

MATEMATICKO-FYZIKÁLNÍ FAKULTA
PRAHA

**TAMIL DEPENDENCY TREEBANK (TAMILTB) - 0.1
ANNOTATION MANUAL**

LOGANATHAN RAMASAMY, ZDENĚK ŽABOKRTSKÝ

ÚFAL/CKL Technical Report
TR-2011-42



UNIVERSITAS CAROLINA PRAGENSIS

Copies of ÚFAL/CKL Technical Reports can be ordered from:

Institute of Formal and Applied Linguistics (ÚFAL MFF UK)

Faculty of Mathematics and Physics, Charles University

Malostranské nám. 25, CZ-11800 Prague 1

Czech Republic

or can be obtained via the Web: <http://ufal.mff.cuni.cz>

Tamil Dependency Treebank (TamilTB) - 0.1
Annotation Manual

Loganathan Ramasamy and Zdeněk Žabokrtský
ÚFAL/CKL Technical Report
TR-2011-42

December 11, 2011

Contents

1	Introduction	1
1.1	Background and Objectives	1
1.2	Data	1
1.3	Text Preprocessing	1
1.3.1	Transliteration	2
1.3.2	Sentence Segmentation	2
1.3.3	Tokenization	3
1.4	Layers of Annotation	4
1.4.1	Morphological Layer (m-layer)	4
1.4.2	Analytical Layer (a-layer)	5
1.5	Obtaining Data	5
1.5.1	Contents of TamilTB.v0.1.tar.gz	6
1.5.2	Prerequisites	6
1.5.3	TectoMT Development Platform	6
1.5.4	Contact Authors	7
1.5.5	License	7
2	Morphological Annotation	8
2.1	Positional Tagging	8
2.1.1	Position 2 - Sub POS	9
2.1.2	Position 3 - Case	15
2.1.3	Position 4 - Tense	15
2.1.4	Position 5 - Person	15
2.1.5	Position 6 - Number	16
2.1.6	Position 7 - Gender	16
2.1.7	Position 8 - Voice	16
2.1.8	Position 8 - Negation	16
2.2	Annotation of Pronouns	16
2.2.1	Singular Referential Definite (Personal) Pronouns	17
2.2.2	Non Referential (Interrogative) Pronouns	17
2.2.3	General Referential Pronouns	17
2.2.4	Specific Indefinite Referential Pronouns	18
2.2.5	Non-specific Indefinite Referential Pronouns	18
3	Syntactic Annotation	19
3.1	Identifying the Structure	19
3.2	Dependency Relations	19
3.3	Detailed Description of <i>afuns</i>	21
3.3.1	Afun: <i>AAdjn</i>	21
3.3.2	Afun: <i>AComp</i>	21
3.3.3	Afun: <i>Apos</i>	22
3.3.4	Afun: <i>Atr</i>	22
3.3.5	Afun: <i>AdjAtr</i>	23
3.3.6	Afun: <i>AuxA</i>	24
3.3.7	Afun: <i>AuxC</i>	25
3.3.8	Afun: <i>AuxG</i>	25

3.3.9	Afun: <i>AuxK</i>	25
3.3.10	Afun: <i>AuxP</i>	25
3.3.11	Afun: <i>AuxS</i>	26
3.3.12	Afun: <i>AuxV</i>	27
3.3.13	Afun: <i>AuxX</i>	27
3.3.14	Afun: <i>AuxZ</i>	27
3.3.15	Afun: <i>CC</i>	28
3.3.16	Afun: <i>Comp</i>	30
3.3.17	Afun: <i>Coord</i>	30
3.3.18	Afun: <i>Obj</i>	31
3.3.19	Afun: <i>Pnom</i>	32
3.3.20	Afun: <i>Pred</i>	32
3.3.21	Afun: <i>Sb</i>	33

List of Tables

1.1	The data for annotation	2
1.2	List of suffixes and words for tokenization	3
2.1	SubPOS	10
2.2	Tamil case	15
2.3	Tense	15
2.4	Person	15
2.5	Number	16
2.6	Gender	16
2.7	Voice	16
2.8	Negation	16
2.9	Personal pronouns	17
2.10	Interrogative pronouns	17
2.11	General referential pronouns	18
2.12	Specific indefinite referential pronouns	18
2.13	Non-specific indefinite referential pronouns	18
3.1	afuns	20

List of Figures

1.1	Transliteration scheme	2
1.2	Tokenization example	3
1.3	Tamil sentence example	4
1.4	An example for morphological layer annotation	5
1.5	An example for analytical layer annotation	5
2.1	Positional tag	8
2.2	Positional tagset	9
3.1	<i>AAdjn</i> : Adverbial adjunct	21
3.2	<i>AAdjn</i> : Adverbial adjunct	22
3.3	<i>AComp</i> : Adverbial complement	22
3.4	<i>Apos</i> : Apposition	23
3.5	<i>Atr</i> : Attribute	23
3.6	<i>AdjAtr</i> : Adjectival attribute	24
3.7	<i>AdjAtr</i> : Adjectival attribute	24
3.8	<i>AuxA</i> : Determiners	25
3.9	<i>AuxC</i> : Subordinating conjunctions	26
3.10	<i>AuxC</i> : Subordinating conjunctions	26
3.11	<i>AuxG</i> : Symbols	27
3.12	<i>AuxK</i> : Sentence termination symbol	27
3.13	<i>AuxP</i> : Postpositions	28
3.14	<i>AuxP</i> : Postpositions	28
3.15	<i>AuxV</i> : Auxiliary verb	29
3.16	<i>AuxX</i> : Commas	29
3.17	<i>AuxZ</i> : Emphasis	30
3.18	<i>AuxZ</i> : Emphasis	30
3.19	<i>CC</i> : Marker for a multi word sequence	31
3.20	<i>Comp</i> : Complement (other than verbs)	31
3.21	<i>Coord</i> : Coordination head (English style)	32
3.22	<i>Coord</i> : Coordination head (Comma)	32
3.23	<i>Coord</i> : Coordination head (Morphological marker <i>-um</i>)	33
3.24	<i>Obj</i> : Object	33
3.25	<i>Pnom</i> : Pnom	34
3.26	<i>Pnom</i> : Pred	34
3.27	<i>Sb</i> : Subject	34
3.28	<i>Sb</i> : Subject	35

Abstract

Tamil Dependency Treebank (TamilTB) - 0.1 [Ramasamy and Žabokrtský, 2011] is an attempt to develop a syntactically annotated corpora for Tamil. TamilTB contains 600 sentences enriched with manual annotation of morphology and dependency syntax in the style of Prague Dependency Treebank. TamilTB has been created at the Institute of Formal and Applied Linguistics, Charles University in Prague. This report serves the purpose of how the annotation has been done at morphological level and syntactic level. Annotation scheme has been elaborately discussed with examples.

Chapter 1

Introduction

This chapter will briefly explain the Tamil Dependency Treebank (TamilTB) in general and various tasks involved in developing the TamilTB.

1.1 Background and Objectives

Treebank is an important resource in many Natural Language Processing (NLP) applications including parsers, language understanding, Machine Translation (MT) and so on. Treebanks are often manually developed and are available for only handful of languages such as English, German and Czech. Most of the world's other languages (irrespective of contemporariness and the number of speakers) do not have treebanks due to various reasons: (i) high development cost (ii) long development time (iii) lack of expertise to name a few. In this project, our main aim is develop a dependency treebank for Tamil language similar to Prague Dependency Treebank (PDT) annotation style.

We briefly introduce the work carried out on the subject prior to this work. There is an active research on dependency parsing [Bharati et al., 2009], [Nivre, 2009] and developing annotated treebanks for other Indian languages such as Hindi and Telugu. One such effort is, developing a large scale dependency treebank [Begum et al., 2008] (aimed at 1 million words) for Telugu, as of now the development for which stands [Vempaty et al., 2010] at around 1500 annotated sentences. For Tamil, previous works which utilised Tamil dependency treebanks are: [Dhanalakshmi et al., 2010] which developed dependency treebank (around 25000 words) as part of the grammar teaching tools, [Selvam et al., 2009] which developed small dependency corpora (5000 words) as part of the parser development. Other works such as [Janarthanam et al., 2007] focused on parsing the Tamil sentences. The current work differs from previous works with respect to the following objectives.

The main objectives of this project include,

- Annotate data at word level and syntactic level
- In each level of annotation, try for maximum level of linguistic representation
- Building large annotated corpora using automatic annotation process

1.2 Data

The data used for the TamilTB annotation comes from the news domain. We decided to use the news data for two reasons: (i) huge amount of data is available in digital format and can be easily downloadable and (ii) the news data can be considered as representative of written Tamil. At present, the data for the annotation comes from www.dinamani.com, and we downloaded pages randomly covering various news topics.

1.3 Text Preprocessing

Before the actual annotation takes place, the raw text data is preprocessed in three steps in sequential order,

No	Description	Value
1	Source	www.dinamani.com
2	Source transliterated	Yes
3	Number of words	9581
4	Number of sentences	600
5	Morphological annotation (sen)	600
6	Syntactic annotation (sen)	600
7	Tectogrammatical annotation	–

Table 1.1: The data for annotation

- Transliteration
- Sentence segmentation
- Tokenization

Each preprocessing step is explained in the following subsections.

1.3.1 Transliteration

The UTF-8 encoded Tamil raw text was transliterated to Latin for ease of processing during all levels of annotation. The raw UTF-8 encoded text can be obtained by applying reverse transliteration to the Latin-transcribed text. The transliteration scheme for Tamil script is given in the Figure 1.1. The Figure shows the transliteration for vowels, consonants, Sanskrit characters and an example transliterated sequence of consonant ('k')-vowel combination. Tamil has separate character representation for each consonant-vowel combination.¹

Tamil Vowels Transliteration	அ ஆ இ ஈ உ ஊ எ ஏ ஐ ஒ ஓ ஔ ஃ a A i I u U e E ai o O au q
Tamil Consonants Transliteration	க ங ச ஞ ட ண த் ந ப் ம ய் ர் ல் வ் ழ் ள் ற் ன் k ng c nj t N T w p m y r l v z L R n
Sanskrit Characters Transliteration	ஷ் ஹ் ஜ் ஸ் sh h j sri S
'k' + vowel combination Transliteration	க கா கி கீ கு கூ கெ கே கை கொ கோ கௌ ka kA ki kI ku kU ke kE kai ko kO kau

Figure 1.1: Transliteration scheme

Tamil is a phonetic language and the transliteration scheme is designed to match the Tamil sounds as much as possible. In this documentation, wherever possible, both Tamil script and transliterated form are used in examples. In few places, only transliterated format is used due to difficulties (mainly rendering problems) in embedding Tamil scripts.

1.3.2 Sentence Segmentation

Before the annotation takes place, the raw corpus downloaded from the source is sentence segmented automatically. This step can be performed before or after the transliteration. Like English, Tamil can also be ambiguous at various places that may look like sentence boundaries. But in reality they may not be sentence boundaries. Those ambiguous sentence boundaries (such as dots at decimal numbers, initials in names and dates) are detected through heuristics and the sentences are segmented only at appropriate places. As an example, in Tamil, the proper names are written using ('*surname name*') format. But in usage the format is shortened to ('**initial.** *name*') where a dot is placed between *initial* and *name*. *Initial* is the first letter of a *surname*. Thus, the fullname '*Palaniappan Chidambaram*' in English will

¹For full list of transliteration map, please refer this URL: https://ufal.mff.cuni.cz/~ramasamy/tamiltb/0.1/utf8_to_latin_map.txt (change the encoding of the browser to UTF-8 to view the contents properly)

be written as ‘*P. Chidambaram*’ & ‘*pa. ciTambaram*’ in English and Tamil respectively. The heuristics for this problem is straightforward. We listed down all possible Tamil characters (letters) and whenever a sentence termination symbol (*dot*) occurs after an initial will not be treated as a sentence boundary.

1.3.3 Tokenization

Tokenization is one of the important module that helps the annotation task. This module splits the sentence into words. Tamil uses spaces to mark word boundaries. But yet, a lot of Tamil wordforms are agglutinative in nature, meaning they glue together atleast two words (in majority of cases). Those cases can be identified as determiners+nouns, nouns+postpositions, verbs+particles, nouns+particles and etc. Except the first pattern (determiners+nouns), in all other cases, the second part of the wordforms are restricted and can be listed. So it is possible to split certain Tamil agglutinative wordforms into separate tokens. Certain particles (also called clitics) such as *um/also*, *O/or* are not treated as separate tokens in Tamil. But for the purpose of annotation we treat them as separate tokens. The same module will be used for tokenization when parsing the raw Tamil text.

Before splitting	Tamil	புதிய சட்டத்தின்படி , பாதுகாக்கப்பட்ட நினைவுச் சின்னத்திலிருந்து 1000 அடி வரை எந்த கட்டுமானமும் கட்ட அனுமதி இல்லை .
	Trans.	puTiya cattaTTin pati , pATukAkkapp atta winaivuc cinnaTT iliruwTu 1000 ati varai ewTa kattumAnam um katta anumaTi illai .
After splitting	Tamil	புதிய சட்டத்தின் படி , பாதுகாக்கப் பட்ட நினைவுச் சின்னத்தி லிருந்து 1000 அடி வரை எந்த கட்டுமான மும் கட்ட அனுமதி இல்லை .
	Trans.	puTiya cattaTTin pati , pATukAkkap atta winaivuc cinnaTT iliruwTu 1000 ati varai ewTa kattumAnam um katta anumaTi illai .

Figure 1.2: Tokenization example

In the Figure 1.2, *pati* (‘manner’) and *iliruwTu* (‘from’) are postpositions, *atta* is an auxiliary verb and *um* is a clitic. This kind of agglutination is very prevalent in Tamil, and it would be useful to tagging process if we are able to reduce the vocabulary size by splitting the known combinations as separate tokens.

Clitics	um, E, EyE, AvaTu
Postpositions	<i>kUta, utan, pati, kuRiTtu, iliruwTu, anRu, uL, ARu, Tavira, pOTu, pOla, pinnar, pin, arukE, aRRa, inRi, illATa, mITu, kIz, mEl, munpE, otti, paRRi, paRRiya, pOnRa, mUlam, vaziyAka</i> etc.
Auxiliary Verbs	patta, pattu, uLLa, pata, mAttATu, patuvArkaL, uLLAr, uLLanar, illai, iruwTAr, iruwTaTu, pattaTu, pattana, mutiyum, kUtATu, vENTum, kUtum, iruppin, uLLana, mutiyATu, patATu, koNtu, ceyTu etc.
Particles	<i>Aka, Ana</i> and their spelling variants <i>Akac, AkaT, Akap, Akak</i>
Demonstrative pronouns	<i>ap, ac, ic, aw, iw</i> etc. as prefixes

Table 1.2: List of suffixes and words for tokenization

Some of the most commonly occurring (from the corpus) words and suffixes which participate in agglutination is given in the Table 1.2 above. Except demonstrative pronouns, all other words and suffixes are added after the stem. Among the categories in the Table 1.2 above. Clitics and Particles are the most participated in the agglutination. The tokenizer will make use of this list and try to separate these words from the original wordform. Even after the tokenization it would be possible to reconstruct the original sentence by making use of the attribute called ‘*no_space_after*’. The ‘*no_space_after*’

will be set to 1 if the following token is not separated from the current token. Whenever the splitting takes place this attribute will be set to 1 for the first token. For example, The ‘no_space_after’ attribute for *pATukAkkap* will be 1. Whereas the ‘no_space_after’ attribute for *um* will be 0. The splitting for the corpora has been done semi-automatically using some of the most commonly occurring combination from the above list and edited manually during the annotation process. At present, the tokenizer includes only few commonly occurring combinations from the Table 1.2 such as Clitics, Particles and very few postpositions.

We evaluated how much such combinations have been splitted from the original corpora. We found that 953 splits took place out of 9581 words. We simply did this by counting how many ‘no_space_after’ attributes have been set to 1. We can say that almost 10% of the additional corpus size is due to splitting some wordforms into separate tokens.

1.4 Layers of Annotation

The annotation scheme used for TamilTB.v0.1 is similar to that of Prague Dependency Treebank 2.0 (PDT 2.0) [UFAL, 2006]. PDT 2.0 uses the notion ‘layers’ to distinguish annotation at various levels (linguistic) such as word level and syntactic level. Precisely, PDT 2.0 is annotated at 3 levels or layers: (i) morphological layer (m-layer), (ii) analytical layer (a-layer) and (iii) tectogrammatical layer (t-layer). At present, TamilTB is annotated at only two layers: m-layer and a-layer.

Tamil	பண்பாட்டு அடையாளங்களைப் பாதுகாக்க தொல்பொருள் ஆய்வுத் துறை உருவாக்கப் பட்டு . தனிச் சட்டங்கள் இயற்றப் பட்டு உள்ளன .
Trans	paNpAttu ataiyALangkaLaip pATukAkka TolporuL AyvuT TuRai uruvAkkap pattu , Tanic cattarengkaL iyaRRap pattu uLLana .
English	Having created Archeological Department, separate laws have been enacted to protect cultural symbols .

Figure 1.3: Tamil sentence example

In the the example shown in the Figure 1.3, the actual setence is given in Tamil script (indicated as Tamil:) in the 1st row, the transliterated (indicated as Tr:) version in the 2nd row, and the actual English translation in the 3rd row. The same format is used to illustrate sentence examples elsewhere in the document. There are 15 words in the Tamil sentence (including punctuations), each word will be treated as a node in each annotation layer. Please note that the term *node* will be used interchangeably with other terms *wordform* or *word* to represent a vertex in a-layer or m-layer.

Each node will have general attributes and layer specific attributes. For ex: a node in morphological layer will have attributes such as, ‘lemma’, ‘form’, ‘tag’ and ‘no_space_after’ corresponding to lemma, wordfom, POS tag of a particular wordform and whether the following wordform is part of the current wordform . A node in analytical layer will have attributes such as dependency label (‘afun’) of the current node, whether the current node is an element in the coordination conjunction (‘is_member’) etc. These attributes will be set automatically during parsing or editing attributes manually using TrEd. Also, the lower layer (m-layer) attributes are visible to upper layers (a-layer or t-layer).

Only transliterated version of the text will be used in all layers of annotation for the ease of processing. Examples in Tamil script are shown only for display purposes.

The following subsections briefly describe the annotation layers of TamilTB with an example.

1.4.1 Morphological Layer (m-layer)

The purpose of m-layer is to assign Parts of Speech (POS) tag or more refined morphological tag to each word in the sentence. This is accomplished by setting the *tag* attribute of the node (corresponds to word) to the POS or morphological tag. The *lemma* attribute will store the conceptual root or the word listed in dictionary as the lemma of the wordform. The Figure 1.4 illustrates m-layer annotation.

The Figure shows, there are three text values that are displayed at each node. The text at the top of the node is the *form* or the exact word which appeared in the text. The text at the middle (for ex: *paNpAttu*) of the node is the *lemma*, and the text at the bottom (for ex: NO - - 3SN - -) of the node is the morphological tag of the wordform. The length of each morphological tag is 9 characters and each character position will correspond to some feature of a wordform. The first 2 positions in the morpholigical

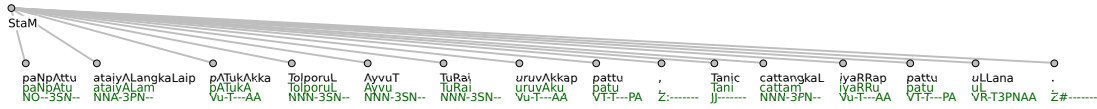


Figure 1.4: An example for morphological layer annotation

tag corresponds to main POS and refined POS. Both together will represent fine details of a wordform. Thus it is possible to train the POS tagger for a fine grained tagset or coarse grained tagset. This kind of tagging is known as positional tagging. Positional tagging is suitable for morphologically rich languages and has been successfully applied to languages such as Czech. The Chapter 2 gives a detailed description about positional tagging and the tagset used to perform annotation for TamilTB.

1.4.2 Analytical Layer (a-layer)

Analytical layer (a-layer) is used to annotate the sentence at syntactic level. There are two phases in a-layer annotation: (i) capture the dependency structure of a sentence in the form of tree and (ii) identify the relationship between words or nodes in the tree. From m-layer, we know that each wordform corresponds to a node in the tree but they are without their parents assigned. The dependency structure is captured by hanging the dependent nodes (words) under their governing nodes (words). Visually, dependent nodes will hang as children of their governing nodes. There will be one extra node called *technical root* to which the predicate node and the terminal node (end of the sentence) will be attached. The sole purpose of the *technical root* is to have some tree level attributes such as tree identifier. The Figure 1.5 illustrates the a-layer annotation of a sentence shown in Figure 1.3.

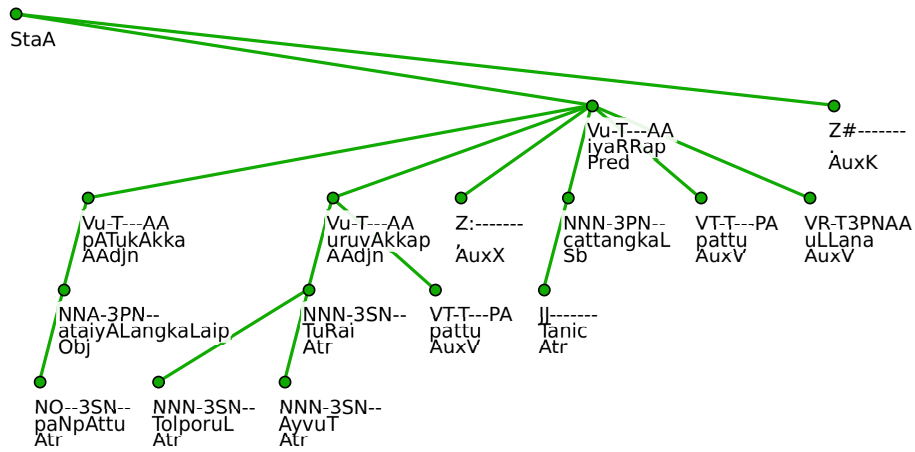


Figure 1.5: An example for analytical layer annotation

Edges between nodes indicate the relationship with which they are connected. In linguistic terms, it is called syntactic relation between governor and dependent. The relationship between two nodes are stored in the attribute called *afun*. In PDT style annotation, for technical reasons, the edges do not have *afun* attribute, instead the dependent nodes will store the *afun* attribute. For example, the *afun* value of the word *cattangkaL* ('laws') is *Sb* meaning Subject, is connected to the verb *iyaRRap* ('to enact').

The more detailed treatment of various syntactic relationships and a-layer annotation scheme is given in Chapter 3.

1.5 Obtaining Data

The annotated data is available in three formats,

1. TMT format - XML-based format used in the TectoMT system
2. CoNLL format - tabular separated format in the CoNLL shared task style
3. TnT style POS tagged format - tabular separated columns with word forms, POS tags, and lemmas.

The single package containing the data and the documentation can be downloaded from the following link,

<http://ufal.mff.cuni.cz/~ramasamy/tamilTB/0.1/download.html>

The package can be uncompressed using the following command,

```
[shell]$ tar -zxvf TamilTB.v0.1.tar.gz
```

1.5.1 Contents of TamilTB.v0.1.tar.gz

- TamilTB.v0.1/ (top level directory)
 - TamilTB.v0.1.tmt (in TMT format)
 - TamilTB.v0.1.conll (in CoNLL format)
 - TamilTB.v0.1.tt (only POS tagged corpora in TnT style)
- doc
 - index.html (main page)
 - ...
 - ...
 - ...
- README.txt

1.5.2 Prerequisites

The annotated data in TMT format requires tree editor TrEd to be installed. TrEd can be downloaded from,

<http://ufal.mff.cuni.cz/~pajas/tred/>

After installing TrEd, you need to install TMT plugin to browse the treebank data in the tmt file. The plugin can be installed by following,

- Start TrEd
- In menu, [Setup -> Manage Extensions -> TMT files support]
- Then treebank can be browsed by opening TamilTB.v0.1.tmt in TrEd

1.5.3 TectoMT Development Platform

The whole annotation work including preprocessing tasks (sentence segmentation, tokenization), manual annotation using tree editor and simple rule based parsing to bootstrap annotation have been implemented in TectoMT framework [Žabokrtský et al., 2008]. TectoMT is a highly modular NLP (Natural Language Processing) software system implemented in Perl programming language under Linux. It can be used for implementing various NLP tasks such as Machine Translation, POS tagging, parsing etc. in a modular way. For more about TectoMT, please refer the following URL,

<http://ufal.mff.cuni.cz/tectomt/>

1.5.4 Contact Authors

We welcome comments and suggestions for improvement of our current release. Please contact us for any suggestions or trouble in working the data.

Author Name	Email	Homepage
Loganathan Ramasamy	ramasamy@ufal.mff.cuni.cz	http://ufal.mff.cuni.cz/~ramasamy/
Zdeněk Žabokrtský	zabokrtsky@ufal.mff.cuni.cz	http://ufal.mff.cuni.cz/~zabokrtsky/

1.5.5 License

TamilTB.v0.1 by Institute of Formal and Applied Linguistics (UFAL) is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 3.0 Unported License. More information about the license can be obtained from,

<http://creativecommons.org/licenses/by-nc-sa/3.0/> (Read this first)
<http://creativecommons.org/licenses/by-nc-sa/3.0/legalcode>



Chapter 2

Morphological Annotation

This chapter gives a detailed description of annotation at word level. The annotation at this layer is roughly equivalent to Part of Speech (POS) tagging. The chapter begins with the introduction of positional tag, a format for tagging wordforms. Then the chapter introduces positional tagset for tagging Tamil data with examples.

2.1 Positional Tagging

For m-layer annotation (aka POS tagging in general), we chose to annotate separate word tokens with morphological features in addition to single main POS. Having morphological features would be ideal and necessary for morphologically rich languages such as Tamil. Just by knowing those features it may be possible to identify certain syntactic phenomena (such as case markers can identify syntactic relations). So to include morphological information (such as case, person, number and etc.) in addition to main POS, we decided to adopt the positional tagging scheme which has been successfully applied for the Czech language.

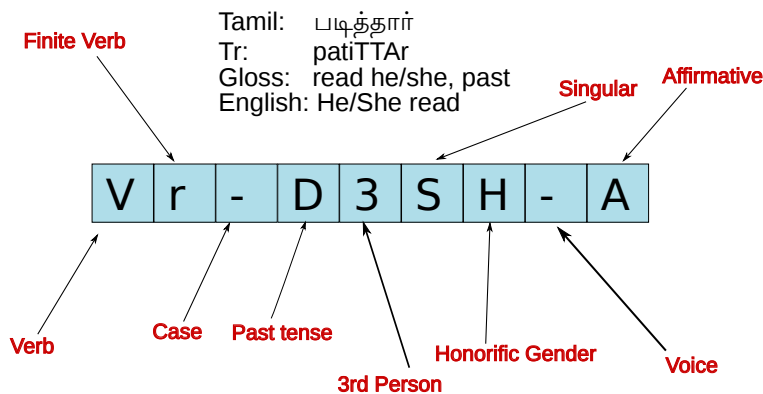


Figure 2.1: Positional tag

In positional tagging, each token (word) is tagged with a fixed length string. Each position or character in the tag represents a particular morphological feature of the token. The first position or character of the tag indicates the broad word category (such as noun, verb etc) to which the token belongs and the 2nd position indicates the detailed POS (for ex, what kind of verb? whether finite or non finite). The original Czech positional tagset includes 15 positions, and we designed a 9 positional tagset for Tamil. The Figure 2.1 illustrates the tagging of a Tamil word using positional tagging scheme.

The Figure 2.2 illustrates our positional tagset. The table in Figure 2.2(a) describes what each position occupies and the number of possible values they can take. The table in Figure 2.2(b) enumerates the possible values for *main POS* (first position). The tagset has been inspired from [Lehmann, 1989]. Our tagset includes separate entries for pronouns, numeral, interjections, particles and punctuations.

Position	Feature	#Possible Values
1	POS	14
2	Sub POS	42
3	Case	10
4	Tense	05
5	Person	04
6	Number	03
7	Gender	06
8	Voice	02
9	Negation	02

(a) Each position & num of possible values

Value	Description
A	Adverbs
C	Conjunctions
D	Determiners
I	Interjections
J	Adjectives
N	Nouns
P	Postpositions
Q	Quantifiers
R	Pronouns
T	Particles
U	Numerals
V	Verbs
X	Unknown
Z	Punctuations

(b) Main POS tags

Figure 2.2: Positional tagset

2.1.1 Position 2 - Sub POS

The *sub POS* corresponds to more finer version of *main POS* i.e. it can give some more information about the *main POS*. For ex, the *main POS* may just indicate that the wordform belongs to verb, the *sub POS* will further indicate that whether the verb is finite or non-finite. The *main POS* and *sub POS* together indicate the finer version of a part of speech. The other positions will further describe other details such as *tense, case, person, gender* and etc.

General note:

- For the tag of a given wordform, very rarely the entire positions of the tag is used. The value ‘-’ has been used wherever the position is not applicable.
- For ex: In the tag *NNN-3SN-*, the positions *4,8 and 9* are marked with ‘-’ since they are not applicable for noun tag.

Table 2.1: SubPOS

POS	Sub POS	Description	Example	Tag	Comments
A	A	Adverbs, general	<i>maRupatiyum</i> <i>vEkam Aka</i>	AA- - - - - NNN-3SN- -, AA- - - - -	<i>maRupatiyum</i> ('again') <i>vEkam</i> ('fast, quick'). Most adverbs in Tamil can be formed by adding adverbial suffix <i>Aka</i> to the nouns. For example, <i>vEkamAka</i> ('quickly') can be formed by adding adverbial suffix <i>Aka</i> to the noun <i>vEkam</i> ('quick'). Though English doesn't split the adverb 'quickly' into 'quick' and 'ly', Tamil tokenizer splits those words with <i>Aka</i> suffix, because the suffix has other occurrences too where it appears as a separate word. Thus, only the adverbial suffix will receive the AA- - - - - tag and the noun will receive the noun tag.
C	C	Conjunctions	<i>maRRum</i> <i>allaTu</i>	CC- - - - - CC- - - - -	<i>maRRum</i> ('and') <i>allaTu</i> ('or')
D	D	Determiners	<i>ivTa</i> <i>awTa</i> <i>ewTa</i>	DD- - - - - DD- - - - - DD- - - - -	<i>ivTa</i> ('this') <i>awTa</i> ('that') <i>ewTa</i> ('which thing')
I	I	Interjections	<i>AhA</i>	II- - - - -	<i>AhA</i> ('oh')
J			<i>azak Ana</i>	NNN-3SN- -, JJ- - - - -	<i>azaku</i> ('beauty'). Adjectives can also be formed by adding adverbial suffix <i>Ana</i> to nouns. In that case, the nouns will receive the noun tag and the adverbial suffix will receive the adjective tag JJ- - - - -
d		Adjective participle, adjectival verbs	<i>OtukiRa</i>	Jd-P- - - - A	<i>OtukiRa</i> ('... which/who is running')
N	N	Common nouns	<i>Otiya</i> <i>Otum</i> <i>kARRu</i> <i>kARRai</i> <i>kARRAl</i> <i>kautamA</i> <i>pirAk</i>	Jd-D- - - - A Jd-F- - - - A NNN-3SN- - NNA-3SN- - NNI-3SN- - NEN-3SH- - NEN-3SH- -	<i>Otiya</i> ('... which/who ran') <i>Otum</i> ('... which will run') <i>kARRu</i> ('air') Accusative Instrumental <i>kautamA</i> - 'Gautama' <i>pirAk</i> - 'Prague'
E	E	Proper nouns			

- continued on next page

Table 2.1 – continued from previous page

POS	Sub POS	Description	Example	Tag	Comments
	P	Participial nouns	<i>vAzwTavarkaL</i>	NPN-3PA--	<i>vAzwTavarkaL</i> - ‘those who have lived’
	O	Oblique nouns	<i>TirumaNa</i>	NO- -3SN--	When two nouns co occur the first noun will undergo certain change depending on the noun ending. They will act as a modifier to the second noun. Though they act as modifiers they will not be tagged as adjectives. They are tagged as oblique nouns. For ex: the noun <i>TirumaNam</i> / ‘wedding’ becomes <i>TirumaNa</i> / ‘wedding’ when it occurs as modifier to another noun say ‘invitation’. So it becomes, <i>TirumaNa azaippiTaz</i> / ‘wedding invitation’
P	P	Postpositions	<i>mItu</i> <i>itamiruwTu</i>	PP----- PP-----	<i>mItu</i> - ‘on’ <i>itamiruwTu</i> - ‘from’ Other examples: pati - ‘according’, pinnar - ‘later, after’, mUlamAka - ‘through’
Q	Q	quantifiers	<i>aTikam</i> , <i>konjcam</i> , <i>mikavum</i> , <i>mazu</i>	QQ-----	(‘a lot’), <i>konjcam</i> (‘a bit’) etc.
R	P	Personal nouns	<i>avar</i>	RpN-3SH--	<i>avar</i> - ‘he or she honorific’
	h	Reflexive nouns	<i>iTu</i> <i>wAm</i> <i>TanaTu</i>	RpN-3SN-- RpN-IPA-- RhG-3SA--	<i>iTu</i> - ‘this’ <i>wAm</i> - ‘we’ <i>TanaTu</i> - “oneself’s”
	B	General referential pronouns	<i>Tannai</i> <i>yArum</i>	RhA-3SA-- RBN-3SA--	<i>Tannai</i> - ‘about oneself’ <i>yArum</i> - ‘anyone’
	F	Specific indefinite referential pronouns	<i>eTuvum</i> <i>yArukkum</i> <i>yArO</i>	RBN-3SN-- RBD-3SA-- RFN-3SA--	<i>eTuvum</i> - ‘any thing’ <i>yArukkum</i> - ‘to anyone’ <i>yArO</i> - ‘someone’
	I	Inerrogative pronouns	<i>eTuvO</i> <i>yAr</i>	RFN-3SN-- RiN-3SA--	<i>eTuvO</i> - ‘something’ <i>yAr</i> - ‘who’

– continued on next page

Table 2.1 – continued from previous page

POS	Sub POS	Description	Example	Tag	Comments
	G	Non specific indefinite nouns	<i>eTu</i> <i>yArAvaTu</i>	RiN-3SN-- RGN-3SA--	<i>eTu</i> - ‘which thing’ <i>yArAvaTu</i> - ‘at least someone’
T	b	Comparative particle	<i>eTAvaTu</i>	RGN-3SN--	<i>eTAvaTu</i> - ‘at least something’
	d	Adjectival particle, adjectivalized verbs	<i>kAttikum, vita</i> <i>enkiRa, enRa</i>	Tb----- Td-P----A	<i>kAttikum, vita</i> (‘than’) <i>enkiRa</i> (‘the so called’)
	e	Interrogative particle	<i>A</i>	Te-----	<i>A</i> (‘is it ...?’)
	g	Adverb and Adjectival suffix	<i>Aka, Ana</i>	Tg-----	<i>Aka, Ana</i>
	k	Intensifier particle	<i>E, EyE, TAn</i>	Tk-----	<i>E, TAn</i> (‘indeed’)
	l	Clitic	<i>AvaTu</i>	Tl-----	<i>AvaTu</i> (‘at least’) or it can represent indefiniteness in combination with pronouns.
	m	Clitic (<i>limit</i>)	<i>mattum</i>	Tm-----	<i>mattum</i> (‘limit’)
	n	Complementizing nouns	<i>pOTu</i> <i>utan</i>	Tn----- Tn-----	<i>pOTu</i> - ‘when, during, while’ <i>utan</i> - ‘as soon as’
	o	Particle of doubt	<i>rAmanQ</i>	To-----	Other examples: <i>piRaku</i> - ‘after’, <i>mun</i> - ‘before’, <i>varai</i> - ‘till, as long as’, <i>mATiri</i> - ‘manner’, <i>pati</i> - ‘manner’, <i>rAmanQ</i> (‘whether it is Raman ...’)
	q	Emphatic particle	<i>E, EyE</i> and <i>TAn</i>	Tq-----	<i>E, EyE</i> and <i>TAn</i> . For ex: <i>ennaigyE</i> (‘just me’)
	Q	Complementizer	<i>enpaTu</i>	TQ-----	<i>enpaTu</i>
	s	Concessive particle	<i>Otiyum</i>	Ts-----	<i>um</i> (‘although, even if’). For ex: <i>Otiyum</i> (‘although ran’)
	S	Immediacy particle	<i>vawTaTum</i>	TS-----	<i>um</i> (‘as soon as’). For ex: <i>vawTaTum</i> (‘as soon as came’).

– continued on next page

Table 2.1 – continued from previous page

POS	Sub POS	Description	Example	Tag	Comments
	t	Complementizer in verbal participle	<i>ena, enRu</i>	Tt - T - - - - A	<i>ena, enRu.</i>
	v	inclusive participle	<i>um, kUta</i>	Tv - - - - - - - -	<i>rAmanum</i> ('also Ram')
	w	complementizer in conditional form	<i>enRAL</i>	Tw - T - - - - A	<i>enRAL</i>
U	x	cardinals	<i>onRu, iraNtu</i>	Ux - - - - - - - -	<i>onRu, iraNtu</i> ('one, two').
	y	ordinals	<i>onRAM, iraN-tAm</i>	Uy - - - - - - - -	<i>onRAM, iraNtAm</i> ('first, second')
	=	digits	<i>10, 20</i>	U = - - - - - - - -	numbers using digits.
V	j	imperative verb (lexical)	<i>Otu, uTavu</i>	Vj - T2PAAA	<i>irungkaL</i> ('you/[polite] wait')
	r	finite verb (lexical)	<i>OtukiRAN, uTavukiRAL</i>	Vr - P3SMAA	<i>OtukiRAN</i> ('he runs/running')
	t	verbal participle (lexical)	<i>Oti, vawTu</i>	Vt - T - - - - AA	<i>Oti</i> ('having run')
	u	infinitive (lexical)	<i>Ota, uTava</i>	Vu - T - - - - AA	<i>Ota, uTava</i> ('to run, to help')
	w	conditional verb (lexical)	<i>OtinAl, vawTAL</i>	Vw - T - - - - AA	<i>OtinAl, vawTAL</i> ('if runs/ran, if comes/came')
	z	verbal nouns (lexical)	<i>pitippaTu</i>	VzNF3SNAA	<i>pitippaTu</i> ('capturing')
R		finite verb (auxiliary)	<i>iruwTAr</i>	VR - D3SHAA	<i>iruwTAr</i>
T		verbal participle (auxiliary)	<i>koNtu</i>	VT - T - - - - AA	<i>koNtu</i>
U		infinitive (auxiliary)	<i>patta</i>	VU - T - - - - AA	<i>patta</i>
W		conditional (auxiliary)	<i>pattAl</i>	VW - T - - - - AA	<i>pattAl</i>

– continued on next page

Table 2.1 – continued from previous page

POS	Sub POS	Description	Example	Tag	Comments
	Z	verbal nouns (auxiliary)	<i>iruppaTu</i>	VZNF3SNAA	<i>iruppaTu</i>
Z	# :	terminal symbol other symbols	.(dot) , , ? , -, (,) , ! etc.	Z# - - - - - Z: - - - - -	sentence terminating symbol (usually <i>period</i>) all other symbols other than sentence termination symbol

2.1.2 Position 3 - Case

Tamil case occupies third position in the positional tag. The Table 2.2 lists Tamil *case markers* and the corresponding values for the 3rd position of the positional tag. There is a disagreement over the number of *case markers* in Tamil. [Lehmann, 1989] defines 9 cases. Some of the *case markers* can be considered as a *bound postpositions*. For that reason, we have not included the cases such as *ablative* and *benefactive* which are treated as postpositions.

#	Value	Description	Example	Tag
1	A	Accusative	<i>katciyai</i> ('party')	NNA - 3SN - -
2	D	Dative	<i>vIttukku</i> ('to/for the house')	NND - 3SN - -
3	I	Instrumental	<i>muyaRciyAl</i> ('by the efforts')	NNI - 3SN - -
4	G	Genitive	<i>aracin</i> ('government's')	NNG - 3SN - -
5	L	Locative	<i>pOril</i> ('in the war')	NNL - 3SN - -
6	N	Nominative	<i>ANtu</i> ('year')	NNN - 3SN - -
7	S	Sociative	<i>TuNaiyOtu</i> ('with the help')	NNS - 3SN - -

Table 2.2: Tamil case

2.1.3 Position 4 - Tense

#	Value	Description	Example	Tag
1	D	past	<i>kattinAr</i> ('built he')	Vr - D3SHAA
2	F	future	<i>uTavum</i> ('it will help')	Vr - F3SNAA
3	P	present	<i>celkiRAr</i> ('he is going')	Vr - P3SHAA
4	T	tenseless	<i>illai</i> ('exist not')	Vr - T3PNAA

Table 2.3: Tense

Note: Tenseless verbs

1. Tenseless verbs are tagged with **T**.
2. Tenseless verbs are unmarked in the wordform for tense.
3. Examples: *infinitives*, *verbal participles*, *imperatives* and some exceptional verbs.

2.1.4 Position 5 - Person

#	Value	Description	Example	Tag
1	1	1 st person	<i>mERkoNtEn</i> ('I undertook')	Vr - D1SAAA
2	2	2 nd person	<i>anjcukiRIrkaL</i> ('you fear')	Vr - P2PAAA
3	3	3 rd person	<i>vivATikkum</i> ('it will discuss')	Vr - F3SNAA
4	X	<i>unused</i>	<i>unused</i>	<i>unused</i>

Table 2.4: Person

The value *X* is reserved for future purposes. At present, *X* has not been used, rather the tag - is used the category is not applicable.

2.1.5 Position 6 - Number

#	Value	Description	Example	Tag
1	P	plural	<i>vivarangkaL</i> ('details')	NNN - 3PN - -
2	S	singular	<i>nyUSilAwTu</i> ('New Zealand')	NEN - 3SN - -
3	X	<i>unused</i>	<i>unused</i>	<i>unused</i>

Table 2.5: Number

2.1.6 Position 7 - Gender

#	Value	Description	Example	Tag
1	F	feminine	<i>varuvAL</i> ('she will come')	Vr - F3SFAA
2	M	masculine	<i>Atavanin</i> ('man's')	NNG - 3SM -
3	N	neuter	<i>etuTTaTu</i> ('it took')	Vr- D3SNAA
4	H	honorific (both masc. and fem.)	<i>avar</i> ('he/she [polite]')	RpN - 3SH - -
5	A	animate (humans)	<i>yAr</i> ('who?')	RiN - 3SA - -
6	I	inanimate (non humans)	<i>unused</i>	<i>unused</i>
7	X	<i>unused</i>	<i>unused</i>	<i>unused</i>

Table 2.6: Gender

2.1.7 Position 8 - Voice

#	Value	Description	Example	Tag
1	A	active	<i>etuTTaTu</i> ('it took')	Vr- D3SNAA
2	P	passive	<i>patukiRaTu</i> ('being [verb]...')	VR - P3SNPA

Table 2.7: Voice

2.1.8 Position 8 - Negation

#	Value	Description	Example	Tag
1	A	affirmative	<i>pinpaRRa</i> ('to follow')	Vu - T - - - AA
2	N	negation	<i>mutiyATu</i> ('cannot')	VR - T3SN -A

Table 2.8: Negation

2.2 Annotation of Pronouns

Tamil pronouns are one of the closed but in combination with clitics produce various derived pronouns. In this section, we list all possible pronouns (in their combination) with their tags. More information about Tamil pronouns can be obtained from [Lehmann, 1989]. The listing of tags for all possible pronouns will be useful in annotation task.

2.2.1 Singular Referential Definite (Personal) Pronouns

#	Person/Number	Pronoun	Tag
1	1 st /singular	<i>wAn</i> ('I')	RpN - 1SA - -
2	1 st /plural	<i>wAm</i> ('we, exclusive')	RpN - 1PA - -
3	"	<i>wAngkaL</i> ('we, inclusive')	RpN - 1PA - -
4	2 nd /singular	<i>wI</i> ('you')	RpN - 2SA - -
5	"	<i>wIngkaL</i> ('you, honorific, singular')	RpN - 2SH - -
6	2 nd /plural	<i>wIngkaL</i> ('you, plural')	RpN - 2PA - -
7	3 rd /singular	<i>avan</i> ('that one - he')	RpN - 3SM - -
8	"	<i>ivan</i> ('this one - he')	RpN - 3SM - -
9	"	<i>avaL</i> ('that one - she')	RpN - 3SF - -
10	"	<i>ivaL</i> ('this one - she')	RpN - 3SF - -
11	"	<i>aTu</i> ('that one - it')	RpN - 3SN - -
12	"	<i>iTu</i> ('this one - it')	RpN - 3SN - -
13	"	<i>avar</i> ('that one - he/she hon.')	RpN - 3SH - -
14	"	<i>ivar</i> ('this one - he/she hon.')	RpN - 3SH - -
15	3 rd /plural	<i>avai/avaikaL</i> ('those ones')	RpN - 3PN - -
16	"	<i>ivai/ivaikaL</i> ('these ones')	RpN - 3PN - -
17	"	<i>avarkaL</i> ('those people')	RpN - 3PA - -
18	"	<i>ivarkaL</i> ('these people')	RpN - 3PA - -

Table 2.9: Personal pronouns

2.2.2 Non Referential (Interrogative) Pronouns

#	Pronoun	Tag
1	<i>yAr</i> ('who')	RiN - 3SH - -
2	<i>enna</i> ('what')	RiN - 3SN - -
3	<i>evan</i> ('which male person')	RiN - 3SM - -
4	<i>evaL</i> ('which female person')	RiN - 3SF - -
5	<i>eTu</i> ('which thing')	RiN - 3SN - -
6	<i>evar</i> ('which male/female person')	RiN - 3SH - -
7	<i>evarkaL</i> ('which persons')	RiN - 3PA - -
8	<i>evai(kaL)</i> ('which things')	RiN - 3PN - -

Table 2.10: Interrogative pronouns

2.2.3 General Referential Pronouns

In the case of general referential pronouns, the particle *-um* is added to interrogative pronouns which results in the addition of referential property to interrogatives. The resultant words will have a meaning of 'anyone, anybody or anything' in English. The Table 2.11 lists general referential pronouns.

Note

- Though the *inclusive* particle *-um* is added to pronouns, the tokenizer does not split the suffix from the pronoun unlike other situations.
- In most other cases, the particle *-um* is splitted from wordforms.

#	Pronoun	Tag
1	<i>yAr<u>um</u></i> ('anyone')	RBN - 3SH - -
3	<i>evan<u>um</u></i> ('anyone, male person')	RBN - 3SM - -
4	<i>evaL<u>um</u></i> ('anyone female person')	RBN - 3SF - -
5	<i>eTuv<u>um</u></i> ('anything')	RBN - 3SN - -
6	<i>evar<u>um</u></i> ('anyone, male/female person')	RBN - 3SH - -
7	<i>evarkaL<u>um</u></i> ('any persons')	RBN - 3PA - -
8	<i>evaikaL<u>um</u></i> ('anything, plural')	RBN - 3PN - -
8	<i>evaiy<u>um</u></i> ('anything, plural')	RBN - 3PN - -

Table 2.11: General referential pronouns

2.2.4 Specific Indefinite Referential Pronouns

#	Pronoun	Tag
1	<i>yAr<u>O</u></i> ('someone')	RFN - 3SH - -
3	<i>evan<u>O</u></i> ('some male person')	RFN - 3SM - -
4	<i>evaL<u>O</u></i> ('some female person')	RFN - 3SF - -
5	<i>eTuv<u>O</u></i> ('something')	RFN - 3SN - -
6	<i>evar<u>O</u></i> ('some male/female person')	RFN - 3SH - -
7	<i>evarkaL<u>O</u></i> ('someone, plural')	RFN - 3PA - -
8	<i>evaikaL<u>O</u></i> ('something, plural')	RFN - 3PN - -
8	<i>evaiy<u>O</u></i> ('something, plural')	RFN - 3PN - -

Table 2.12: Specific indefinite referential pronouns

2.2.5 Non-specific Indefinite Referential Pronouns

#	Pronoun	Tag
1	<i>yAr<u>AvaTu</u></i> ('someone or other')	RGN - 3SH - -
3	<i>evan<u>AvaTu</u></i> ('some male person or other')	RGN - 3SM - -
4	<i>evaL<u>AvaTu</u></i> ('some female person or other')	RGN - 3SF - -
5	<i>eT<u>AvaTu</u></i> ('something or other')	RGN - 3SN - -
6	<i>evar<u>AvaTu</u></i> ('some male/female person or other')	RGN - 3SH - -
7	<i>evarkaL<u>AvaTu</u></i> ('someone, or other (plural)')	RGN - 3PA - -
8	<i>evaikaL<u>AvaTu</u></i> ('something, or other (plural)')	RGN - 3PN - -
8	<i>evaiy<u>AvaTu</u></i> ('something, other (plural)')	RGN - 3PN - -

Table 2.13: Non-specific indefinite referential pronouns

Chapter 3

Syntactic Annotation

This section will give a detailed description about how annotation takes place at the syntactic level. The syntactic annotation consists of two phases: (i) identifying the structure (dependency) of the sentence in the form of dependency tree and (ii) identifying the dependency relations and assigning those relations to edges in the dependency tree structure.

3.1 Identifying the Structure

The structure of the sentence is identified manually by attaching the dependent nodes to the governing nodes. In the sentential structure, the head of the sentence will be predicate, and the predicate will have arguments (noun phrases, adverbials) as their children. The objective of this step would be, identifying the predicate rooted structure and attaching to the technical root (AuxS, defined below) of the tree. The end of the sentence will also be attached to the technical root. Once the structure is identified, all the edges have to be labeled with their relations. For technical reasons, the relation between the dependent and the governing node is stored as an a-layer attribute of the dependent node. The attribute is called '*afun*'. The following sections explain each dependency relation in detail.

3.2 Dependency Relations

According to PDT naming convention, dependency relations are also called as *analytical functions* or *afuns*. The documentation uses these names interchangeably.

#	Afun	Description	Comments
1	AAdjn	Adverbial Adjunct	Optional adverbs, optional PP phrases attaching to verb
2	AComp	Adverbial Complement	Obligatory adverbs, obligatory PP phrases attaching to verb
3	Apos	Apposition	Heads of the apposition clauses - clauses attaching to <i>enRa</i>
4	Atr	Attribute	Noun modifiers
5	AdjAtr	Adjectival participial	Adjectivalized verbs, or relative clauses
6	AuxA	Determiners	Demonstrative pronouns <i>iwTa</i> ('this'), <i>awTa</i> ('that')
7	AuxC	Subordinating conjunctions	Subordinating Conjunctions - <i>enRu</i> , <i>ena</i> , <i>Aka</i>
8	AuxG	Symbols other than comma	-, ' , ; \$, rU, (,],] etc.
9	AuxK	Sentence termination symbols	;, ., ?
10	AuxP	Postpositions	<i>mITu</i> ('on'), <i>paRRi</i> ('about'), <i>kiZ</i> ('under')
11	AuxS	Technical root	Technical root
12	AuxV	Auxiliary verb	<i>uL</i> , <i>koNtu</i> , <i>iru</i> etc.
13	AuxX	Comma (not coordination)	,
14	AuxZ	Emphasis words or particles	<i>TAn</i> (emphasis), <i>um</i> (also, even), <i>E</i> ('even')
15	CC	Part of a word	<i>kiLaruTu ezwuTu</i> - ('rising'). <i>kiLaruTu</i> will be labeled as <i>CC</i>
16	Comp	Complement other than attaching to verbs	Obligatory attachments to non verbs. Ex: "belongs to the batch of 1977"
17	Coord	Coordination node	<i>maRRum</i> ('and'), <i>um</i>
18	Obj	Object (both direct and indirect)	usually nouns with accusative case
19	Pnom	Predicate nominal	Nominals that act as main verbs
20	Pred	Predicate	Main verb (usually finite) of a sentence.
21	Sb	Subject	Subject usually nominals

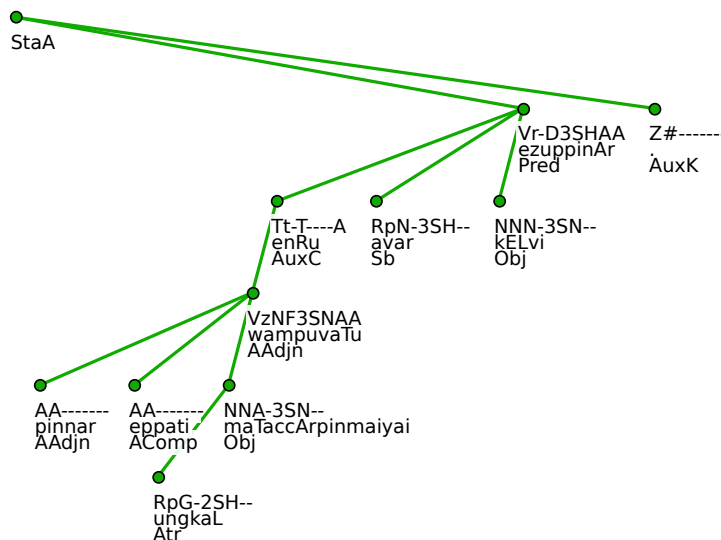
Table 3.1: Dependency relations (*afuns*)

3.3 Detailed Description of *afuns*

In this section, we describe each dependency relation in detail with example annotation of that relation. For the list of dependency relations or *afuns*, please refer the Table 3.1. Each dependency relation is explained with an analytical tree and the sentence it represents. The sentence is shown in Tamil script and it's transliteration, gloss and the actual English translation in that order. The word shown in bold will receive the dependency relation that is being explained.

3.3.1 *Afun: AAdjn*

The *AAdjn* relation is used to mark adverbial adjuncts. Adverbial adjuncts are optional adverbial phrases, prepositional phrases, clauses or simple adverbs modifying the verbs. Figures 3.1 & 3.3 show examples for *AAdjn* relation. In the Figure 3.1, the adverb *pinnar* ('later, after') has been labeled with *AAdjn* relation.



Tamil: பின்னர் எப்படி உங்கள் மதச்சார்பின்மையை நம்புவது என்று கேள்வி எழுப்பினார் .
 Tr: **pinnar** eppati ungkaL maTaccArpinmayai wampuvaTu enRu kELvi ezuppinAr .
 Gloss: **then** how your secularism to-believe - question raised-he .
 English: "Then how to believe your secular credentials ", he raised a question .

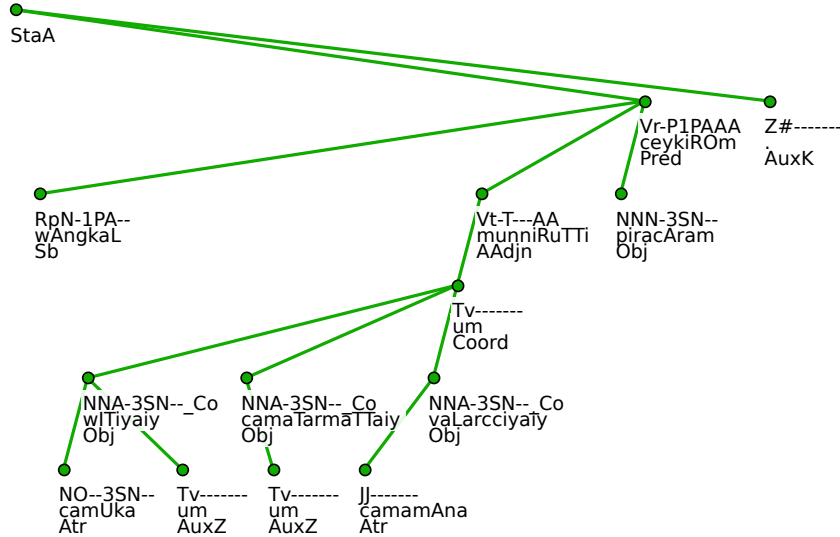
Figure 3.1: *AAdjn*: Adverbial adjunct

In the Figure 3.3, the adverbial phrase *munniRuTTi* ('by putting forward ...') has been labeled with *AAdjn* relation. The *AAdjn* label is determined by whether excluding the adverbial adjunct affects the meaning of the sentence. If it does not (only provides extra information about the sentence), the head of the phrase is assigned *AAdjn* relation.

3.3.2 *Afun: AComp*

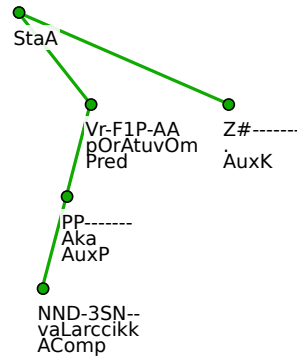
The *AComp* relation is used to mark the obligatory adverbials or adverbial complements in the structure. The context of occurrence of *AComp* relation is same as that of *AAdjn* relation. The only difference is that the *adverbial adjuncts* (*AAdjn*) are optional elements in the sentential structure whereas *adverbial complements* (*AComp*) are obligatory elements in the sentence structure. While doing annotation, this relationship is determined by whether the adverbial structure is required to complete the sentence. If the removal of the adverbial structure does not affect the sentence as a whole, then it is labeled as *AAdjn* otherwise it will be labeled as *AComp*.

- Our annotation treats the adverbials as *adverbial complements* (*AComp*) and *adverbial adjuncts* (*AAdjn*). Whereas the original PDT treats the both under *adverbials* (*Adv*).



Tamil: நாங்கள் சமூக நீதியை உம் சமதர்மத்தை உம் சமமான வளர்ச்சியை உம் முன்னிறுத்தி பிரசாரம் செய்கிறோம் .
 Tr: wAngkaL camUka wITiyai um camaTarnaTTai um camamAna vaLarcciyai um **munniRuTTi** piracAram ceykiROm .
 Gloss: we social justice - equality - equal development - **by-putting-forward** campaign doing-we .
 English: We are doing campaign **by putting forward** social justice, equality and equal development .

Figure 3.2: *AAdjn*: Adverbial adjunct



Tamil: வளர்ச்சிக்க ஆக போராடுவோம் .
 Tr: **vaLarccikk** Aka pOrAtuvOm .
 Gloss: **development** for struggle-we .
 English: We will struggle for **development** .

Figure 3.3: *AComp*: Adverbial complement

3.3.3 Afun: *Apos*

The adjectival clauses headed by *enRa* are appositional clauses. The entire finite clause will be attached to *enRa* which will act as a modifier to the following noun phrase. The clausal head which is attached to *enRa* will receive Apos label. The following Figure 3.4 illustrates the labeling of *Apos*.

3.3.4 Afun: *Atr*

Attribute [UFAL, 2006] is a sentence member which depends on noun and closely determines its meaning. Original PDT annotation differentiates “agreeing” and “non-agreeing” attribute. But, since Tamil does not have any agreement between nouns and their modifiers, all noun modifiers will receive the afun label *Atr*. The noun modifiers include nouns (except the head noun) in noun compounds, adjectives, numerals and adjectival participles. The Figure 3.5 shows the usage of afun *Atr*.

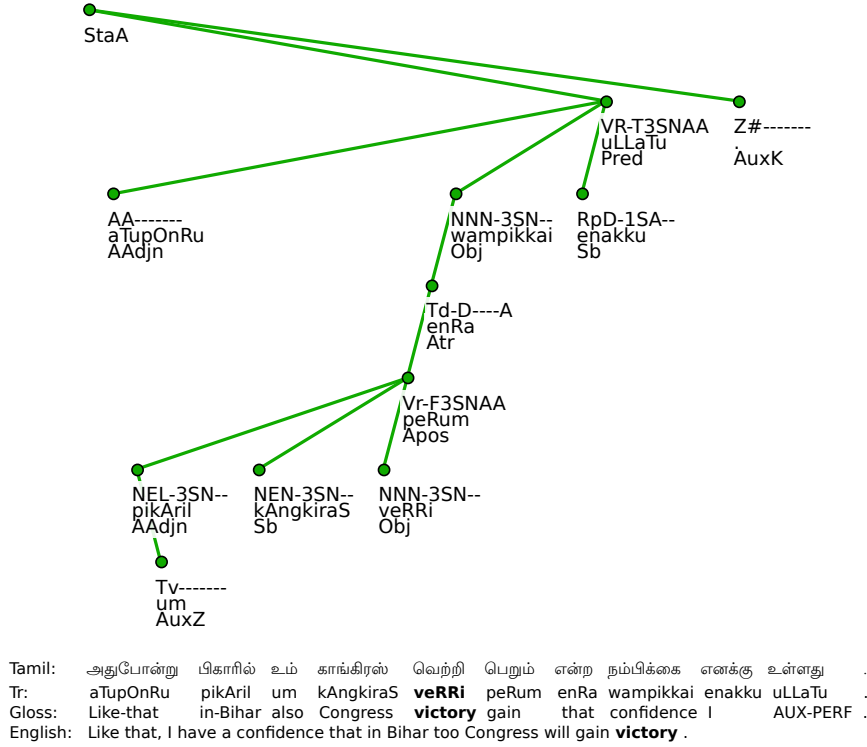


Figure 3.4: Apos: Apposition

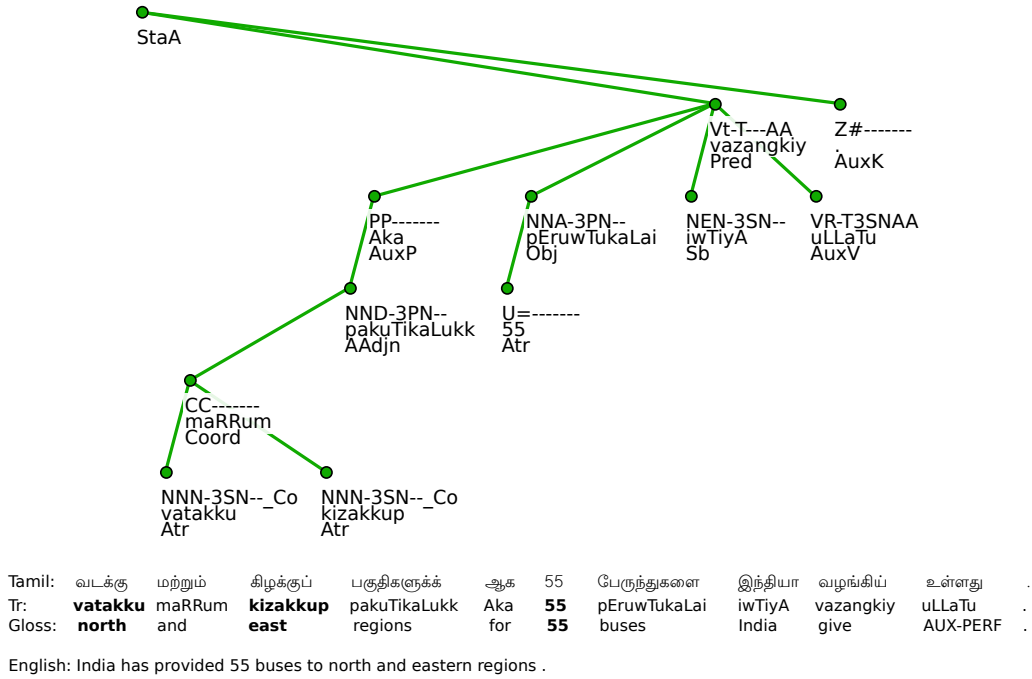


Figure 3.5: Atr: Attribute

3.3.5 Afun: AdjAtr

AdjAtr label is used to mark the adjectival clauses, adjectivalized verbs or adjectival participials. They are equivalent to -ing, -ed (singing girl, departed train) forms in English. Verbs in Tamil can be adjectivalized for all three tenses, and they take appropriate tags depending on the word form features. Figures 3.6 & 3.7 show examples for *AdjAtr* labeling.

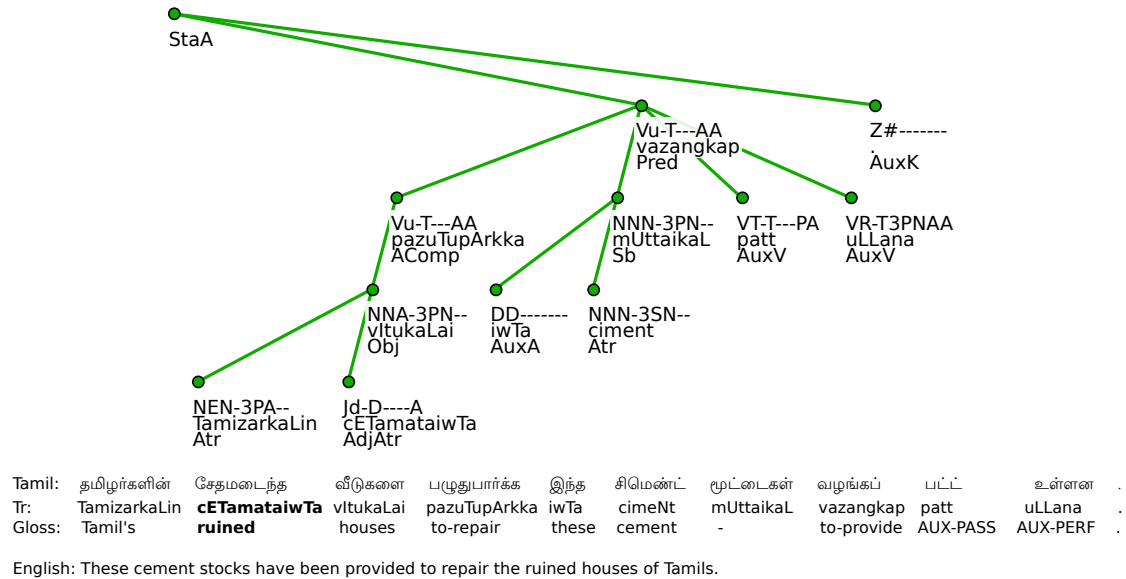


Figure 3.6: *AdjAtr*: Adjectival attribute

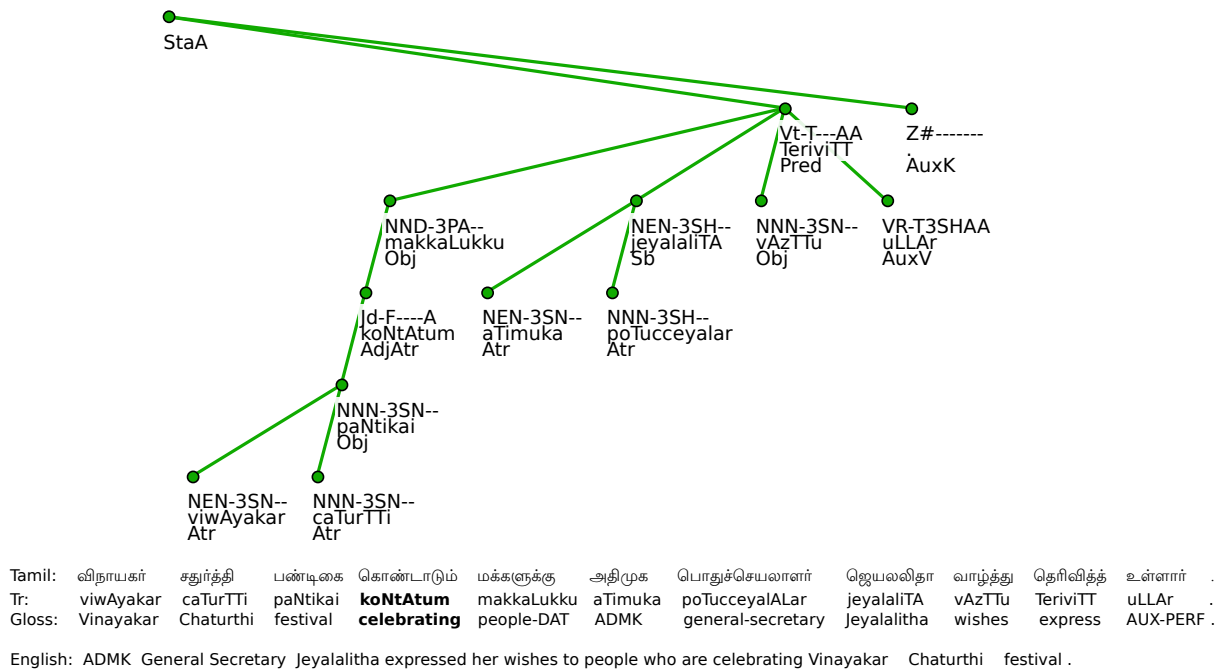
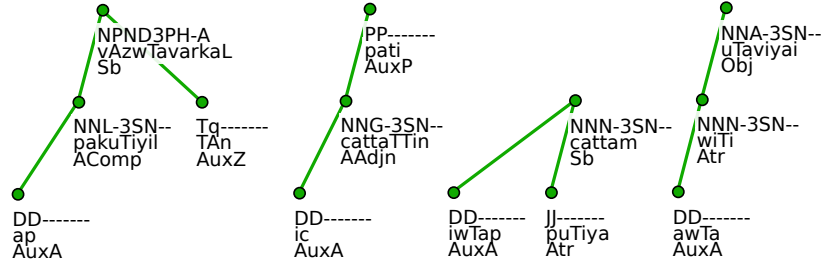


Figure 3.7: *AdjAtr*: Adjectival attribute

3.3.6 Afun: *AuxA*

Tamil has two demonstrative determiners corresponding to ‘this’ and ‘that’ in English. At present, the question word *ewTa* (‘which’) is also tagged as determiner. Sometimes the determiners occur as a prefix to the following noun or noun phrase in a contracted form. During the tokenization phase, these demonstrative suffixes will be separated from the noun or noun phrase, and will be considered as a separate token. The determiners are *iwTa* (‘this’), *awTa* (‘that’) and *ewTa* (‘which’), and the corresponding contracted determiners are *i*, *a* and *e*. Due to orthographic rules, the first letter of the following noun phrase will be added to the contracted form. The Figure 3.8 illustrates the usage of all determiners.



Tamil: அப் பகுதியில் வாழ்ந்தவர்கள் தான் ...
 Tr: **ap** pakuTiyil vAwTavarkaL TAn ...
 Gloss: **that** in-the-region who-had-lived indeed ...
 English: Indeed who had lived in **that** region

Tamil: இச் சட்டத்தின் படி
 Tr: **ic** cattaTTin pati
 Gloss: **this** law according-to
 English: according to **this** law

Tamil: இந்தப் புதிய சட்டம்
 Tr: **iwTap** puTiya cattam
 Gloss: **this** new law
 English: this new law

Tamil: அந்த நிதி உதவியை
 Tr: **awTa** wiTi uTaviyai
 Gloss: **that** financial help
 English: that financial help

Figure 3.8: *AuxA*: Determiners

3.3.7 Afun: *AuxC*

In Tamil, embedding or adjoining of clauses are performed either by morphologically marking the clause or by using separate words. When separate words are used, they function similar to that of subordinating conjunction words in other languages such as English. These separate words are called complementizers in Tamil. Complementizers can be verbs, nouns or postpositions after nominalized clauses. There are three complementizing verbs - *en* ('say'), *pOl* ('seem') and *Aku* ('become'). They have grammatical function during embedding of clauses, otherwise they retain their lexical meanings. The following list provides some of the noun complementizers - *pOTu* ('time'), *mun* ('before'), *piRaku* ('after'), *utan* ('immediacy'), *varai* ('as long as') and etc. The postpositions can also be interpreted as subordinating conjunction words when they are preceded by nominalized clauses. Refer [Lehmann, 1989] for detailed treatment of how complementizers work in Tamil.

3.3.8 Afun: *AuxG*

The symbols other than sentence boundary and comma are labeled with *AuxG* afun. The *AuxG* symbols include pairs of symbols such as (, }, [, ", ' and other symbols such as !, , #, \$ and etc. The Figure 3.11 shows an example for *AuxG*.

3.3.9 Afun: *AuxK*

The *AuxK* afun is assigned to sentence termination symbols. The symbols ., ? and : are considered as sentence terminals and they are expected at the end of a sentence.

3.3.10 Afun: *AuxP*

AuxP is used to mark the postpositions (heads) of the postpositional phrases. The postposition will receive the *AuxP* label and the element attached to *AuxP* will receive the afun (*AAdjn*, *AComp*, *Atr*, *Comp*) according to their context of occurrence. The Figures 3.13 & 3.14 show examples for *AuxP*.

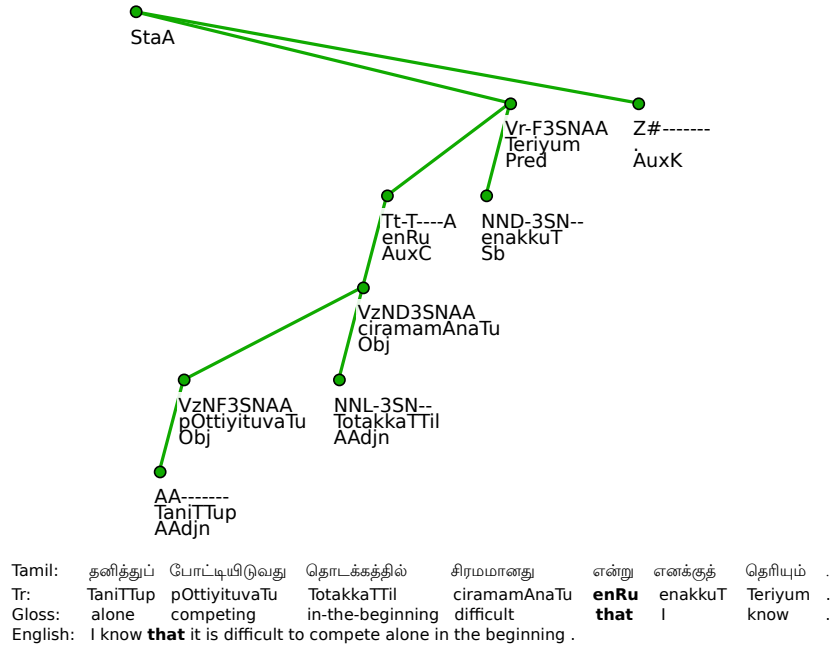


Figure 3.9: *AuxC*: Subordinating conjunctions

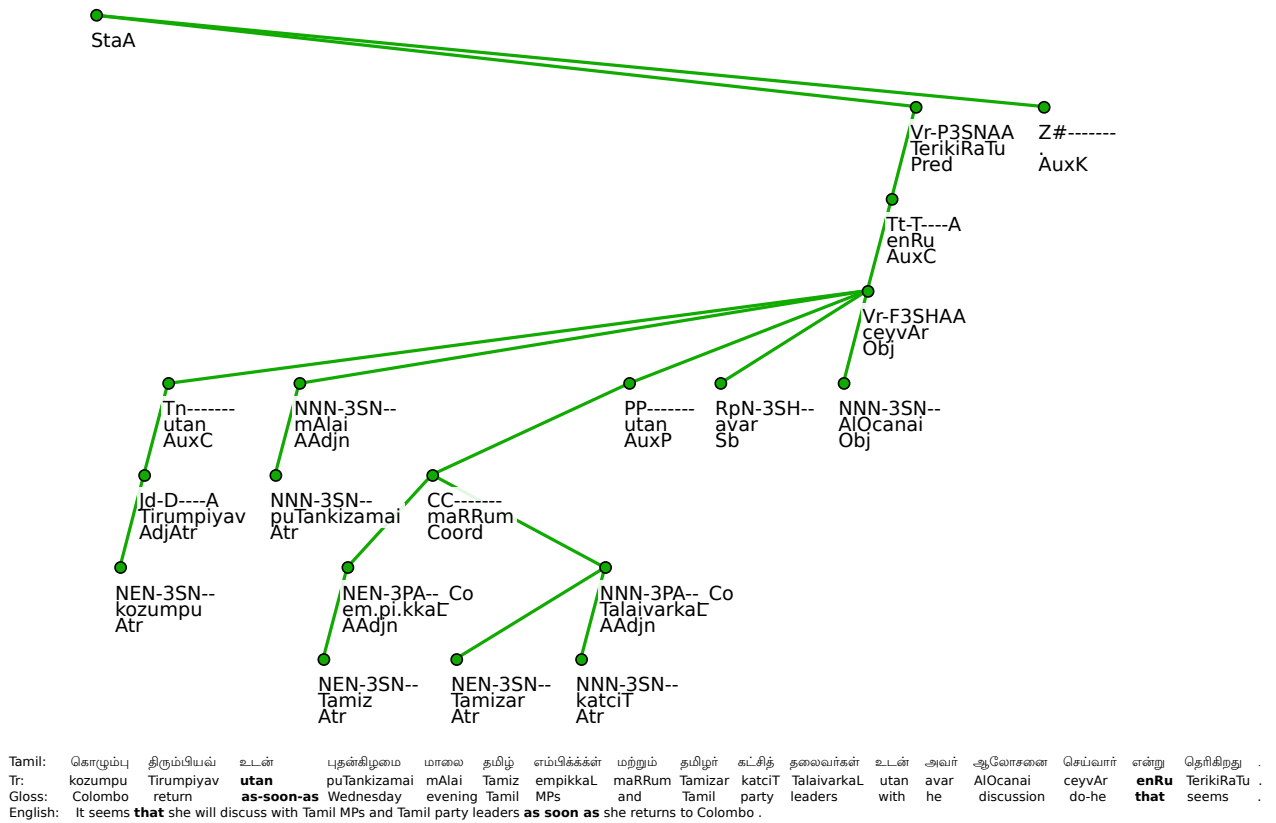


Figure 3.10: *AuxC*: Subordinating conjunctions

3.3.11 Afun: *AuxS*

AuxS is used to label the technical root of a tree. In the Figure 3.14, the technical root (shown as *StaA*) is the root of the tree structure. To this node, the predicate and the sentencing ending node will be attached.

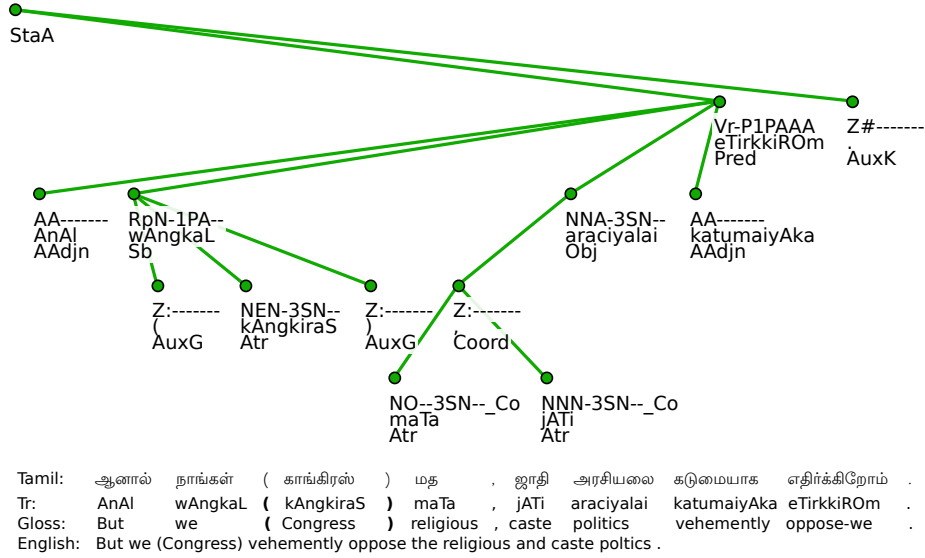


Figure 3.11: *AuxG*: Symbols

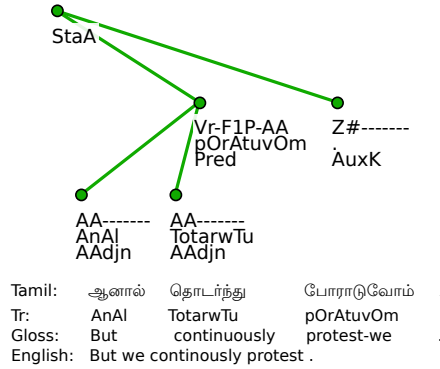


Figure 3.12: *AuxK*: Sentence termination symbol

3.3.12 Afun: *AuxV*

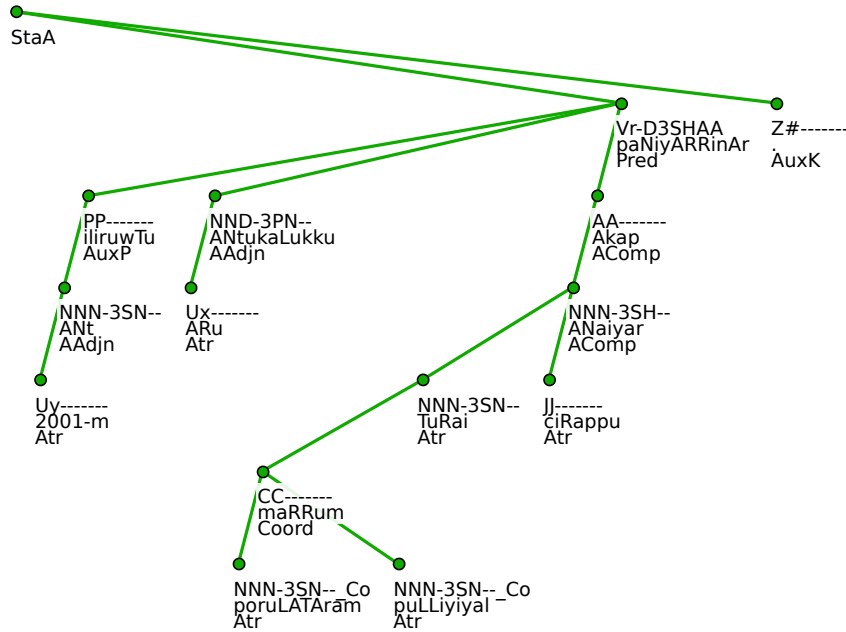
Auxiliary verbs are assigned the afun *AuxV*. In compound verb constructions, the auxiliary verb will be hanged under lexical verb. All auxiliary words including passive constructions will receive the afun *AuxV*. In the Figure 3.15, there are two auxiliary verbs *patu* (‘experience’) and *uL* (‘exist’). The auxiliary *patu* (‘experience’) and *uLLaTu* (‘exist’) are labeled with *AuxV*. The lexical verb receive the label *Pred* if there is only one clause in the sentence, otherwise the label of the lexical verb depends on the upper clauses.

3.3.13 Afun: *AuxX*

AuxX afun is used to label commas. All commas except the commas which act as coordination head is labeled with *AuxX*. The Figure 3.16 shows how *AuxX* has been annotated in the coordination structure.

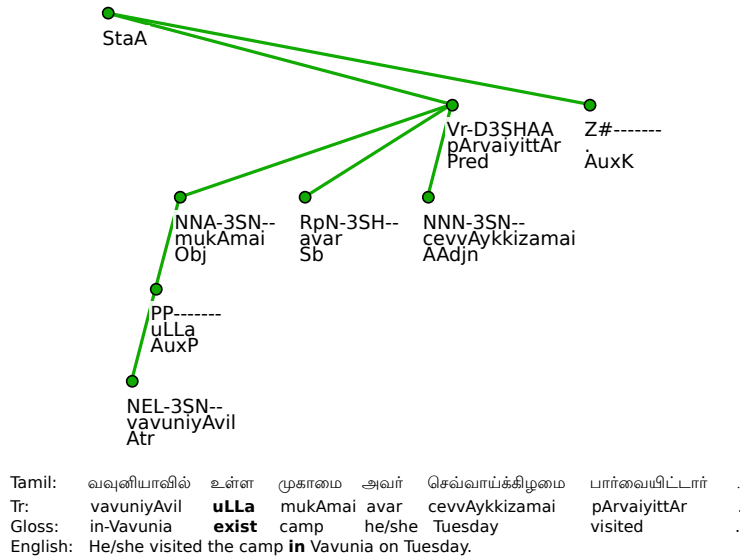
3.3.14 Afun: *AuxZ*

The clitics *um*, *E*, *TAn*, *mattum* will receive the afun *AuxZ*. The Figures 3.17 & 3.18 illustrate the labeling of *AuxZ*. The afun *AuxZ* corresponds to emphasizing words in English such as (‘just’, ‘only’, ‘indeed’) etc. These are not separate words but clitics in Tamil. Among the clitics, *um* has various semantics functions including acting as *coordination head*. In the Figure 3.17, the clitic *um* (‘also’) adds an inclusive meaning to the sentence. In the Figure 3.18, the clitic *TAn* (‘only’) emphasizes the entire clause.



Tamil: 2001-ம் ஆண்ட் இலிருந்து ஆறு ஆண்டுகளுக்கு பொருளாதாரம் மற்றும் புள்ளியியல் துறை சிறப்பு ஆணையர் ஆகப் பணியாற்றினார் .
 Tr: 2001-m ANT **iliruwTu** ARu ANTukaLukku poruLATArAm maRRum puLLiyiyal TuRai ciRappu ANaiyar Akap paNiyARRinAr .
 Gloss: 2001 year **from** six years-for economics and statistics department special commissioner as worked-he .
 English: He/she worked as special commissioner of economics and statistics department for six years from 2001 .

Figure 3.13: *AuxP*: Postpositions



Tamil: வவுனியாவில் உள்ள முகாமம் அவர் செவ்வாய்க்கிழமை பார்வையிட்டார் .
 Tr: vavuniyAvil **uLLa** mukAmai avar cevvAykkizamai pArvaiyittAr .
 Gloss: in-Vavunia **exist** camp he/she Tuesday visited .
 English: He/she visited the camp in Vavunia on Tuesday.

Figure 3.14: *AuxP*: Postpositions

3.3.15 Afun: CC

The label *CC* is used to mark in places where a single lexical unit is composed of multiple words. The *CC* relation of a word would indicate that the current word together with its parent word form a single lexical unit. In Tamil, a single action verb can be split into multiple words. But during annotation, the problem arises as to which part of the word sequence the arguments of that lexical unit should be attached. To resolve this issue, we treat the first word as having the lexical meaning and the remaining to be children of the first word. In that case, the first word would receive the afun corresponding to the entire lexical unit and the remaining words would receive the label *CC*. For ex, the Figure 3.19 indicates that the Tamil verb *veRRipeRu* ('to win') is splitted into 2 words - *veRRi peRu*. The first word *veRRi* receives the afun which is meant for the entire lexical unit comprising the following words with afun *CC*.

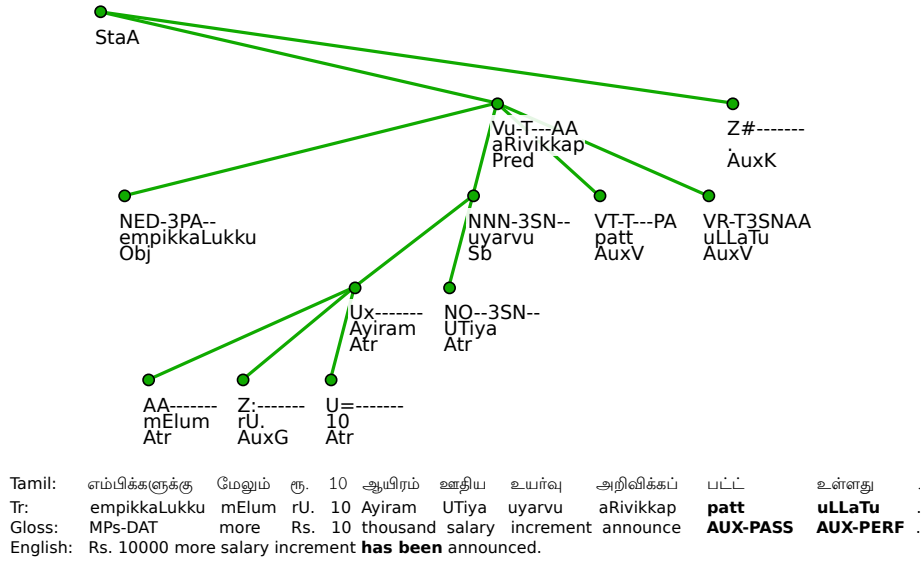


Figure 3.15: *AuxV*: Auxiliary verb

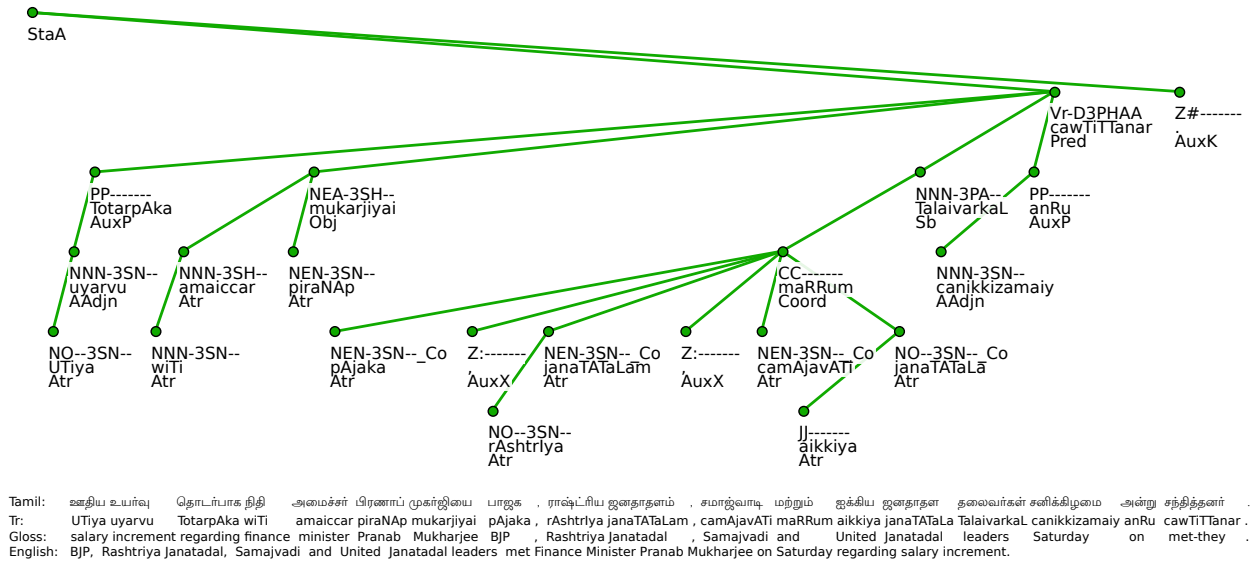


Figure 3.16: *AuxX*: Commas

Note:

- The afun *CC* is not a relation for Multi Word Expressions (MWE).
- In many cases, Tamil verbs are formed using *noun + auxiliary* combination which acts as a single verb. In this case, the auxiliary would receive *CC* and would be attached under the noun.
- The word sequence can be written together as a single word. This implies that the word sequence is not a MWE.

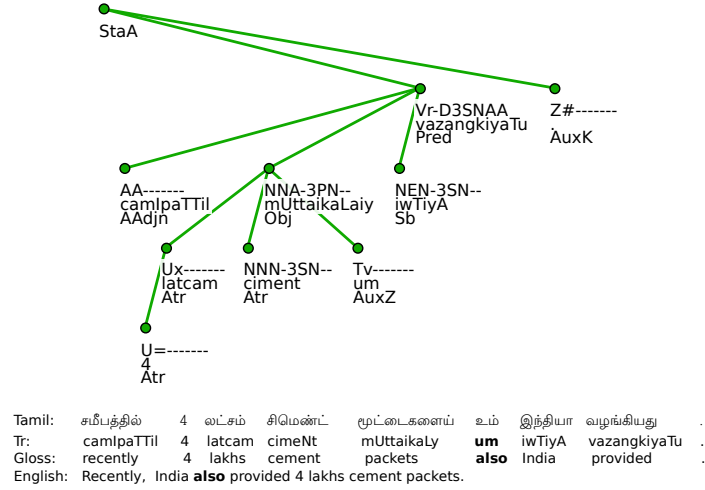


Figure 3.17: *AuxZ*: Emphasis

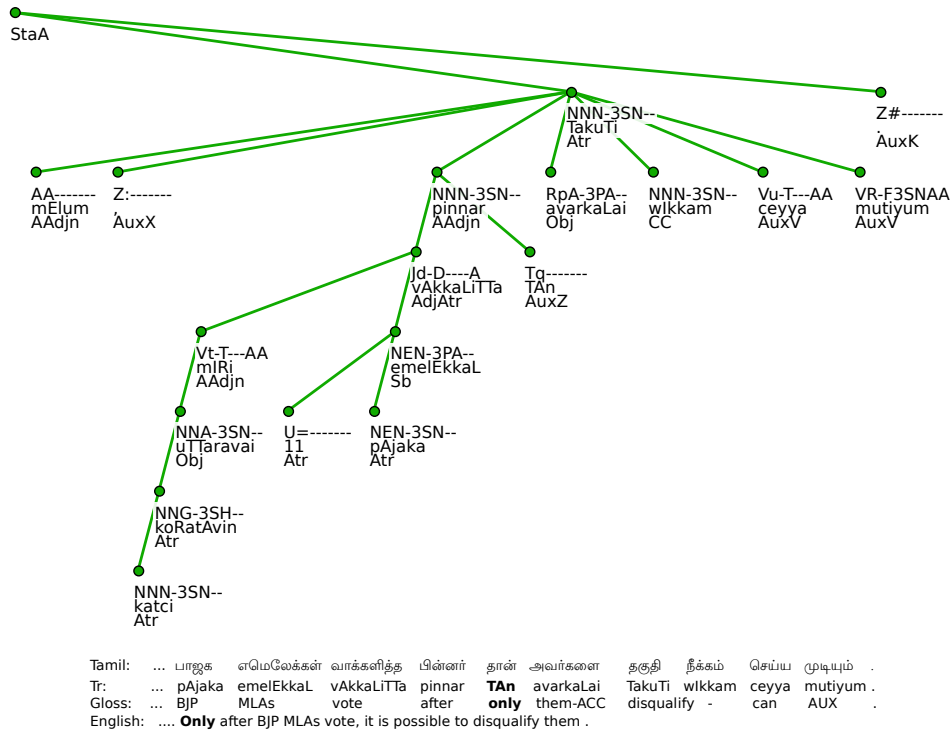


Figure 3.18: *AuxZ*: Emphasis

3.3.16 Afun: Comp

Comp label is used to mark the obligatory element not attaching to verbs. For example, consider the phrase *1200kk um mERpatta poTumakkaL uyirizawT uLLanar* (“more than 1200 people have been died”), in that phrase, *1200kk* occurs as an obligatory argument to *mERpatta* (‘more than’). So, *1200kk* will be labeled with *Comp*. Even nouns (not modifiers) which obligatorily attach to other nouns are labeled with *Comp* afun. Other occurrences of *Comp* is when postpositional phrase (PP phrase) attaches to a noun phrase. The postpositional head will receive *AuxP* label whereas the head noun phrase of the PP phrase will receive *Comp* label. The Figure 3.20 illustrates the labeling *Comp* relation.

3.3.17 Afun: Coord

Coordination is one of the complex phenomena in Tamil. Coordination conjunction in Tamil can be performed using at least 2 different ways. In the first method (for ‘and’ coordination), all conjoining

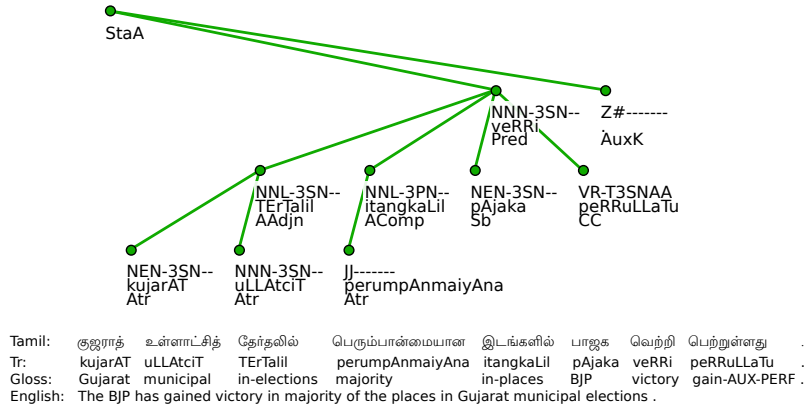


Figure 3.19: *CC*: Marker for a multi word sequence

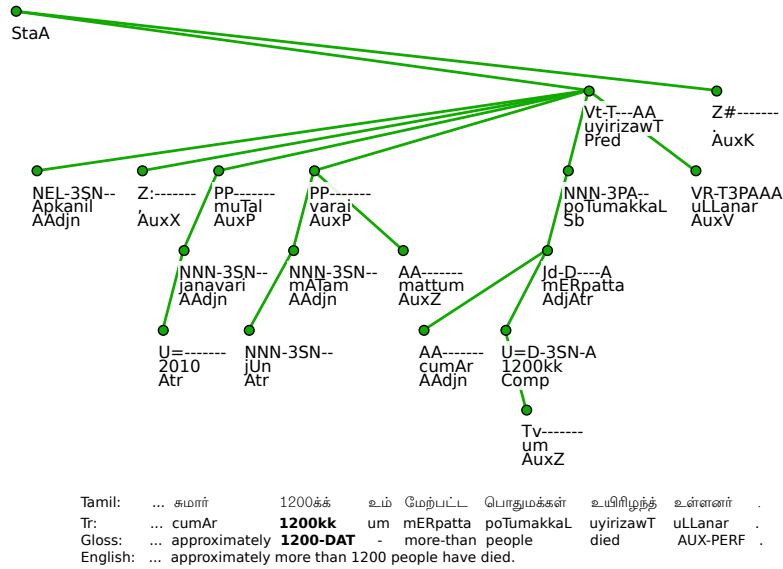


Figure 3.20: *Comp*: Complement (other than verbs)

elements adds the inclusive particle *um* ('also') at the end of the word form. Thus in this method, all conjoining elements possess the suffix (um) which would indicate the coordination is taking place. Moreover, the 'is_member' attribute of the conjoining elements will be set to 1. The separator (comma) between elements is optional. It is perfectly legitimate if there is no comma between any of the conjoining elements.

Second method for Tamil 'and' coordination is similar to English style 'and' coordination. The conjunction word *maRRum* ('and') is added between conjoining elements. If there are more than 2 elements, then *maRRum* ('and') will be added just before the last conjunct. The other elements will be separated by comma. Again, the 'is_member' attribute of the conjoining elements will be set to 1.

The 'or' coordination is performed in a similar way for both the methods. For the first method, the suffix *-O* is added to all conjuncts, and for the second method, the conjunction word *allaTu* ('or') is added between the last 2 conjoining elements. The remaining elements will be separated using comma.

Apart from the above two main methods, the coordination can be done via with just commas. In that case, the comma between conjuncts will act as coordination head. The Figures 3.21, 3.22 and 3.23 illustrates how coordination head is marked in the above mentioned scenarios.

3.3.18 Afun: *Obj*

Direct and indirect objects receive the afun *Obj*. The Figure 3.24 shows an example annotation of sentence with *Obj* relation. If there are both direct and indirect objects in a same sentence, then both will be labeled with afun *Obj*.

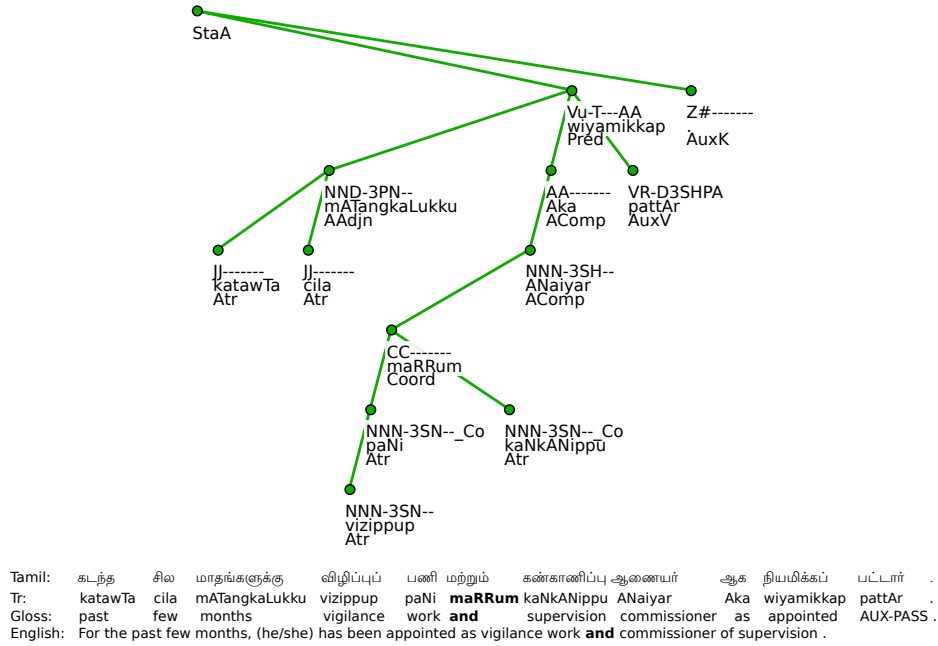


Figure 3.21: *Coord*: Coordination head (English style)

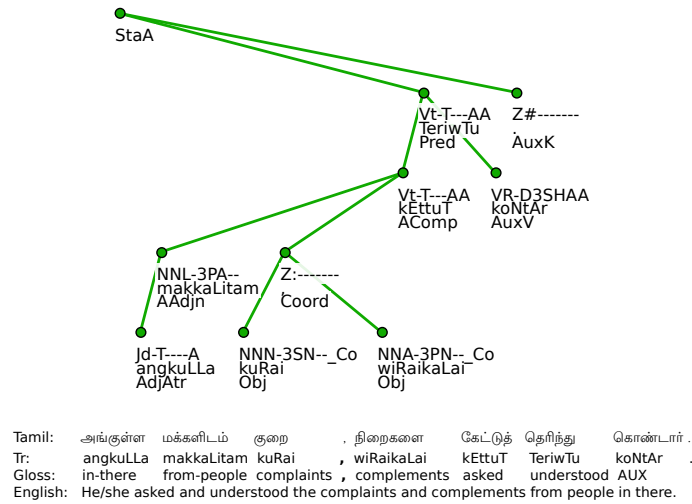


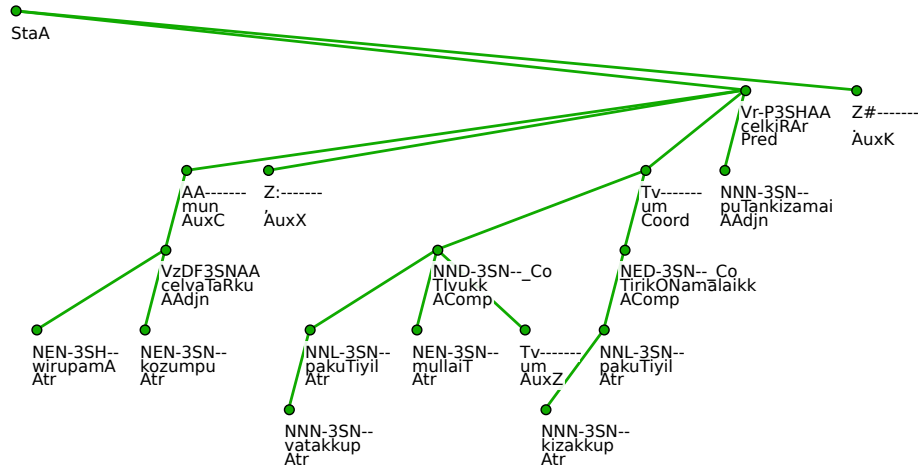
Figure 3.22: *Coord*: Coordination head (Comma)

3.3.19 Afun: *Pnom*

Nominal predicate occurs in the copula (be) constructions or verbless constructions. In these constructions, the sentence will not have any lexical verb. Instead, the predicate will contain only noun phrase. The noun phrase will be the predicate, and the afun label *Pnom* will be assigned to the noun phrase. The Figure 3.25 shows the usage of *Pnom*.

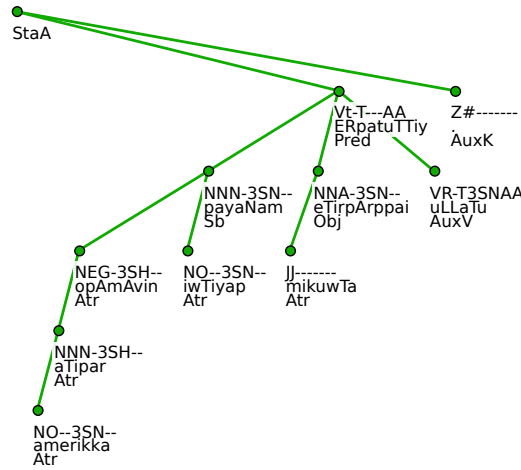
3.3.20 Afun: *Pred*

Predicate of the main clause will be given *Pred*. Only finite verbs in Tamil can be the predicate of the sentence. In Tamil, finite verbs at the end of the sentences are main predicates. So they receive *Pred* afun. In some cases, finite verbs will be absent at the end of the sentence. In that case, if there is a nominal predicate, it will receive *Pnom* afun or there won't be any *Pred* for that sentence. The Figure 3.26 shows an example annotation for *Pred* relation.



Tamil: ... வடக்குப் பகுதியில் மும்லைத் தீவுக்க உம் கிழக்குப் பகுதியில் திரிகோணமலைக்க உம் புதன்கிழமை செல்கிறார் .
 Tr: ... vatakkup pakuTiyil mullaiT Tlvukk um kizakkup pakuTiyil TirikONamalaikk **um** puTankizamai celkiRAR .
 Gloss: ... north region-LOC Mullai Tivu - east region-LOC Trincomalee - Wednesday going-he/she .
 English: ... He/She is going to Mullaitivu on the north **and** Trincomalee on the east .

Figure 3.23: *Coord*: Coordination head (Morphological marker *-um*)

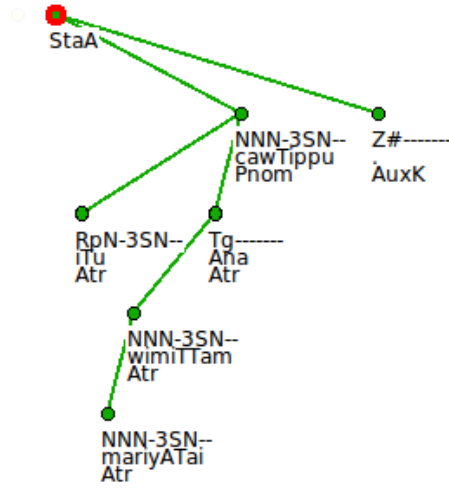


Tamil: அமெரிக்க அதிபர் ஓபாமாவின் இந்தியப் பயணம் மிகுந்த எதிர்பார்ப்பை ஏற்படுத்திய உள்ளது .
 Tr: amerikka aTipar opAmAvin iwTiyap payaNam mikuwTa **eTirpArppai** ERpatuTTiy uLLaTu .
 Gloss: American president Obama's India visit lot-of expectations create AUX-PERF .
 English: American president Obama's India visit has created lot of expectations.

Figure 3.24: *Obj*: Object

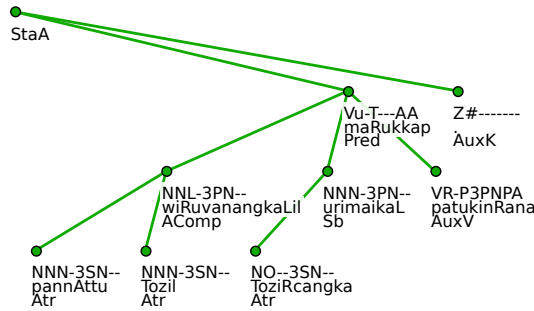
3.3.21 Afun: *Sb*

The label *Sb* is assigned to the subject of the sentence. If there are more than one subjects in the sentence (i.e in the case of multiple clauses), then the label *Sb* will be assigned to all of them. The subject of passive verbs will also be labeled with *Sb* relation. The Figure 4.21 shows the example usage of *Sb*.



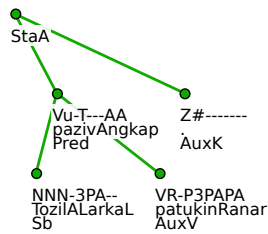
Tamil: இது மரியாதை நிமித்தம் ஆன சந்திப்பு .
 Tr: itU marryATai wimiTTam Ana cawTippu .
 Gloss: This courtesy - - call .
 English: This is a courtesy call .

Figure 3.25: Pnom: Pnom



Tamil: பன்னாட்டு தொழில் நிறுவனங்களில் தொழிற்சங்க உரிமைகள் மறுக்கப் படுகின்றன .
 Tr: pannAttu Tozil wiRuvanangkaLil ToziRcangka urimaikaL maRukkap patukinRana .
 Gloss: multinational business industries-LOC union rights deny-INF AUX-PASS .
 English: Union rights are denied in multinational business industries .

Figure 3.26: Pnom: Pred



Tamil: தொழிலாளர்கள் பழிவாங்கப் படுகின்றனர் .
 Tr: ToziALarkaL pazivAngkap patukinRana .
 Gloss: employees revenge-INF AUX-PASS .
 English: The employees are being revenged .

Figure 3.27: Sb: Subject

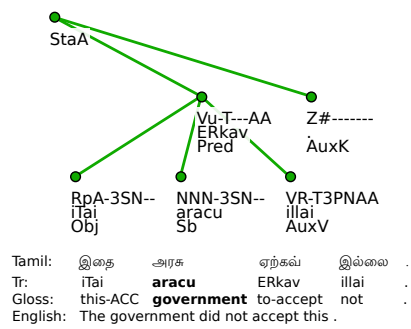


Figure 3.28: *Sb*: Subject

Acknowledgements

This project has been supported by,

- The European Commission's 7th Framework Program (FP7) under grant agreement n° 238405 (CLARA)¹ and



- The Grant MSM 0021620838 of the Czech Ministry of Education.



¹CLARA Homepage: <http://clara.uib.no/>

Bibliography

- [Begum et al., 2008] Begum, R., Husain, S., Dhvaj, A., Sharma, D. M., Bai, L., and Sangal, R. (2008). Dependency Annotation Scheme for Indian Languages. In *Proceedings of the Third International Joint Conference on Natural Language Processing, IJCNLP 2008*, pages 721–726. Asian Federation of Natural Language Processing (AFNLP).
- [Bharati et al., 2009] Bharati, A., Gupta, M., Yadav, V., Gali, K., and Sharma, D. M. (2009). Simple parser for Indian languages in a dependency framework. In *Proceedings of the Third Linguistic Annotation Workshop, ACL-IJCNLP '09*, pages 162–165, Stroudsburg, PA, USA. Association for Computational Linguistics.
- [Dhanalakshmi et al., 2010] Dhanalakshmi, V., Anand Kumar, M., Rekha, R., Soman, K., and Rajendran, S. (2010). Grammar Teaching Tools for Tamil Language. In *Technology for Education Conference (T4E 2010)*.
- [Janarthanam et al., 2007] Janarthanam, S., Nallasamy, U., Ramasamy, L., and Santhoshkumar, C. (2007). Robust Dependency Parser for Natural Language Dialog Systems in Tamil. In *Proceedings of the 5th Workshop on Knowledge and Reasoning in Practical Dialogue Systems, IJCAI KRPDS-2007*, pages 1–6.
- [Lehmann, 1989] Lehmann, T. (1989). *A Grammar of Modern Tamil*. Pondicherry Institute of Linguistics and Culture.
- [Nivre, 2009] Nivre, J. (2009). Parsing Indian Languages with MaltParser. In *Proceedings of the ICON09 NLP Tools Contest: Indian Language Dependency Parsing, ICON 2009*, pages 12–18.
- [Ramasamy and Žabokrtský, 2011] Ramasamy, L. and Žabokrtský, Z. (2011). Tamil Dependency Parsing: Results Using Rule Based and Corpus Based Approaches. In *Proceedings of the 12th international conference on Computational linguistics and intelligent text processing - Volume Part I, CICLing'11*, pages 82–95, Berlin, Heidelberg. Springer-Verlag.
- [Selvam et al., 2009] Selvam, M., Natarajan, A. M., and Thangarajan, R. (2009). Structural Parsing of Natural Language Text in Tamil Language Using Dependency Model. *Int. J. Comput. Proc. Oriental Lang.*, 22(2-3):237–256.
- [UFAL, 2006] UFAL (2006). The Prague Dependency Treebank 2.0.
- [Vempaty et al., 2010] Vempaty, C., Naidu, V., Husain, S., Kiran, R., Bai, L., Sharma, D. M., and Sangal, R. (2010). Issues in Analyzing Telugu Sentences towards Building a Telugu Treebank. In Gelbukh, A. F., editor, *CICLing*, volume 6008 of *Lecture Notes in Computer Science*, pages 50–59. Springer.
- [Žabokrtský et al., 2008] Žabokrtský, Z., Ptáček, J., and Pajas, P. (2008). TectoMT: highly modular MT system with tectogramatics used as transfer layer. In *Proceedings of the Third Workshop on Statistical Machine Translation, StatMT '08*, pages 167–170, Stroudsburg, PA, USA. Association for Computational Linguistics.

ÚFAL

ÚFAL (Ústav formální a aplikované lingvistiky; <http://ufal.mff.cuni.cz>) is the Institute of Formal and Applied linguistics, at the Faculty of Mathematics and Physics of Charles University, Prague, Czech Republic. The Institute was established in 1990 after the political changes as a continuation of the research work and teaching carried out by the former Laboratory of Algebraic Linguistics since the early 60s at the Faculty of Philosophy and later the Faculty of Mathematics and Physics. Together with the “sister” Institute of Theoretical and Computational Linguistics (Faculty of Arts) we aim at the development of teaching programs and research in the domain of theoretical and computational linguistics at the respective Faculties, collaborating closely with other departments such as the Institute of the Czech National Corpus at the Faculty of Philosophy and the Department of Computer Science at the Faculty of Mathematics and Physics.

CKL

As of 1 June 2000 the Center for Computational Linguistics (Centrum počítačnické lingvistiky; <http://ckl.mff.cuni.cz>) was established as one of the centers of excellence within the governmental program for support of research in the Czech Republic. The center is attached to the Faculty of Mathematics and Physics of Charles University in Prague.

TECHNICAL REPORTS

The ÚFAL/CKL technical report series has been established with the aim of disseminate topical results of research currently pursued by members, cooperators, or visitors of the Institute. The technical reports published in this Series are results of the research carried out in the research projects supported by the Grant Agency of the Czech Republic, GAČR 405/96/K214 (“Komplexní program”), GAČR 405/96/0198 (Treebank project), grant of the Ministry of Education of the Czech Republic VS 96151, and project of the Ministry of Education of the Czech Republic LN00A063 (Center for Computational Linguistics). Since November 1996, the following reports have been published.

- ÚFAL TR-1996-01** Eva Hajičová, *The Past and Present of Computational Linguistics at Charles University*
Jan Hajič and Barbora Hladká, *Probabilistic and Rule-Based Tagging of an Inflective Language – A Comparison*
- ÚFAL TR-1997-02** Vladislav Kuboň, Tomáš Holan and Martin Plátek, *A Grammar-Checker for Czech*
- ÚFAL TR-1997-03** Alla Bémová at al., *Anotace na analytické rovině, Návod pro anotátory (in Czech)*
- ÚFAL TR-1997-04** Jan Hajič and Barbora Hladká, *Tagging Inflective Languages: Prediction of Morphological Categories for a Rich, Structural Tagset*
- ÚFAL TR-1998-05** Geert-Jan M. Kruijff, *Basic Dependency-Based Logical Grammar*
- ÚFAL TR-1999-06** Vladislav Kuboň, *A Robust Parser for Czech*
- ÚFAL TR-1999-07** Eva Hajičová, Jarmila Panevová and Petr Sgall, *Manuál pro tektogramatické značkování (in Czech)*
- ÚFAL TR-2000-08** Tomáš Holan, Vladislav Kuboň, Karel Oliva, Martin Plátek, *On Complexity of Word Order*
- ÚFAL/CKL TR-2000-09** Eva Hajičová, Jarmila Panevová and Petr Sgall, *A Manual for Tectogrammatical Tagging of the Prague Dependency Treebank*
- ÚFAL/CKL TR-2001-10** Zdeněk Žabokrtský, *Automatic Functor Assignment in the Prague Dependency Treebank*
- ÚFAL/CKL TR-2001-11** Markéta Straňáková, *Homonymie předložkových skupin v češtině a možnost jejich automatického zpracování*
- ÚFAL/CKL TR-2001-12** Eva Hajičová, Jarmila Panevová and Petr Sgall, *Manuál pro tektogramatické značkování (III. verze)*

- ÚFAL/CKL TR-2002-13 Pavel Pecina and Martin Holub, *Sémanticky signifikantní kolokace*
- ÚFAL/CKL TR-2002-14 Jiří Hana, Hana Hanová, *Manual for Morphological Annotation*
- ÚFAL/CKL TR-2002-15 Markéta Lopatková, Zdeněk Žabokrtský, Karolína Skwarská and Vendula Benešová, *Tektogramaticky anotovaný valenční slovník českých sloves*
- ÚFAL/CKL TR-2002-16 Radu Gramatovici and Martin Plátek, *D-trivial Dependency Grammars with Global Word-Order Restrictions*
- ÚFAL/CKL TR-2003-17 Pavel Květoň, *Language for Grammatical Rules*
- ÚFAL/CKL TR-2003-18 Markéta Lopatková, Zdeněk Žabokrtský, Karolína Skwarska, Václava Benešová, *Valency Lexicon of Czech Verbs VALLEX 1.0*
- ÚFAL/CKL TR-2003-19 Lucie Kučová, Veronika Kolářová, Zdeněk Žabokrtský, Petr Pajas, Oliver Čulo, *Anotování koreference v Pražském závislostním korpusu*
- ÚFAL/CKL TR-2003-20 Kateřina Veselá, Jiří Havelka, *Anotování aktuálního členění věty v Pražském závislostním korpusu*
- ÚFAL/CKL TR-2004-21 Silvie Cinková, *Manuál pro tektogramatickou anotaci angličtiny*
- ÚFAL/CKL TR-2004-22 Daniel Zeman, *Neprojektivity v Pražském závislostním korpusu (PDT)*
- ÚFAL/CKL TR-2004-23 Jan Hajič a kol., *Anotace na analytické rovině, návod pro anotátory*
- ÚFAL/CKL TR-2004-24 Jan Hajič, Zdeňka Urešová, Alevtina Bémová, Marie Kaplanová, *Anotace na tektogramatické rovině (úroveň 3)*
- ÚFAL/CKL TR-2004-25 Jan Hajič, Zdeňka Urešová, Alevtina Bémová, Marie Kaplanová, *The Prague Dependency Treebank, Annotation on tectogrammatical level*
- ÚFAL/CKL TR-2004-26 Martin Holub, Jiří Diviš, Jan Pávek, Pavel Pecina, Jiří Semecký, *Topics of Texts. Annotation, Automatic Searching and Indexing*
- ÚFAL/CKL TR-2005-27 Jiří Hana, Daniel Zeman, *Manual for Morphological Annotation (Revision for PDT 2.0)*
- ÚFAL/CKL TR-2005-28 Marie Mikulová a kol., *Pražský závislostní korpus (The Prague Dependency Treebank) Anotace na tektogramatické rovině (úroveň 3)*
- ÚFAL/CKL TR-2005-29 Petr Pajas, Jan Štěpánek, *A Generic XML-Based Format for Structured Linguistic Annotation and Its application to the Prague Dependency Treebank 2.0*
- ÚFAL/CKL TR-2006-30 Marie Mikulová, Alevtina Bémová, Jan Hajič, Eva Hajičová, Jiří Havelka, Veronika Kolařová, Lucie Kučová, Markéta Lopatková, Petr Pajas, Jarmila Panevová, Magda Razímová, Petr Sgall, Jan Štěpánek, Zdeňka Urešová, Kateřina Veselá, Zdeněk Žabokrtský, *Annotation on the tectogrammatical level in the Prague Dependency Treebank (Annotation manual)*
- ÚFAL/CKL TR-2006-31 Marie Mikulová, Alevtina Bémová, Jan Hajič, Eva Hajičová, Jiří Havelka, Veronika Kolařová, Lucie Kučová, Markéta Lopatková, Petr Pajas, Jarmila Panevová, Petr Sgall, Magda Ševčíková, Jan Štěpánek, Zdeňka Urešová, Kateřina Veselá, Zdeněk Žabokrtský, *Anotace na tektogramatické rovině Pražského závislostního korpusu (Referenční příručka)*
- ÚFAL/CKL TR-2006-32 Marie Mikulová, Alevtina Bémová, Jan Hajič, Eva Hajičová, Jiří Havelka, Veronika Kolařová, Lucie Kučová, Markéta Lopatková, Petr Pajas, Jarmila Panevová, Petr Sgall, Magda Ševčíková, Jan Štěpánek, Zdeňka Urešová, Kateřina Veselá, Zdeněk Žabokrtský, *Annotation on the tectogrammatical level in the Prague Dependency Treebank (Reference book)*
- ÚFAL/CKL TR-2006-33 Jan Hajič, Marie Mikulová, Martina Otradovcová, Petr Pajas, Petr Podveský, Zdeňka Urešová, *Pražský závislostní korpus mluvené češtiny. Rekonstrukce standardizovaného textu z mluvené řeči*
- ÚFAL/CKL TR-2006-34 Markéta Lopatková, Zdeněk Žabokrtský, Václava Benešová (in cooperation with Karolína Skwarska, Klára Hrstková, Michaela Nová, Eduard Bejček, Miroslav Tichý) *Valency Lexicon of Czech Verbs. VALLEX 2.0*
- ÚFAL/CKL TR-2006-35 Silvie Cinková, Jan Hajič, Marie Mikulová, Lucie Mladová, Anja Nedolužko, Petr Pajas, Jarmila Panevová, Jiří Semecký, Jana Šindlerová, Josef Toman, Zdeňka Urešová, Zdeněk Žabokrtský, *Annotation of English on the tectogrammatical level*
- ÚFAL/CKL TR-2007-36 Magda Ševčíková, Zdeněk Žabokrtský, Oldřich Krůza, *Zpracování pojmenovaných entit v českých textech*
- ÚFAL/CKL TR-2008-37 Silvie Cinková, Marie Mikulová, *Spontaneous speech reconstruction for the syntactic and semantic analysis of the NAP corpus*

- ÚFAL/CKL TR-2008-38** Marie Mikulová, *Rekonstrukce standardizovaného textu z mluvené řeči v Pražském závislostním korpusu mluvené češtiny. Manuál pro anotátory*
- ÚFAL/CKL TR-2008-39** Zdeněk Žabokrtský, Ondřej Bojar, *TectoMT, Developer's Guide*
- ÚFAL/CKL TR-2008-40** Lucie Mladová, *Diskurzní vztahy v češtině a jejich zachycení v Pražském závislostním korpusu 2.0*
- ÚFAL/CKL TR-2009-41** Marie Mikulová, *Pokyny k překladu určené překladatelům, revizorům a korektorům textů z Wall Street Journal pro projekt PCEDT*
- ÚFAL/CKL TR-2011-42** Loganathan Ramasamy, Zdeněk Žabokrtský, *Tamil Dependency Treebank (TamilTB) - 0.1 Annotation Manual*