## Final phase Report of the project on "Computing Tools for Tamil Language teaching and learning"

Dr. Anand Kumar M, Rajendran S, Soman K P

CEN, Amrita Vishwa Vidyapeetham, Coimbatore

## I. Tool for learning Verb Conjugation and Noun Declension

A Tool for learning Verb Conjugation and Noun Declension is getting developed as one of the components of the present project. Tamil is a morphologically rich language. Being agglutinative language most of the grammatical information are expressed by suffixes. For example, nouns are inflected for number and cases and verb are inflected for tenses, moods and aspects and subject agreement markers. A morphological generator capturing conjugation of verbs have been developed as one of components of the present project. The Morphological Generator takes lemma and grammatical information as input and gives inflected forms of the given word. It is a reverse process of Morphological Analyzer. Morphological generator system implemented here is a rule based system which makes use of morpheme concatenating rules and gives us the all the conjugated forms of a given verb and declension forms of a given noun. In the first phase only "verb conjugation" has been completed. We will take up noun declension in the second phase. We hope to develop a full-fledged tool which helps Tamil learners in understanding verb conjugation and noun declension in Tamil.
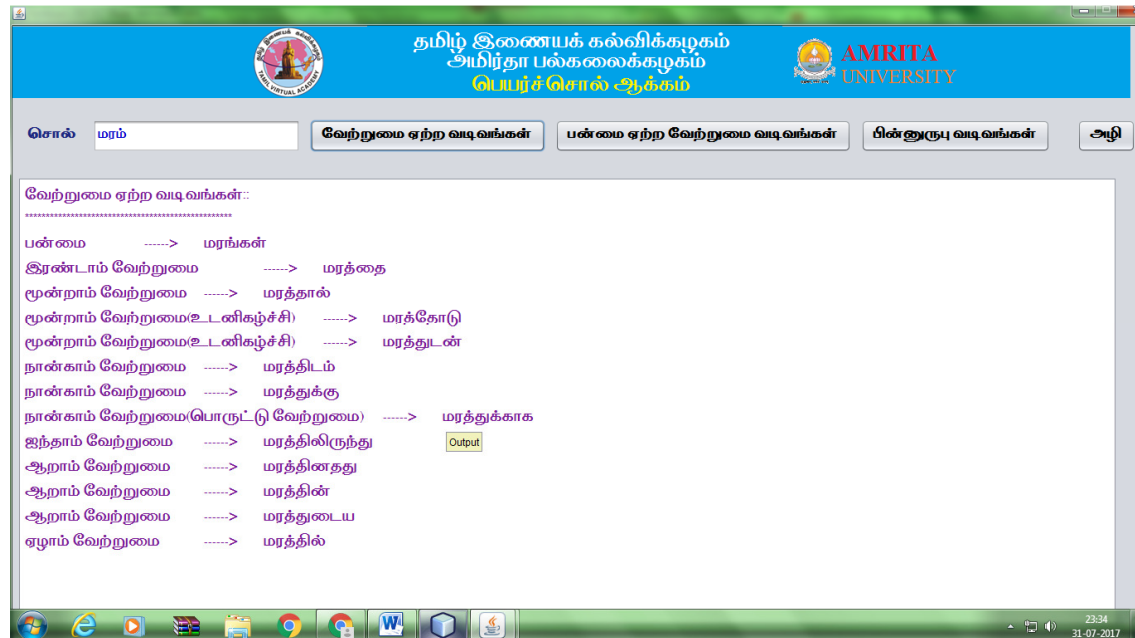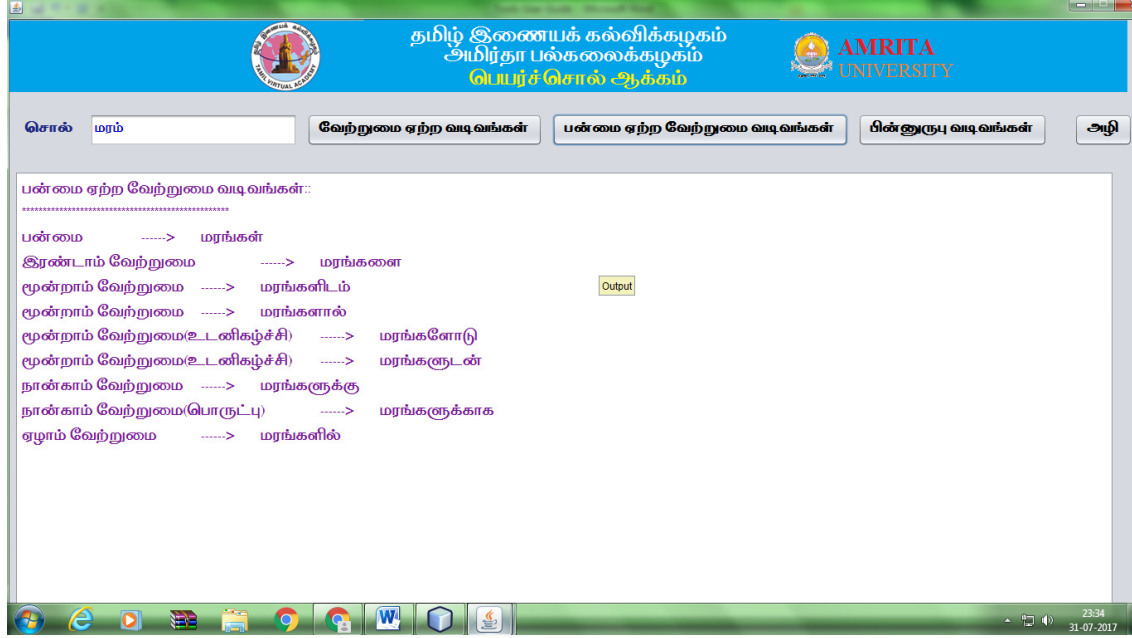


Fig 1.1பெயர்ச்சொல்ஆக்கம்- வேற்றுமைஏற்ற வடிவங்கள்
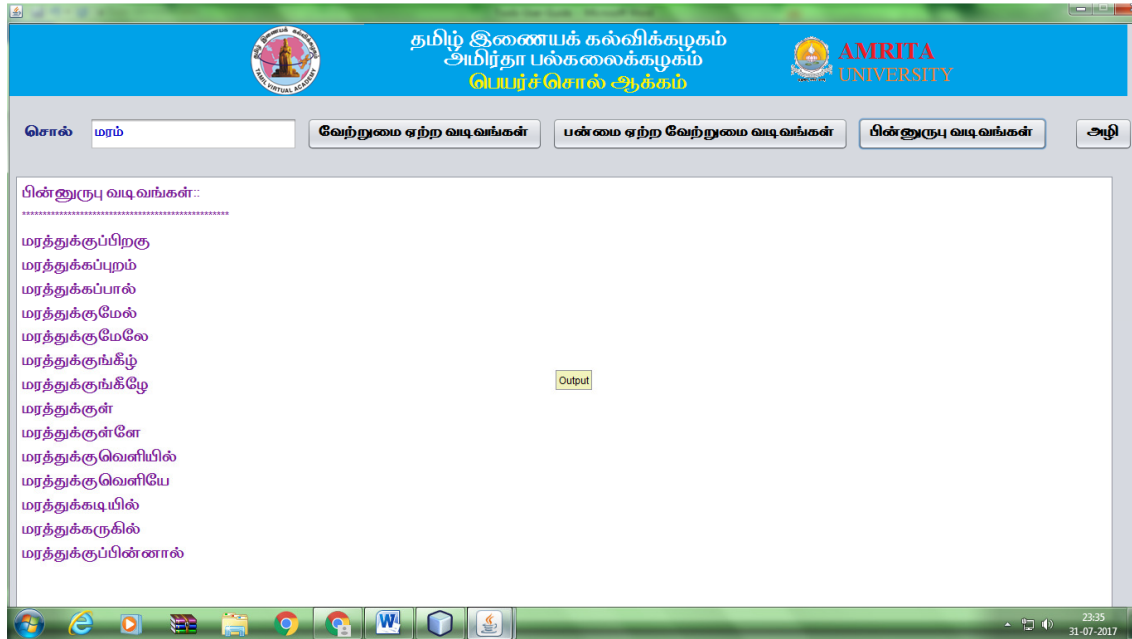
Fig 1.2பெயர்ச்சொல்ஆக்கம்-  பன்மைஏற்றவேற்றுமை வடிவங்கள்



Fig 1.3பெயர்ச்சொல்ஆக்கம்-  பின்னுருபு வடிவங்கள்
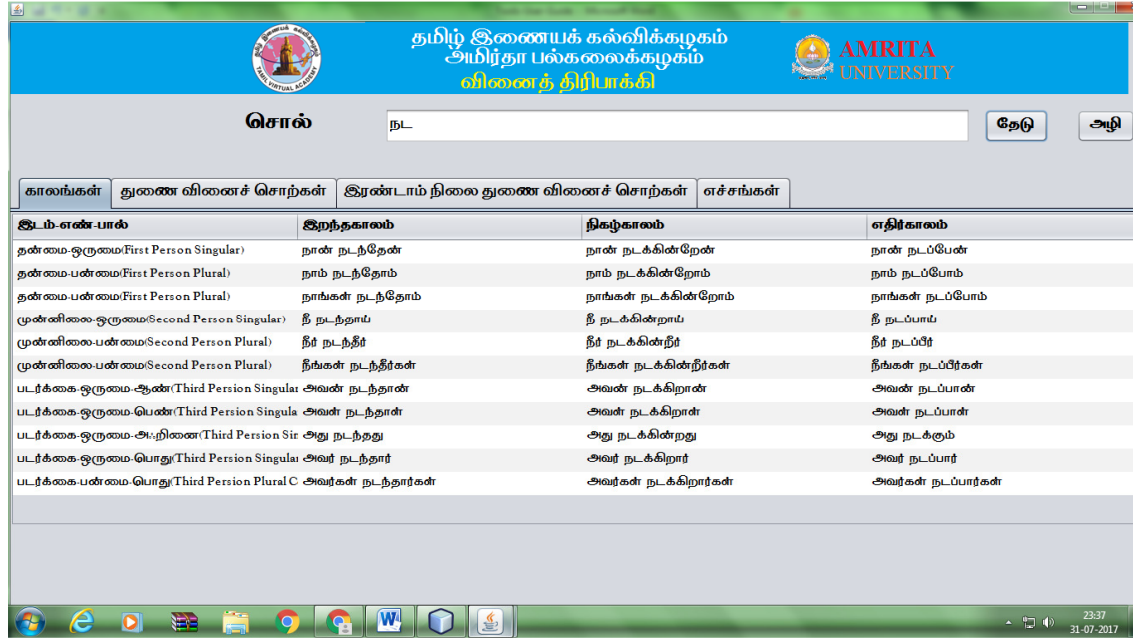
சொல் | நட_ | தேடு | அழி

காலங்கள் | துணை வினைச் சொற்கள் | இரண்டாம் நிலை துணை வினைச் சொற்கள் | எச்சங்கள்

| இடம்-எண்-பால் | இறந்தகாலம் | நிகழ்காலம் | எதிர்காலம் |
| --- | --- | --- | --- |
| தன்மை-ஒருமை(First Person Singular) | நான் நடந்தேன் | நான் நடக்கின்றேன் | நான் நடப்பேன் |
| தன்மை-பன்மை(First Person Plural) | நாம் நடந்தோம் | நாம் நடக்கின்றோம் | நாம் நடப்போம் |
| தன்மை-பன்மை(First Person Plural) | நாங்கள் நடந்தோம் | நாங்கள் நடக்கின்றோம் | நாங்கள் நடப்போம் |
| முன்னிலை-ஒருமை(Second Person Singular) | நீ நடந்தாய் | நீ நடக்கின்றாய் | நீ நடப்பாய் |
| முன்னிலை-பன்மை(Second Person Plural) | நீர் நடந்தீர் | நீர் நடக்கின்றீர் | நீர் நடப்பீர் |
| முன்னிலை-பன்மை(Second Person Plural) | நீங்கள் நடந்தீர்கள் | நீங்கள் நடக்கின்றீர்கள் | நீங்கள் நடப்பீர்கள் |
| படர்க்கை-ஒருமை-ஆண்(Third Persion Singular | அவன் நடந்தான் | அவன் நடக்கிறான் | அவன் நடப்பான் |
| படர்க்கை-ஒருமை-பெண்(Third Persion Singula | அவள் நடந்தாள் | அவள் நடக்கிறாள் | அவள் நடப்பாள் |
| படர்க்கை-ஒருமை-அஃறிணை(Third Persion Sir | அது நடந்தது | அது நடக்கின்றது | அது நடக்கும் |
| படர்க்கை-ஒருமை-பொது(Third Persion Singular | அவர் நடந்தார் | அவர் நடக்கிறார் | அவர் நடப்பார் |
| படர்க்கை-பன்மை-பொது(Third Persion Plural C | அவர்கள் நடந்தார்கள் | அவர்கள் நடக்கிறார்கள் | அவர்கள் நடப்பார்கள் |

Fig 1.4வினைத்திரிபாக்கி

## II. Tool for developing vocabulary skill using Ontology of Tamil vocabulary

Tamil Onto-thesaurus is an outcome of a very long research activity went on in the field of lexical semantics of Tamil vocabulary. It went through several stages before being culminated into Tamil onto-thesaurus. It depicts our travel from Tamil thesaurus to Tamil word net. It is a lexical resource which amalgamates all sorts of information available in a dictionary, thesaurus and word net.A paper thesaurus for Tamil was prepared in 1990 based on the principles of componential analysis of meaning propounded by Nida (1975) and was published in 2001 (Rajendran, 2001), nearly after a decade. Following the paper thesaurus, an Electronic thesaurus for Tamil was attempted and a book on Tamil electronic thesaurus was published in 2006 (Rajendran and Baskaran, 2006) The preparation of wordNet for Tamil was undertaken (2001-2003) with the financial assistance from Tamil Virtual University (renamed now as Tamil virtual academy) and a crude version of it based on the ontology developed by Rajendran (Rajendran, 2001) was submitted to the institute in 2003. After that, from 2009 onwards with the fund received from MHRD and Department of electronics and information Technology of Govt. of India the building of Dravidian wordNet was executed based on Hindi wordNet; nearly 3000 synsets (concepts) have been completed. Still we have to go a long way to achieve the desired target. At present a team from CEN, Amrita University is involved in building onto-thesaurus for Tamil as a part of the project entitled "Computing Tools for Tamil Language teaching and learning".The project is funded by Tamil Virtual Academy, Chennai.

**Tamil Onto-thesaurus**

Thesaurus is a in its wider sense is a classification of words by concepts, topics, or subjects. The present Tamil Onto-thesaurus is the extended version of Electronic thesaurus of Tamil focusing more on the ontological features. Two kinds of issues arise in the preparation of Tamil onto-thesaurus:

- Linguistic issues
- Computational issues

**Linguistic Issues**

It involves mainly the following four tasks:

       1. Developing an ontology for Tamil based on structural semantic principles.

       2. Establishing semantic domains and sub domains based on distinguishing semantic or Componential features of lexical items.

       3. Classifying Tamil vocabulary to fit into the ontology developed.

       4. Linking words by various semantic or lexical relations such as synonymy, hyponymy-hyperonymy, meronymy-holonymy, compatibility, and incompatibiliity.

**Computational Issues**

       It involves mainly the following three tasks:

1. Conversion of linguistic data base into computer accessible format.

2. Preparation of a tool to provide the facilities for augmenting, entering and editing the raw data, and classifying the lexical items in a semi-automatic way.

3. Creation of user friendly interfaces for accessing the onto-thesaurus in simple manner.

**Ontology of Tamil vocabulary**

The ontology available in Rajendran (2001), which is founded on the theory of componential analysis of meaning propounded by Nida (1975) is enhanced to suit the present purpose. The following is the skeletal structure of the Tamil ontology adopted in Onto-thesaurus.

*Ontology Relationships*

    Relationship between nodes in the ontology tree has been classified in to two types namely IS-A relationship and
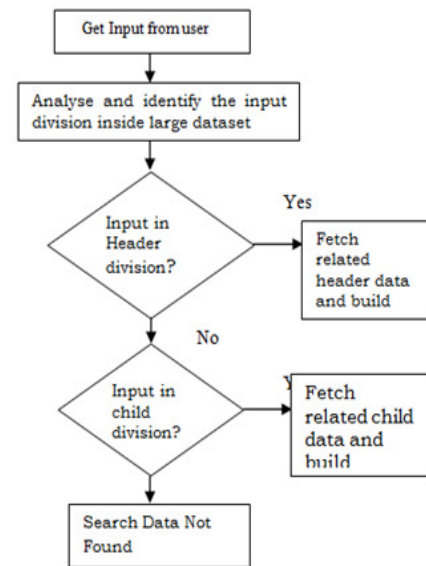
HAS-A relationship.

Both of the relationships are frequently used in hierarchical to show the link between the nodes.

- *IS-A Relation*

  This types specifies the relationship between parent nodes and child nodes in the ontology tree. *Eg. Crow is a Bird, Dog is a Animal.* Here using this is-a relationship, one can easily identify Crow is a child node of Bird node and Dog is a child node of Animal node.

- *HAS-A Relation*

  It shows the properties of a node or entity. Each node has its own properties and it can be expressed with has-a relationship notations. E.g.: Crow has wings.

- *Sibling Relation*

  When a node entity A has a direct child B, then B is sibling for the entity A. i.e. Sibling refers t a direct child for a node.

- *Transitive Relation*

  Transitivity between the nodes also possible in the ontology tree. It comes like a node A->B->C, So A-->C. This type of relations is called transitive.



In the above hierarchy, there exist both transitive and sibilings relation. The node விண்வெளி has three sibilings namely (வான்வெளி, ஆகாயவெளி, வானம்), since all these three nodes are direct child. There is not sibilings exist for the nodes (வான்வெளி, ஆகாயவெளி, வானம்). Transitive relations comes between the nodes விண்வெளி ->வானம் ->ஆகாயம்,வின்,வான். Node entity விண்வெளி has a direct child as வானம் and வானம் has a three direct child ஆகாயம்,வின்,வான். So using transitive rule, it can be structured as விண்வெளி ->ஆகாயம்வின்வான்.

**Building Onltology Framework**

      Generating ontology tree from a given huge data corpus based on the user query is a backbone activities for this intelligent information retrieval and visualization process. This type of tool provides a powerful way for data representation and knowledge mining process. User can input search word to the system, then the system will analyze the dataset based on the word and identify the location of the given input word inside the dataset. If the input is located in the header division, then the system fetch the related header node, else if the input is located in the non-header division then the related data from the non-header division is fetched to generate the simple hierarchy and ontology tree. The following are the flow steps and proposed algorithm to build ontology and to visualize it effectively.

The following are the data representation label used in the ontology tree generation process.

| Notations | Tamil Equivalent | Symbols |
|---|---|---|
| Synonym | இணைமொழியம் | <> |
| Hyponym | வகைமொழியம் | { } |
| Meronym | பகுதிமொழியம் | / \| |
| Gender | பால்மொழியம் | $# |
| Antonym | எதிர்மொழியம் | * |
| Derived | பெயராக்கமொழியம் | ~ ` |
| Adv | வினையுரிச்சொல் | ^ ! |

- *Synonym<>*

  It represents the different words which gives same meaning. It shows one to many relationship in the hierarchy. Eg:Sysnonym<>சூரியன், <ஞாயிறு, ஆதவன், பகலவன், திவாகரன், ஆதித்தன், கதிரவன், கதிரோன், கிரணன், செங்கதிரோன், செங்கதிர், வெங்கதிரோன், வெங்கதிர், வெஞ்சுடர்ஆயிரங்கதிரோன், ஆயிரங்கிரணன், உலகநேந்திரன், உலகப்பாந்தவன்>; { எழுஞாயிறு, ஏறுஞாயிறு; சாய்ஞாயிறு, இறங்குஞாயிறு; இளஞ்சூரியன், <பாலசூரியன்>; }.

- *Hyponym { }*

  சப்தமண்டலம்2: { வாயுமண்டலம்; ருணமண்டலம்; சந்திரமண்டலம்; சூரியமண்டலம்; ட்சத்திரமண்டலம்; அக்கினிமண்டலம்; திரிசங்குமண்டலம்; }

- *Meronym / |*

  காற்றுமண்டலம்:/வெளிமண்டலம்;வெப்பமண்டலம், உஷ்ணமண்டலம்;படலம்:அயனிப்படலம்;ஓசோன்படலம்; |

- *Gender $ #*

  இளைஞர்கள்: {இளைஞர் $ இளைஞன் #, வாலிபர் $ வாலிபன் #, வாலிபப்பையன், இளவல், இளந்தாரி, இளைஞோர் $ இளைஞோன் #, காளை, விடலை, இளவட்டம், வயசுப்பையன்,வயசுப்பிள்ளை, வயசுப்பிள்ளையாண்டன், இளம்வயதினன்-இளம்வயதினர், இளம்பருவத்தார் $ இளம்பருவத்தான்#, இளம்பிராயத்தார் $ இளம்பியாயத்தான்#, பதின்பருவத்தினர் $ பதின்பருவத்தினன் #, பதின்வயதினர் $ பதின்வயதினன் #, வயதுவந்தோர், <வயதுவந்தவர்> $ வயதுவந்தவன் #, இளந்தை, இளந்தாரி, வாலியன், பாலியன்; இளரத்தம்>}

- *Derived ~ `*

  நிகழ்வுகள்குறித்தவை: { நிகழ்; ~ நிகழ்தல்; நிகழ்ச்சி; நிகழ்வு; ` <சம்பவி; ~ சம்பவித்தல்; சம்பவம்; ` விளை; > ~ விளைதல்; விளைவு; ` >நேர்; ~ நேர்தல்; ` பலி; ~ பலித்தல்; ` ஈடேறு; ~ ஈடேறுதல்; ஈடேற்றம்; ` நிலவு; ~ நிலவுதல்; ` }

- *Antonym% *  
  It gives you opposite meaning for a word. It shows one to one relationship in the hierarchy.  
  *Eg:ஆண்மரம் - >பெண்மரம்.*

- *Adv/Adj^ !*  
  நீண்டகாலம்தொடர்பானவை: { நீண்டகாலம்; ^ நீண்டகாலமாக !: நீடூழி, <நெடுங்காலம் ^ நெடுங்காலமாக; !  
  ரெம்பகாலம்; ^ ரெம்பகாலமாக; ! நீண்டநேரம், <நெடுநேரம், ரெம்பநேரம், நீண்டநேரம், நிறையநேரம்; >நீண்டநாள்,  
  <நெடுநாள்; ரெம்பநாள், நிறையநாள்>; பல்லாண்டு;நூறாண்டு, <நூற்றாண்டு> }

Onto-thesaurus is a knowledge representation and these knowledge are visualized in the form of clusters instead of hierarchical tree. Each parent entity and its all available sibling entities re-grouped on to a cluster and all clusters generated during search patterns are linked with the label entities.

**CONCLUSION**

Presently we developed a real time system for Ontology based intelligent system for information retrieval. This system works with huge datasets which contains 50,000 word count for more general domain. So it a word based information retrieval which gives hierarchy and ontology tree for the user. We don't have any mathematical related computational works right now in this system and it can be accommodated in future research works. Based on requirement, the same system can also be used for domain based dataset and any other use cases. Since the implementation has been done generic manner which support any languages. Onto-thesaurus can be a very efficient tool for visualizing the dataset from large data corpus. There data corpus are represented using XML format.In future, such implementation for ontology based intelligent system can be integrated with probability graphical model (PGM) to improve the existing features with mathematical modeling to shows the relationship between the nodes in the hierarchy based on probability distributions. Apart from that, distance between the node can also be calculated from the generated ontology tree. These information's can be pretty much useful for further ontology research and improve more knowledge mining during representation.
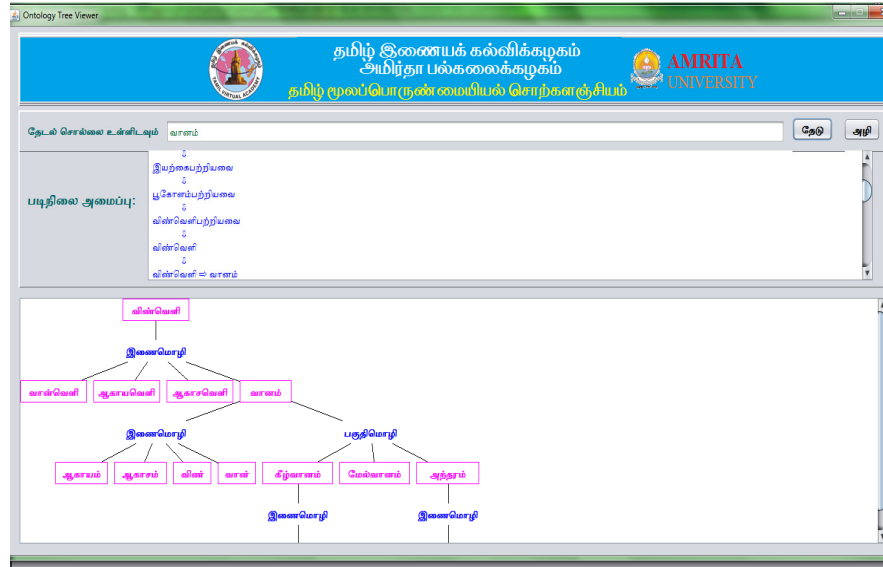


Fig 2.1மூலப்பொருண்மையியல்
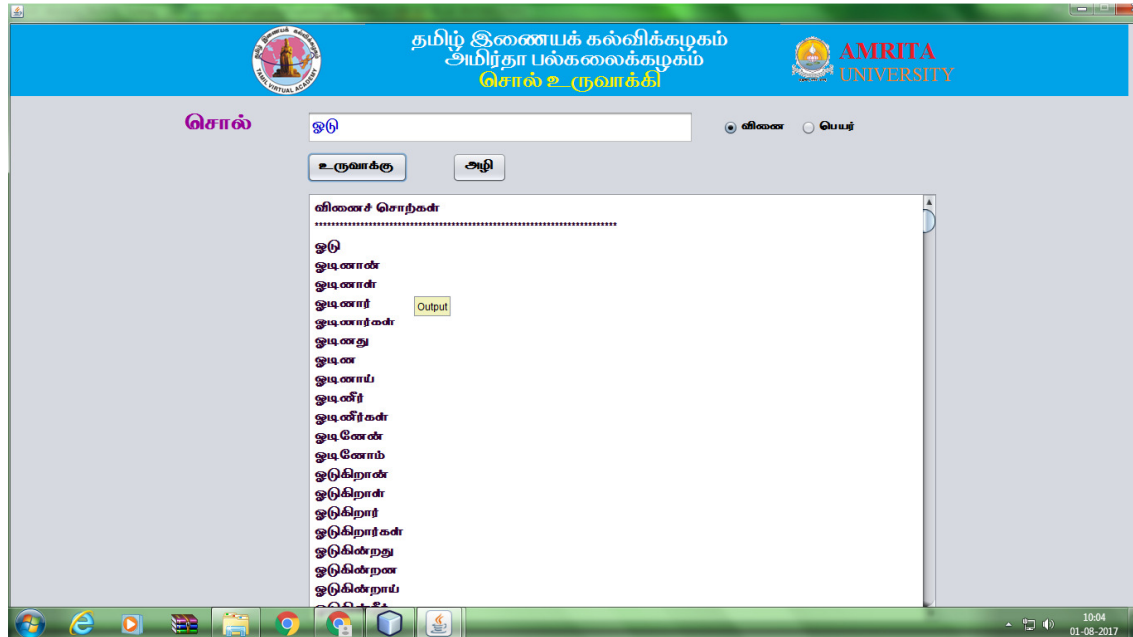
## III. Tool for learning Morphology generator



Fig 3.1சொல் உருவாக்கி

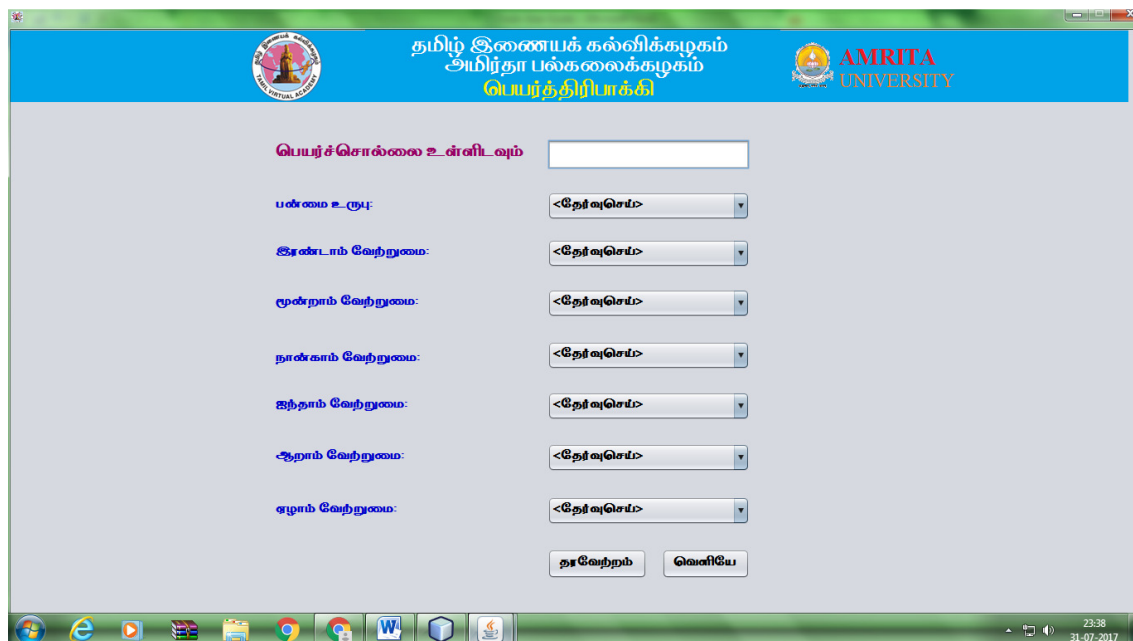## IV. Computing Tools for Teaching and Learning
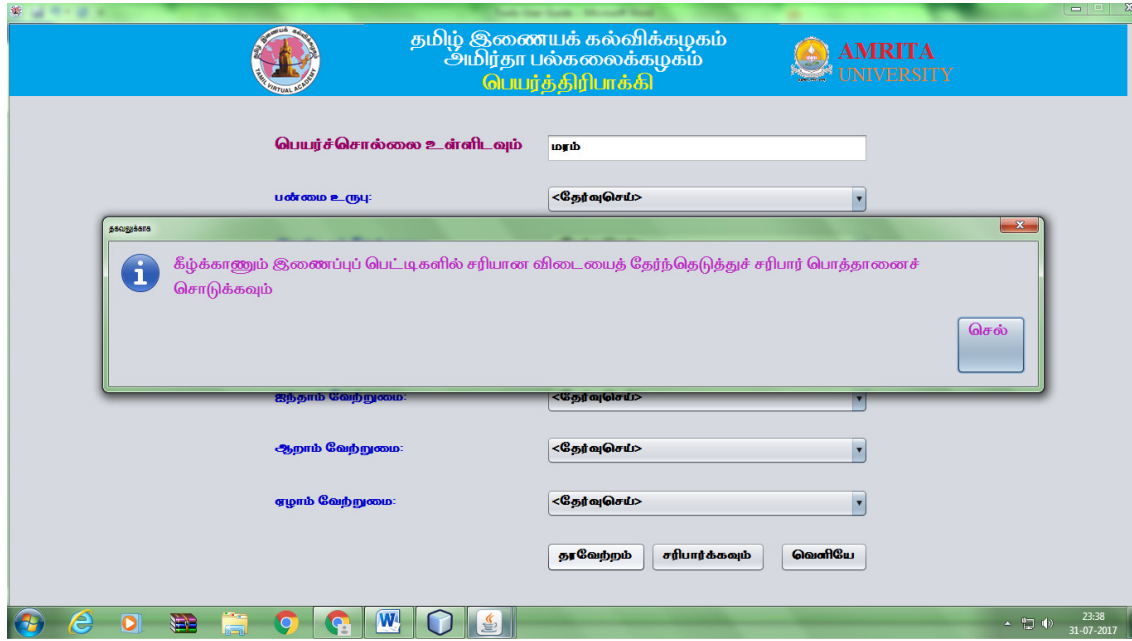


Fig 4.1பெயர்த்திரிபாக்கி

Fig 4.2பெயர்த்திரிபாக்கி



Fig 4.3பெயர்த்திரிபாக்கி
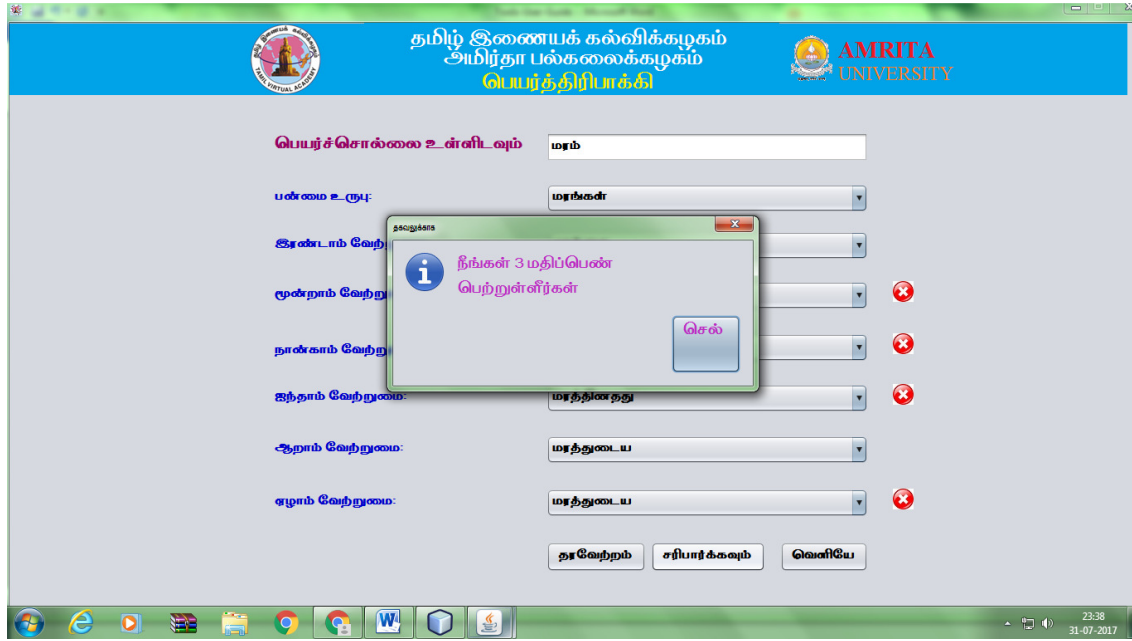
Fig 4.4வினைத்திரிபாக்கி



Fig 4.5வினைத்திரிபாக்கி

Fig 4.6வினைத்திரிபாக்கி



Fig 4.7வினைத் திணை -எண் பால்

Fig 4.8வினைத் திண  -எண் பால்



Fig 4.9வினைத் திண  -எண் பால்

Fig 4.10 கால விகுதி



Fig 4.11 கால விகுதி

**A Smart English-Tamil Electronic dictionary**

A Smart English-Tamil Electronic dictionary has been developed as the third component of the project. It has nearly one lakh entries for general vocabulary. The technical vocabulariesof various branches of knowledge have been integrated into the present super dictionary. The different domains of technical terms included in the super dictionary are the following: Administration, Agriculture, Arts and Humanities,Building and Construction, Economics, Engineering and Technology, Home Sciences, Law, Linguistics, Metallurgy, general Science, Social Work, Veterinary science, Chemical Engineering, Civil and Mechanical Engineering, Electrical and Electronics Engineering, Engineering and Technology(General), Information Technology, Mechanical Engineering, Textile Technology, Child Development, Cookery, First Aid Sick Nursing, Home Management, Home Science (General), Laundry, Needle Work , Botany, Chemistry, Geology, Mathematics, Physics, Statistics, and Zoology. A very useful GUI has been developed through which the user can get the information he expects from a English-Tamil bilingual dictionary. We hope to improve this tool with many user friendly and user useful components and develop it into a real smart English-Tamil dictionary.
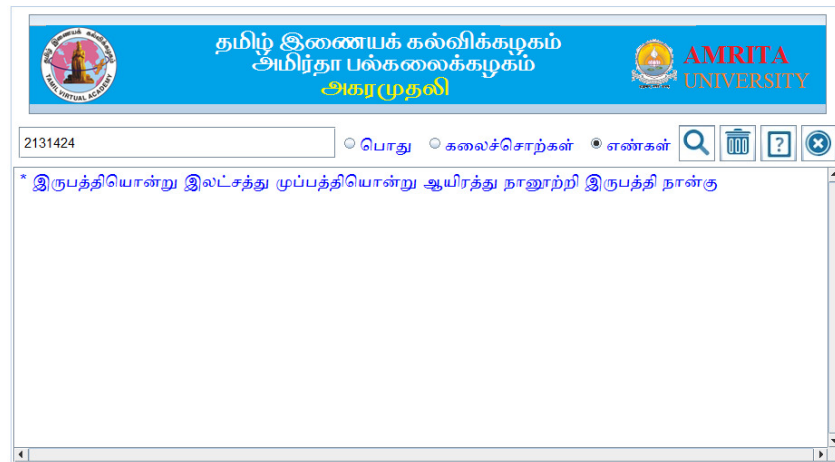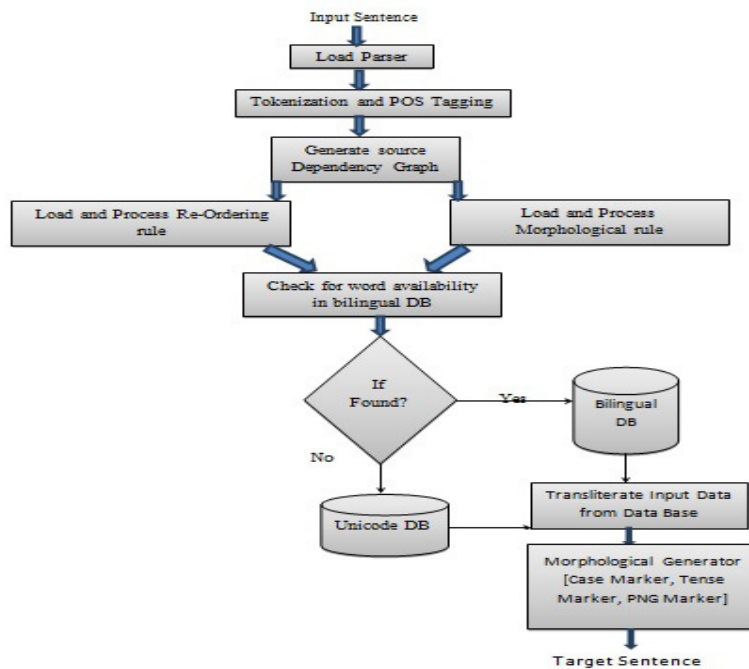


Fig 5.1 அகரமுதலி



Fig 5.2 அகரமுதலி

**English to Tamil Rule Based Machine Translation System (RBMTS)**

```
                    Input Sentence
                         │
                   ┌─────▼─────┐
                   │ Load Parser │
                   └─────┬─────┘
              ┌──────────▼──────────────┐
              │ Tokenization and POS Tagging │
              └──────────┬──────────────┘
                   ┌─────▼─────┐
                   │  Generate source  │
                   │ Dependency Graph  │
                   └──────────────────┘
  ┌───────────────────┐              ┌───────────────────┐
  │ Load and Process Re-Ordering │   │ Load and Process │
  │        rule        │          │ Morphological rule │
  └───────────────────┘              └───────────────────┘
              ┌───────────────────────┐
              │ Check for word availability │
              │      in bilingual DB      │
              └───────────┬───────────┘
                      ◇ If Found? ◇ ──Yes──▶ Bilingual DB
                          │
                          No
                          │
                     Unicode DB ──────▶ Transliterate Input Data
                                              from Data Base
                                                  │
                                      Morphological Generator
                                      [Case Marker, Tense
                                       Marker, PNG Marker]
                                                  │
                                           Target Sentence
```

**Flow Description:-**

❖ The system accept input sentence from the user and passes to the Stanford parser to tokenize the sentence in to multiple word tokens and identify the POST(Parts of Speech Tagging) for each words from input sentence.

❖ After POST, Stanford parser generates dependency graph for source input sentence and then load both Re-Ordering rule and Morphological rule.

❖ After Morphological rule implementation, system check for each word token availability in Bilingual Database. If the word found, then transliteration process executed.

❖ If the word token not available in Bilingual Data Base, then system execute transliterate process from Unicode Data Base.

❖ After Transliteration, Morphological Generator generates and provides morphological information for input sentence. This process contains Case Marker, Tense Marker and PNG (Person, Number, and Gender) marker.

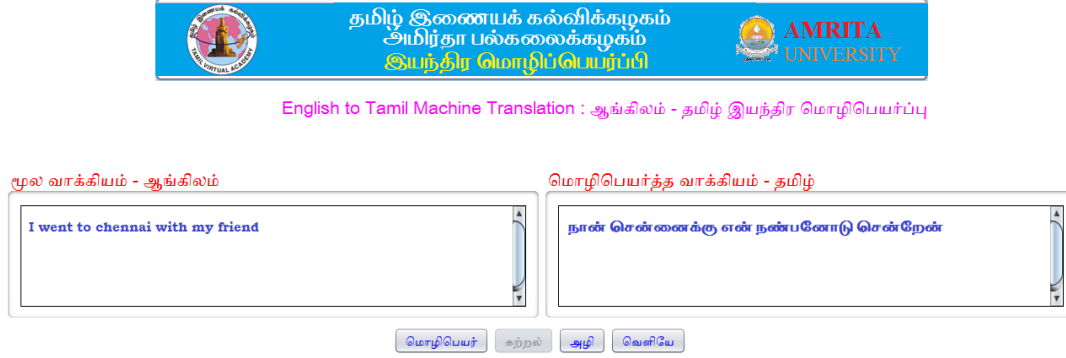❖ Finally Morphological Generator generates Target Sentence.

**Screenshot:-**

English to Tamil Machine Translation : ஆங்கிலம் - தமிழ் இயந்திர மொழிபெயர்ப்பு

மூல வாக்கியம் - ஆங்கிலம்

மொழிபெயர்த்த வாக்கியம் - தமிழ்

I went to chennai with my friend

நான் சென்னைக்கு என் நண்பனோடு சென்றேன்

மொழிபெயர்   கற்றல்   அழி   வெளியே

Fig 5.3 இயந்திர மொழிப்பெயர்ப்பி

### Dictionary Statistics:-

- ✓ No. of Database          : 1
- ✓ No. of Table             : 1
- ✓ No.of records            : 9500
- ✓ Total Re-Ordering Rules  : 211
- ✓ Morphological Rules      : 29
- ✓ Parser Used              : stanford-parser3.5.2

### Test Sentences:-

| S.No | Source Sentence | Target Sentence |
|---|---|---|
| 1 | She **cannot** reveal it | அவள் அதை வெளிப்படுத்தமுடியாது |
| 2 | I **does not** simplify the equation | நான் சமன்பாட்டை தெளிவாக்கவில்லை |
| 3 | Ram **should not** wait | ராம் காத்திருக்கக்கூடாது |
| 4 | I simulate my work | நான் என் வேலையை வடிவமைக்கிறேன் |
| 5 | I shall not wash my hands | நான் என் கைகளை கழுவமாட்டேன் |
| 6 | Seeta met ravana **in** the airport | சீதா ராவணனை விமான நிலையத்தில் சந்தித்தாள் |
| 7 | My father came **from** chennai | என் தந்தை சென்னையிலிருந்து வந்தார் |
| 8 | This **is** for my friend | இது என் நண்பனுக்காக |
| 9 | I went **to** market with my sister | நான் சந்தைக்கு என் தங்கையோடு போனேன் |
| 10 | I kept my pen **above** the table | நான் என் பேனாவை மேசைக்கு மேலே வைத்தேன் |
| 11 | He is **below** average in the class | அவர் வகுப்பில் சராசரிக்கு கீழே |
| 12 | I have **12345678** rupee | நான் ஒரு கோடியே இருபத்திமூன்று இலட்சத்து நாற்பத்தியைந்து ஆயிரத்து அறுநூற்றி எழுபத்தி எட்டு ரூபாயை வைத்திருக்கின்றேன் |

*CONCLUSION AND FUTURE ENHANCEMENT*

Each above stages are mandatory for translation engine, where even a single stages cannot be skipped to develop a simple system. Since the system approach us rule based, it performs well for all kind of simple sentences. But it may not able to translate exactly for long and complex sentences. So to overcome this difficulty, it should be extended and enhanced to the next level of machine translation methodology.