మాటలను బయటకు చెప్పుట.	"అతను తన అభిప్రాయాలను వెళ్లడం చేశాడు."
5474	తెలపడము, తెలియజేయడము	ఏదేని వస్తువు, సూచన మొదలగు వాటి గురించి పరిచితులను చేయుట.	"అమె నాకు ఈ విషయము గురించి ముందుగానే తెలియజేసింది."
7302	ఆర్థమయ్యేటట్లు చెప్పునచ్చవెప్పు	బోధించు లేక జ్ఞానం కలిగించుట.	"ఆధ్యాపకుడు పిల్లలకు గణితాన్ని ఆర్థమయ్యేటట్లు వివరిస్తున్నారు."
3295	పనిచేయు	ఖాళీగా లేకుండా ఏదోకటి చేయడం	"ఈ పని చేసిన తర్వాత మీ పని చేస్తాను."

Fig 2: Hypernymy Information of అడుగు in Telugu WordNet

4. The Disambiguation Methodology

4.1 Adapted Lesk Algorithm

We proposed use of extended-sense-overlap between all possible senses of an ambiguous word and the context to determine the appropriate sense of an ambiguous word in the given sentence. Telugu WordNet is lacking the required information as it is in its early stage of development. To compensate for this scarcity of information, we are using the following features: we are not using any fixed size context window, rather including every word in the sentence after removing the stop words. Glosses, Synonymous words, and example sentences for the resulting words are retrieved from the Telugu WordNet and used in our algorithm.

Algorithm: Adapted Lesk Algorithm for Telugu Word Sense Disambiguation

Input: Telugu sentence that consists of an ambiguous word w to be disambiguated.

Output: Appropriate sense of w implied from the context.

Procedure:

Step 1: Tokenization is performed on the sentence to get the list of words and stop words are removed if any exists.

Step 2: Lemmatization is performed to get the meaningful root words.

Step 3: We form a context bag C that consists of Gloss and Example sentence for every word in the list which are retrieved from the WordNet.

Step 4: We form a set of strings S_i , $i = 1, 2, \dots, n$, where n is the number of senses of an ambiguous word as per the WordNet.

Step 5: $\max_overlap = 0, j = 1;$

Step 6: For $i : = 1$ to n do



Compute the overlap between C and S_i.

If (max_overlap < overlap) then

Max_overlap = overlap; j = i;

Step 7: Return Sense_j as the result of our WSD System.

4.2 Word Sense Disambiguation for Telugu based on Measure of Semantic Relatedness

Gloss overlaps may also be considered as another possible measure of semantic relatedness. So, we propose a WSD Algorithm for Telugu involving one of the existing measures of semantic relatedness as an extension of Adapted Lesk algorithm. The value of semantic relatedness to be assigned to a pair of concepts is determined using their relative position in a concept hierarchy. The concept hierarchies available from the WordNet are used in defining various semantic relatedness measures such as Leacock- Chodorow measure, Resnik measure, Jiang-Conrath measure, Lin measure, and the Hirst-St. Onge measure.

We are using Leacock – Chodorow measure of semantic relatedness in our proposed Word Sense Disambiguation Algorithm for Telugu. It is computed using path lengths in an is-a hierarchy of noun concepts. The path which involves least number of intermediate concepts is the shortest path between two concepts. It is divided by the value equals the double of the depth D of the hierarchy, which is the length of the longest path from the root node to a leaf node in the hierarchy. So, the measure of relatedness is

$$\text{related}_{\text{Ich}}(w_1, w_2) = [-\log(\text{ShortestLength}(w_1, w_2) / (2D))]]$$

Algorithm: WSD for Telugu using measure of Semantic Relatedness - Leacock – Chodorow

Input: Telugu sentence that consists of an ambiguous word w to be disambiguated.

Output: Appropriate sense of w implied from the context.

Procedure:

Step 1: Remove stop words from the given sentence, and create a context vector C with the nouns.

Step 2: Let Syn_Input = Set of all Synset_Id's of the target word, extracted from the WordNet.

Syn_Output = Set of all Synset_Id's of the nouns in the context vector C.



Step 3: Suppose that $n := \text{len}(\text{Syn_Input})$, $m := \text{len}(\text{Syn_Output})$, and $\text{sense_score} := 0$, $\text{result_sense} = \text{Syn_Inut}[1]$;

Step 4: For $i := 1$ to n do

```
{
    rs :=0
    For j := 1 to m do
        { rs+ = [ -log(ShortestLength(Syn_Input[i], Syn_Output[j]) / (2D))];
            if ( rs > sene_score ) then
                {
                    sense_score = rs;
                    result_sense = Syn_Inut[i];
                }
        }
    } }
```

Step 5: Return result_sense ;

5. Data Set and Results

We have created a Telugu data set of sentences for 30 most frequently used ambiguous words from Telugu language. We have collected sentences from Internet, Telugu news, and Telugu books. It consists of 590 sentences with 5 sentences for each sense. The translation and transliteration of the data set is given in Table A1 in Appendix.

The Evaluation measure of Accuracy used to evaluate our algorithms.

$$\text{Accuracy} = \frac{\# \text{ Test Instances correctly disambiguated}}{\# \text{ Test Instances}}$$

Method	Accuracy %
Adapted Lesk WSD	72%
WSD based on Leacock –Chodorow Measure	78%

6. Conclusion and Future Work:



In this paper, we have proposed two algorithms for Word Sense Disambiguation of Telugu language and compared the results. Due to the scarcity of required information in the WordNet these algorithms were tested on limited data set. In future, we want to develop Unsupervised WSD System based on Word Embeddings to avoid the dependence on lexical resources and free from knowledge acquisition bottleneck.

References

- [1]. G. A. Miller, R. Beckwith, C. Fellbaum, D. Gross and K. J. Miller, "WordNet An on-line lexical database," *International Journal of Lexicography*, Vol. 3, No. 4, pp. 235-244, 1990
- [2]. https://www1.essex.ac.uk/linguistics/external/clmt/w3c/corpus_ling/content/corpora/list/private/brown/brown.html
- [3]. Gaizauskas, "Gold Standard Datasets for Evaluating Word Sense Disambiguation Programs. *Computer Speech and Language*, Vol. 12, No. 3, Special Issue on Evaluation of Speech and Language Technology, pp. 453-472, 2009 .
- [4]. G. Nagaraju, N. Mangathayaru, B. Padmajarani. "Transition based parser for Telugu language". *International Journal of Engineering and Technology*, Vol.7, No.4, pp: 4674-4677, 2018
- [5]. G. Nagaraju, N. Mangathayaru, B. Padmaja Rani. "Dependency Parser for Telugu language". In *Proceedings of the Second International Conference on Information and Communication Technology for Competitive Strategies*, Article No. 138, 2016. <https://doi.org/10.1145/2905055.2905354>.
- [6]. Pushpak Bhattacharyya, "IndoWordnet" Department of Computer Science and Engineering Indian Institute of Technology Bombay <http://www.cfil.t.iitb.ac.in/indowordnet/index.jsp>
- [7]. M. Lesk, "Automatic Sense Disambiguation Using Machine Readable Dictionaries: How to Tell a Pine Cone from an Ice Cream Cone," *Proceedings of SIGDOC*, 1986.
- [8]. S. Banerjee and T. Pedersen, "An adapted Lesk algorithm for word sense disambiguation using WordNet," In *Proceedings of the Third International Conference on Intelligent Text Processing and Computational Linguistics*, Mexico City, February 2002
- [9]. Prity Bala, "Knowledge Based Approach For Word Sense Disambiguation Using Hindi Wordnet," *The International Journal Of Engineering And Science (IJES)*, Volume 2 Issue 4, PP. 36-41, 2013.
- [10]. Alok Ranjan Pal, DigantaSaha, and Sudip Kumar Naskar, " Word Sense Disambiguation in Bengali: a Knowledge based Approach using Bengali WordNet" 2017 IEEE
- [11]. G Nagaraju, N Mangathayaru, B Padmaja Rani. "MST Parser for Telugu Language". *Proceedings of the Third International Conference on Computational Intelligence and Informatics*, pp: 271-279, 2020
- [12]. G. Nagaraju, N. Mangathayaru, B. Padmajarani "Integrating Transition and Graph Based Dependency Parsers Using Ensembled and Stacking Approaches for Parsing Telugu Language", *International Journal of Advanced Research in Engineering and Technology (IJARET)*, Vol.12, No.2, PP.96-105, Feb 2021
- [13]. Roberto NAVIGLI and Paola VELARDI, "Structural Semantic Interconnection: a knowledge-based approach to Word Sense Disambiguation. SENSEVAL-3: Third



- International Workshop on the Evaluation of Systems for the Semantic Analysis of Text,” Barcelona, Spain, Association for Computational Linguistics, 2004.
- [14]. E. Agirre and G. Rigau. Word sense disambiguation using conceptual density. In *Proceedings of the 16th International Conference on Computational Linguistics*, pages 16–22, Copenhagen, 1996.
- [15]. C. Leacock and M. Chodorow. Combining local context and WordNet similarity for word sense identification. In C. Fellbaum, editor, *WordNet: An electronic lexical database*, pages 265–283. MIT Press, 1998.
- [16]. D. Lin. Using syntactic dependency as a local context to resolve word sense ambiguity. In *Proceedings of the 35th Annual Meeting of the Association for Computational Linguistics*, pages 64–71, Madrid, July 1997.
- [17]. Patwardhan, S., Banerjee, S., Pedersen, T. (2003). Using Measures of Semantic Relatedness for Word Sense Disambiguation. In: Gelbukh, A. (eds) *Computational Linguistics and Intelligent Text Processing. CICLing 2003. Lecture Notes in Computer Science*, vol 2588. Springer, Berlin, Heidelberg.
- [18]. Singh, S., Singh, V. K. and Siddiqui, T. J.: Hindi Word Sense Disambiguation using Semantic Relatedness measure. In *Proceedings of 7th Multi -Disciplinary workshop on Artificial Intelligence*, Krabi, Thailand, pp. 247-256 (2013)
- [19]. Suneetha Eluru, Vishala Siddu, “A Knowledge Based Word Sense Disambiguation in Telugu Language”, *International Journal of Engineering and Advanced Technology (IJEAT)* ISSN: 2249-8958 (Online), Volume-10 Issue-1, October 2020



Appendix:

Table A1. Translation, Transliteration and details of Telugu Ambiguous Words

Word	No. of Senses	Sense Number : Translation of senses in English
అడుగు (adugu)	2	Senes1: The foot step, foot print Senes2: The bottom, basis
ఉత్తరం (uttaram)	3	Senes1: An answer, a reply Senes2: A letter Senes3: The North direction
కమ్మ (kamma)	4	Senes1: A kind of ear ornament worn by women Senes2: A letter written upon a palm leaf Senes3: A certain caste Senes4: A branch or bough of any species of palm tree
కారు (kaaru)	4	Senes1: Season, time of the year Senes2: A forest Senes3: Black, dark color Senes4: Tongs
గంట (ganta)	3	Senes1: A bell, a gong Senes2: A stub, the stump of a corn -stalk Senes3: 60 minutes of time
గుంట (gunta)	4	Senes1: A pond or tank Senes2: A pit Senes3: A certain square measure of land Senes4: A little girl, a brat
గుడి (gudi)	3	Senes1: A temple Senes2: A circle, a halo round the Sun or Moon Senes3: Circular mark added to a consonant on the top
గురువు (guruvu)	3	Senes1: A teacher Senes2: The planet Senes3: A long syllable
ఘనము (ghanamu)	5	Senes1: Greatness, dignity, honor Senes2: A cloud Senes3: A cymbal, a bell, a gong Senes4: A cube Sense5: The cube of a number
తంత్రము (tantramu)	5	Senes1: A device, contrivance, means



(thantramu)		Senes2: A trick, A stratagem Senes3: The regular order of ceremonies Senes4: A scientific work or treatise Senes5: A doctrine, rule, theory
తార (thara)	4	Senes1: A star Senes2: The pupil of the eye Senes3: A high tone or note in music Senes4: Name of a women
దండము (dhandamu)	6	Senes1: Salutation Senes2: A rod, stick Senes3: A measure of length equal to 4 Hasthas Senes4: Punishment Sense5: An army Sense6: A collection
పక్షము (pakshamu)	3	Senes1: A side Senes2: A wing, a feather Senes3: Fortnight
పర్వము (parvamu)	5	Senes1: A joint or knot Senes2: A division or section of book Senes3: A festival Senes4: A name given to certain days in the lunar month Sense5: An opportunity
పాదము (paadamu)	4	Senes1: The foot Senes2: A line in a stanza or a verse Senes3: A quarter or fourth part Senes4: A Root
పురి ^o (puri)	4	Senes1: A town, city Senes2: Pack-thread, twist Senes3: A peacock's tail Senes4: A straw basket in which seed grain is preserved
పూజ్యము (poojyamu)	2	Senes1: Blank, empty Senes2: Reverence, respect
ప్రమాణము (pramanamu)	4	Senes1: Measure, size, dimension Senes2: Standard Senes3: Oath, swearing Senes4: Rule, sanction, authority, ground



బంతి (banthi)	3	Senes1: A ball Senes2: A row, line, rank Senes3: Marigold – flower
భంగము (bhangamu)	4	Senes1: Breaking, splitting Senes2: Failure, loss Senes3: Prevention Senes4: Disappointment
మండలము (mandalamu)	6	Senes1: A disk of the sun or moon Senes2: A circle, a wheel, a ring, circumference Senes3: A district, province Senes4: A group, assemblage Senes5: An association, society, company Senes6: A division of the Rigveda
మంత్రోరము (manthramu)	5	Senes1: A charm, incarnation Senes2: Consulting Senes3: Counsel, advice Senes4: Plan, design Senes5: A sacred text or hymn
రసము (rasamu)	5	Senes1: Juice, sap, exudation Senes2: water Senes3: Mercury Senes4: Taste, relish, savor Senes5: Literary or artistic beauty
లవము (lavamu)	5	Senes1: A small quantity Senes2: Half a second Senes3: A degree Senes4: The numerator of a fraction Senes5: Loss, destruction
వర్గము (vargamu)	3	Senes1: A class, a tribe, race Senes2: The square of number or quantity Senes3: A group, series, set, a multitude of similar things
వాసి (vaasi)	4	Senes1: Difference, comparison Senes2: Superiority Senes3: Benefit, advantage, gain Senes4: A dweller, inhabitant, resident
వోధి (vidhi)	4	Senes1: A rule, sacred precept Senes2: Duty Senes3: Name of Brahma



		Sense 4: Fate, destiny, luck
శిక్ష (shiksha)	3	Senes1: Teaching, training Senes2: Punishment Senes3: One of 6 Vedangas
సంధు (sandhu)	5	Senes1: A narrow street, lane Senes2: A nook, a corner Senes3: An opening fissure, crack, hole Senes4: An interval, an intervening space Senes5: An opportunity
హారము (haramu)	3	Senes1: A garland Senes2: Divisor, denominator Senes3: Removing, deprivation, loss