

*Addis Ababa
University*

(Since 1950)



ADDIS ABABA UNIVERSITY
SCHOOL OF GRADUATES STUDIES
DEPARTMENT OF INFORMATION SCIENCE

**DESIGNING A STEMMING ALGORITHM FOR SILT'E
LANGUAGE**

**A thesis submitted to School of Graduate Studies of Addis Ababa University
in partial fulfillment of the Requirement for the Degree of Master of Science
in Information Science**

By

MUZEYN KEDIR ABEDO

**June, 2012
Addis Ababa**

ADDIS ABABA UNIVERSITY
SCHOOL OF GRADUATES STUDIES
DEPARTMENT OF INFORMATION SCIENCE

**DESIGNING A STEMMING ALGORITHM FOR SILT'E
LANGUAGE**

**A thesis submitted to School of Graduate Studies of Addis Ababa University
in partial fulfillment of the Requirement for the Degree of Master of Science
in Information Science**

By

MUZEYN KEDIR ABEDO

Names and signature of Members of Examining Board

<u>Name</u>	<u>Title</u>	<u>Signature</u>	<u>Date</u>
_____	Chairperson	_____	_____
_____	Advisor	_____	_____
_____	Examiner	_____	_____

ACKNOWLEDGMENTS

Allahamdulilah, for everything Allah has made for me and again allahamdulilah.

Here, I would like to take this opportunity to express my most gratitude to those that had helped and inspired me in the completion of this thesis. I owe a great debt of gratitude to my advisor Marta Yifru (PhD). This thesis would not have been possible without her guidance, as well as inspiring and enlightening ideas, comments and suggestions. Thank you very much.

I would like to express my heartfelt thank to Million M. (PhD) for his comment and constructive ideas and I am proud of being his student.

I would like to express my gratitude to my beloved family for their encouragement and prayer (“Dua”).

To my friends, thanks for all the moral supports, especially Abubeker Usman, Enyew Ahmedine Nuredin Nasser, Kedir Mohamed, Dr. Elias Gebru and Mohammedamin Nassir for that they have made some important resources available and for their constructive comments and many other individuals that are unable to be listed who have directly or indirectly help me in completion of my thesis.

Last but not least, I would like to thank the University of Addis Abeba for the financial support without which the study would not have been possible.

Muzeyn Kedir Abedo

June, 2012

TABLE OF CONTENTS

CHAPTER ONE	1
1.1 BACKGROUND OF THE STUDY	1
1.2 STATEMENT OF THE PROBLEM	4
1.3 OBJECTIVES	6
1.3.1 GENERAL OBJECTIVE.....	6
1.3.2 SPECIFIC OBJECTIVES.....	6
1.4 METHODOLOGY AND METHODS.....	6
1.4.1 LITRATURE REVIEW	6
1.4.2 DATA SOURCE.....	7
1.4.3 DESIGNING AND DEVELOPING TOOLS AND TECHNIQUES.	7
1.4.4 TESTING PROCEDURE	7
1.5 SCOPE OF THE STUDY	7
1.6 SIGNIFICANCE OF THE STUDY	8
1.7 ORGANIZATION OF THE THESIS.....	8
CHAPTER TWO	10
REVIEW OF RELATED LITERATURE	10
2.1 CONFILATION TECHNIQUES.....	10
2.2 STEMMING ALGORITHMS	13
2.2.1 DICTIONARY-BASED TECHNIQUES	16
2.2.SUCCESSOR VARIETY	17
2.2.3AFFIX REMOVAL ALGORITHMS	18
2.2.4. STATISTICAL APPROACH	18
2.3 STEMMING ALGORITHMS FOR ETHIOPIAN LANGUAGES.....	19
2.3.1 TIGRIGNA STEMMERS.....	19
2.3.2 AMHARIC STEMMERS	20
2.3.3 OROMO STEMMERS	21
2.3.4 WOLAYTA STEMMER	22
2.4 STEMMING ALGORITHM FOR OTHER LANGUAGES	22
2.4.1 ENGLISH LANGUAGE STEMMERS.....	23
2.4.2 ARABIC STEMMING ALGORITHMS	26
2.5 EVALUATION METHODS FOR STEMMING ALGORITHMS	27

CHAPTER THREE	29
MORPHOLOGY OF SILT'E LANGUAGE.....	29
3.1 INTRODUCTION	29
3.2 THE WRITING SYSTEM OF SILT'E LANGUAGE	29
3.2.1 VOWELS AND CONSONANTS OF SILT'E LANGUAGE	30
3.3 MORPHOLOGY	31
3.3.1 SILT'E MORPHEME.....	32
3.3.2 WORD FORMATION IN SILT'E	32
3.4 INFLECTIONAL AND DERIVATIONAL MORPHOLOGY OF SILT'E.....	33
3.4.1 NOUN'S INFLECTION AND DERIVATION.....	33
3.4.2 ADJECTIVES'S INFLECTION AND DERIVATION.....	38
3.4.3 INFLECTIONAL AND DERIVATIONAL MORPHOLOGIES OF VERBS	40
3.4.3.1 ASPECT.....	41
3.4.3.1.1. VERB INFLECTION	43
3.4.3.1.2. PERFECTIVE.....	44
3.4.3.1.3. SIMPLE PAST.....	45
3.4.3.1.4. PRESENT PERFECT	47
3.4.3.1.5. PLUPERFECT	49
3.4.3.1.6. NEGATIVE	49
3.4.3.2 IMPERFECTIVE (NON-PAST).....	49
3.4.3.2.1. CONTINUOUS PAST.....	50
3.4.3.2.2. PRESENT /FUTURE.....	51
3.4.3.2.3. NEGATIVE	51
3.4.3.2.4. PRESENT/FUTURE (CONTINUOUS PAST) NEGATIVE	52
3.4.3.3 NON – ASPECTUAL.....	53
3.4.3.3.1. INFINITIVE.....	53
3.4.3.3.2. IMPERATIVE AND JUSSIVE	53
3.4.3.3.3. NEGATIVE	54
3.4.3.3.4. MOODS	55
3.4.3.3.5. MORPHOPHONOEMIC CHANGES WITH PREFIXES	55
3.4.3.4 DERIVATION OF VERBS	56

CHAPTER FOUR	62
A STEMMING ALGORITHM DESIGN AND IMPLEMENTATION (PROTO TYPE) FOR SILT'E LANGUAGE TEXT	62
4.1 INTRODUCTION	62
4.2 THE CORPUS	62
4.3 WORD DISTRIBUTION OF SILT'E	63
4.4. DOCUMENT (MORPHOLOGICAL) PREPROCESSING	64
4.4.1. TOKENIZATION.....	65
4.4.2. NORMALIZATION	65
4.5. COMPILATION OF THE STOP WORD LIST	65
4.6. COMPILATION OF AFFIXES.....	70
4.6.1. COMPILATION OF PREFIXES.....	70
4.6.2. COMPILATION OF SUFFIXES.....	71
4.7. THE PROPOSED STEMMER.....	72
4.8. RECODING AND CONTEXT SENSITIVE RULES.....	74
4.8.1. RECODING RULS.....	74
4.8.2. CONTEXTE SENSITIVE CONDITIONS.....	75
4.9. COMPONENTS OF THE STEMMER	76
4.9.1. PREFIX STRIPING.....	78
4.9.2 SUFFIX STRIPPING.....	81
4.10.3. LETTER REDUPLICATIO STRIPING.....	81
4.10.1. TYPE 1 REDUPLICATION.....	81
4.10.2. TYPE 2 REDUPLICATION.....	83
4.11. IMPLIMENTATION OF THE STEMMER.....	84
4.12. EXPERMANTS AND DISCUSSIONS.....	85
4.12.1. EVALUATION OF THE STEMMER.....	85
4.12.2. THE RESULTS.....	86
CHAPTER FIVE	89
CONCLUSION AND RECOMMENDATION.....	89
5.1 CONCLUSIONS.....	89
REFERENCES	92
APENDECES	99

LIST OF TABLES

Table 3.1: shows the seven alphabetical order of Ethiopic (Ge'ez) alphabet. -----	29
Table 3.2: the mapping of seven Ethiopic vowels to ten of Silt'e and order of Silt'e alphabet ---	30
Table 3.3: Possessive suffix for noun inflection -----	33
Table 3.4: Noun paucal formation with suffixation of -ቸረሻል-----	34
Table 3.5: paucal noun formation with reduplication and reduplication plus suffixation -----	35
Table 3.6: noun derivation from perfective stem by adding the suffix ኢሎ -ilol -----	36
Table 3.7: adjective paucal formation with suffixation of -ቸረሻል-----	37
Table 3.8: the derivation of adjectives from nouns-----	38
Table 3.9: examples of verbs from 2 to 4 radicals-----	39
Table 3.10: the nature of final radical of verbs-----	39
Table 3.11: the perfective, imperfective and non-aspectual forms of verb -----	40
Table 3.12: the pattern of consonant and vowel in class one verbs -----	41
Table 3.13: examples for class 1 verbs of Silt'e -----	41
Table 3.14: examples of class 2 verbs of Silt'e -----	41
Table 3.15: suffixes used to form the verb in present tense -----	42
Table 3.16: suffixes added to the perfective stem to mark the subject in the verb -----	43
Table 3.17 Simple past form of verbs -----	44
Table 3.18: Simple past form of palatal verbs -----	45
Table 3.19: Simple past form of Ca-verbs (constant followed by vowel a) -----	45
Table 3.20: present perfect form of verbs -----	46
Table 3.21: the pluperfect of verbs -----	47
Table 3.22: the prefixes and suffixes of imperfective verb-----	48
Table 3.23: the affixes of imperfective verb in continuous past tense -----	48
Table 3.24: the present/future form of Silt'e verb -----	49
Table 3.25: shows the negative prefix of imperfective verb -----	50
Table 3.26: shows the negative forms of present and future (past continuous tense) verbs -----	50
Table 3.27: the prefixes of imperative and jussive verb -----	51
Table 3.28: Shows the imperative and jussive forms of a verb -----	52
Table 3.29: the prefix for negative imperative/jussive verb formation -----	52

Table 3.30: the morphophonemic changes with prefixes-----	53
Table 3.31: morphophonemic changes for verbs which begin with ᵱ wa- and ᵱ • w- -----	54
Table 3.32: vowel initial verbs lengthened after they get prefixed-----	54
Table 3.33: derivation of causative verb -----	55
Table 3.34: derivation of transitive verbs by prifixing ᵱ la- -----	55
Table 3.35: derivation of passive verbs by prefixing ᵱ • ta- -----	55
Table 3.36: type 1 reduplication-----	56
Table 3.37: type two reduplication -----	57
Table 3.38: reduplication, transitive and causative derivation -----	58
Table 4.1: Number of words and their distributions to Silt'e sample text data -----	63
Table 4.2: comparison of word distribution ratio of the sample data set among different languages-----	63
Table 4.3: The most frequent words from the SampleDocument -----	67
Table 4.4: words with lower frequencies but are stop words-----	68
Table 4.5: sample list of stop words -----	69
Table 4.6: sample list of prefixes-----	70
Table 4.7: sample suffixes -----	71
Table 4.8: examples of recoding -----	74
Table 4.9: prefix striping stemming algorithm -----	78
Table4.10: Suffix Stripping Stemming Algorithm-----	79
Table 4.11: type1 redpuplication stripping stemmer algorithm -----	81
Table 4.12: type 2 reduplication stripping stemmer algorithm-----	82
Table 4.13: Sample of prefix striping-----	83
Table 4.14: Sample of suffix striping-----	83
Table 4.15: Sample of type1 the reduplication striping -----	83
Table 4.16: Sample of type2 the reduplication striping -----	84
Table 4.17: examples of stemming errors-----	84
Table 4.18: correct stem -----	85
Table 4.19: distribution of errors -----	85

LIST OF FIGURES

Figure 4.1: Flow chart for the general affix removal algorithm -----	68
--	----

LIST OF APPENDECES

APPENDIX I: List of stop words compiled for the stemmer -----	98
APPENDIX II: List of Silt'e prefixes compiled for the stemmer -----	100
APPENDIX III: List of Silt'e suffixes compiled for the stemmer -----	100
APPENDIX IV: Silt'e alphabet -----	101
APPENDIX V: Ge'ez (Ethiopic) alphabet -----	102
APPENDIX VI: Comparison of Silt'e words before stem and after stem-----	14
Appendix VII words with lower frequencies but are stop words -----	110

LIST OF ABBREVIATIONS AND SYMBOLS

Symbol	Meaning
CV	=Consonant and vowel sequence
‘ ‘	= meaning of words
 	= the English translation of words
1ps	= first person singular
1pp	=first person plural
2sm	=second person singular masculine
2sf	= second person singular feminine
2pm	=second person plural masculine
2pf	= second person plural feminine
3sm	=third person masculine singular
3sf	=third person feminine singular
3pm	=third person masculine plural
3pf	= third person feminine plural

Abstract

Variant word forms that are likely to be encountered in indexing and retrieval are one of the causes of the problems that are involved in the use of freetext retrieval system. The variant word structures used in indexing and searching are to be expected in determining the relevance of a document to a user query that specifies just a single form. Shrinking the variant words into one form advances the performance of IR system and this can be achieved by conflation techniques, which is usually stemming that is established in this work. Stemmers are used in information retrieval to reduce as many related words and word forms as possible to a standard form, which can then be used in the retrieval process. This research explores the possibility of developing a stemmer to conflate variant words of Silt'e language.

Silt'e belongs to the Semitic language group. These languages have a common grammatical system based on a root-pattern structure. Consonants bear the basic meaning while vowels form different patterns. Stems are built from consonantal roots before other word forms are built. Silt'e uses affixation and reduplication to derive different word forms from stems. Common affixations are prefix, suffix, and infix. Silt'e uses extensive concatenation of affixes and can result in relatively long words, which often contain an amount of semantic information equivalent to a whole English phrase, clause or sentence. As a result of this complex morphological structure, a single Silt'e word can have very large variants.

To design the stemmer, a sample text was collected from different sources and research paper that explains the morphology of Silt'e language also used and affixes and stopwords collected from this research paper and the sample text document to develop the stemmer. The stemmer, developed in this study is iterative and uses context sensitive and recoding rules that remove prefix, suffix and reduplication of letters (type 1 and type 2). In this experiment the stripping procedure were applied in order: prefix, suffix and finally letter reduplication. The stemmer was tested on a sample data of 1486 words, which were selected randomly from the sample texts. The result of the experiment shows that, the stemmer performs at accuracy of **85.71%**, and brings a dictionary reduction of 34.99% for stem words. Lastly conclusion and the possible recommendation for future work were reported.

CHAPTER ONE

1.1 BACKGROUND OF THE STUDY

The activity of archiving written information can be traced back to around 3000 BC, when the Sumerians selected special areas to store clay tablets with cuneiform inscriptions. Even then the Sumerians realized that proper organization and access to the archives was critical for efficient use of information. They developed special classifications to identify every tablet and its content.

The need to organize, store and retrieve written information happened to be increasingly important over centuries, particularly with the emergence of paper and printing press. Then computers came into existence and people understood that it could be used for storing and mechanically retrieving large amounts of information [32].

The collection development problem is aggravated by the growth in available information. In early times, the total available information (knowledge) changed slowly. However, by the year 1800 the amount of scientific publication was doubling every 50 years [33]. More recently with the remarkable growth of science and technology, the rate of increase of available information (knowledge) has greatly accelerated. "Between 1800 and 1966, the number of scientific journals has increased from 100 to 100,000". At the present time, no upper limit is apparent in the rate of increase of available information items.

The dramatic growth in size and variety of available information made the utilization and access of the right information difficult. Therefore, proper organization of information resources matters to get the most out of it and effective and efficient utilization of it.

Information retrieval systems (IRs) [45] are designed to facilitate access to stored information. In addition, they are concerned with the representation, storage and organization of information items. The components of an IR system consist of a set of information items, a set of requests and a mechanism to determine which information items are most likely to meet the requirements of the requests. Basically, users of information state their information needs in the form of queries and submit their queries to the information retrieval system. Based on their queries, the system outputs information items that are believed to be relevant to the user query. In other words, the items that have common entries in both the database

and users query are returned to users. This happens when a matching exists between queries and information items. To attain a match between these two components, the items in the documents as well as in the query should be represented in some form.

Frequently, the user specifies a word in a query but variant of this word are present in a relevant document. The variation prevents a perfect match between a query word and word in relevant documents. This problem can be overcome with the substitution of the words by their respective *stems*. A stem is a portion of a word that is left after the removal of its affixes (i.e., prefixes and suffixes). A typical example of a stem can be the word “comput” which is the stem for the variants “computed,” “computing,” “computation,” and “computations.”

Morphology describes how various forms of words are created, and studies structure of words in the language. Suffixing and prefixing are the main and common ways of creating word variations in natural language text [36]. Morphology (the internal structure of words) can be broken down into two sub classes: inflectional and derivational [6]. Inflectional morphology describes predictable changes of words. These are as a result of syntax, such as changes in person, number, tense and gender. These changes have no effect on a word’s part of speech [38] that means; a noun still remains noun after pluralized. For example, kick, kicks, kicking, kicked.

On the other hand, changes of derivational types may affect the word’s meaning in part of speech. For example, affix changes from adjective to nouns, from verb to nouns, from noun to verbs, and so on; like friend, friendly, friendliness and friendship.

From either type of morphologies, depending up on the complexity of the morphological natures of language types, very massive variants of words may be resulted from a single word. This large number of occurrences has strong impact on information retrieval systems. Thus, there is a need of automated procedure that can reduce the size of the various words to manageable level, and also capture the strong correlation existing between different word forms [39]. Even if the complexity of morphological variations may differ from one natural language to others, stemming is widely used in information retrieval (IR), with the basic reason that morphological variants represent similar meaning.

Morphological processing becomes widely used for effective and efficient information retrieval [37], machine translation (MT) and text classification [39]. Hence, morphological processing becomes particularly important for information retrieval (IR) because IR needs to determine the appropriate forms of words as index [24].

Information retrieval particularly automatic information retrieval system is an information processing activity which is carried out with the help of automatic equipment. As [1] defined automatic information retrieval system is a software /hardware package that lets different users to access query and retrieve information from the data base.

Stemming is a natural language processing technique that identifies a stem/root from morphologically conflated words using various techniques on inflectional and derivational affixes [36]. It helps to reduce morphological variants of words in to a common form. As an example, ‘throws’, ‘thrower’, and ‘throwing’ are the common variances of a stemmed word ‘throw. Therefore, stemming can be applied in different languages to increase searching efficiency and reducing vocabulary size of indexed files in information retrievals [38, 6].

In stemming, all of the word’s true prefixes and suffixes are removed to produce the stem or the root word [40]. This is important to classify texts and index. It is also used to make operations fewer dependents on particular form of words rather than enclosing all possible forms. The computational process gathers all words that share the same stem and have some semantic relations. It is also used for document matching and classification to convert all likely forms of a word in the input document to the form in reference document [40].

The types of automatic stemming strategies are four: affix removal, table look up, successor variety, and N-gram [36]. Affix removal stemming is imitative, simple and can be implemented efficiently.

In this type of strategy suffixes and/or prefixes will be removed from the word to get the stem. This type of strategy was used by [39, 36, and 40]. Table look up strategy looks for the stem of a word in a table of the dictionary. This strategy was used by [6] to stem Amharic morphological words. It is a simple procedure and depends on the size of stemmed dictionary before the activities will be performed. It also requires considerable storage space. Successor variety stemming is done based on the determination of morpheme boundaries, uses knowledge from structural linguistics, and is more complex than affix removal. N-gram stemming is used based

on the identification of bi-grams and tri-grams and is more of a term used for clustering than stemming. It is done based on numbers of n-grams. This strategy was used by [41] to stem language independent terms using uni-gram.

Unlike languages like English that has less morphological variants, Semitic languages like Amharic, Arabic, Ge'ez are much rich in morphology. So, like Amharic [38], Arabic [40], Ge'ez [44], Silt'e (which is one of the Semitic languages in Ethiopia) involves dealing with prefixes, infixes and derivatives in addition to suffixes.

1.2 Statement of the problem

According to [46], the EthioSemitic and South-Arabia languages are categorized under south Semitic. The EthioSemitic languages are further classified as North- EthioSemitic and South-EthioSemitic. Tigrigna together with Ge'ez and Tigre form the North- EthioSemitic language and the South- EthioSemitic includes Amharic, Argoba, Hariri, Silt'e, etc.

Silt'e is one of the zonal languages in the southern Nations, Nationalities and peoples region of Ethiopia. According to the 1998 census, the language is spoken in Central Ethiopia by about more than 827,764 people. Besides the Silt'e Zone, Silt'e is also spoken in Ethiopia's major cities, particularly in Addis Ababa. Elementary schools in Silt'e zone give all subjects in Silt'e from grade one up to four. After grade four, it is given as one subject for grades up to ten and also given for grade ten national student exam in Silt'e zone. The language uses Saba Ge'ez (Ethiopic script) [35]. Based on the researcher's preliminary analysis there are a significant amount of electronic documents in Silt'e in schools, government organizations of the zone and on net and there is an increase from time to time. However, there is no tool or system that supports the organizations to access valuable information for the purpose of decision making and other problem solving. Most of the time organizations are requested for information by other organizations and the request may need immediate response. The absence of tools or systems that help them in accessing and searching information delays their reply; as a result, they may fail to achieve their organizational goals.

Retrieval systems enable us to get information easily and timely. In today's shrinking world, it is becoming evident that there is a large body of information and research available only in the

language of the primary researcher, and much information can not be shared between research communities with out considerable translation and time delay [47].

Including Silt'e language, there are more than 86 languages in Ethiopia. But there are few research works done to design stemmer or an information retrieval system for local languages, such as Amharic [42], Tigrigna [40], Afanoromo [43], Wolayta [13], and Ge'ez [44] to simplify the problem of morphological complexity of the respective languages. Since most of the stemmers developed for other languages are language specific, it is not possible to make use of a stemmer developed for other languages. This is because of structural difference of the languages. So it seems quite difficult to follow the same stemming pattern and apply the same technique for all languages. Therefore, the stemmer of one language can not be effectively applied to other languages [23].

In developing an IR system for a particular language, the study of automatic indexing technique is important for the language. Among the process of automatic indexing, stemming and stop words are highly dependent on the document language. The motivation for using stemming is the need to increase effectiveness of retrieval system since stem of a term represents a broader notion than the original term itself [17]. To my knowledge so far no IR computational tools like stemming have been developed for Silt'e language. Therefore, it is felt necessary to undertake a research on the development of stemming algorithm for the language. Hence, this study.

It is, therefore, the aim of this study to explore the possibility of designing stemmer algorithm for Silt'e language. So this research attempts to answer the following questions:-

- What are the properties and word formations in Silt'e language?
- Is it possible to adopt some techniques from stemming algorithms developed for other languages?
- Does the designed prototype system registers good performance based on the sample text documents?

1.3 OBJECTIVES

The present research has the following general and specific objectives.

1.3.1 GENERAL OBJECTIVE

The main objective of this study is to design a stemmer for Silt'e language.

1.3.2 Specific objectives

To achieve the above general objective the following specific objectives will be attempted in the research work.

- To have a clear understanding of the area through conducting relevant literature review and identify different stemming algorithms that have been developed for other languages.
- To review properties of the Silt'e language in order to get familiar with the different aspects of the language and knows how the separation of words is achieved.
- To build a stop word list and compile a list of affixes used for the corpuse.
- To designe appropriate stemming algorithm for Silt'e language.
- To develop a stemmer that identifies inflectional and derivational affixes of Silt'e language.
- To test the performance of the stemmer on the selected Silt'e text and report the finding of the study.
- Forward recommendations for further study.

1.4 METHODOLOGY AND METHODS

1.4.1 LITRATURE REVIEW

Document analysis used to understand the characteristics of the language. As studying the language's morphology constitute an important component in the research, a literature survey was made to gather information and to understand the language. Discussion with appropriate individuals (linguists) was not possible because of different resons.

1.4.2 DATA SOURCE

A text corpus is one of the resources required in stemming algorithm or IR researches. The text was used for compiling stop words, prefixes and suffixes. Moreover, the text corpus has been used for testing the algorithm. A good sized text can show a reasonable language morphological behavior. Selection of text is, therefore, an important component in developing a stemmer. For the purpose of this research texts from school books and other sources were used.

1.4.3 DESIGNING AND DEVELOPING TOOLS AND TECHNIQUES.

Python (3.1) programming language was used to develop the stemmer algorithm for Silt'e. This is because it is relatively easier to manipulate text (string manipulation) and also the researcher has some experience in programs using python programming language. Text documents (texts written in Silt'e, grammar research book of Silt'e language and school books) have been used to study the language's characteristics. For the experiment context sensitive affix removal technique was used. This is because affix removal technique is relatively easy and, efficient to implement than other techniques.

1.4.4 TESTING PROCEDURE

To evaluate the stemmer, manual error counting and vocabulary reduction method were used. Manual error counting method found to be appropriate and was used to test the stemmer. Dictionary reduction method also was used to test the comparison power of the stemmer.

1.5 SCOPE OF THE STUDY

The main aim of this study is designing a stemmer for Silt'e language. This research mainly focused on derivational and inflectional word variants of the language. Compound and irregular words and infix stripping technique were not addressed here because of difficulty to handle and time limitation. The stemmer contains; prefix stripping, suffix stripping and letter reduplication stripping techniques (type 1 and type 2). Prefix-suffix stripping is not included in the stemmer because it can be handled by incorporating the prefix part to prefix list and the suffix to suffix list.

1.6 Significance of the study

One of the tools for IR is stemming. By developing a stemming algorithm for Silt'e variant words can be conflated. In Silt'e a word can have very large variants and conflating these variants increases retrieval performance. It also reduces storage of index files. In addition to the above results, it can also give the following benefits:-

- Can be used as one input to integrate and develop a general Semitic language information retrieval system and to form multi-lingual information retrieval.
- This research can also help to develop tools like spell checker, grammar checker, thesauri, word frequency counter, document summarizers and indexers. It can also be used to reduce word variants and to minimize total number of files.
- It can be used for developing cross lingual retrieval systems. For example one can feed a Silt'e text query to search engine, this entered word may be translated to other languages like English and access documents in English or other languages

1.7 ORGANIZATION OF THE THESIS

The research work consists of five chapters. This chapter introduces the importance of stemming on IR environment and the need to develop a stemming algorithm for conflating variants of a word in Silt'e language. Statement of the problem and its justification, methodology employed, and scope of the study are also discussed in this chapter.

Chapter two reviews the works on conflation techniques in general and stemming algorithms in particular. Detailed discussions are made on approaches to stemming and types of stemmers. Review is also made in this chapter on some stemming algorithms developed for other foreign and local languages.

Silt'e language morphology is reviewed in chapter three. The inflectional and derivational morphologies of the language are the main concerns of this chapter. Word formation processes for Silt'e nouns, adjectives, and verbs will also be presented in detail in the chapter.

Discussions on the development and, experimentation of the stemming algorithm for Silt'e text are given in the fourth chapter. The compilation of stopword lists and affix (prefix-suffix pair, prefix, and suffix) list is presented in this chapter. The approach employed to develop the stemmer and the reasons for its selection a lso parts of the discussions in this chapter.

The last chapter, chapter five, presents conclusions deduced from the findings and recommendations for future research.

CHAPTER TWO

REVIEW OF RELATED LITERATURE

2.1 CONFILATION TECHNIQUES

Queries and documents related with each other by the terms they have in common .The more they have similar terms the more they will be related. Therefore the degree of relationship between them is determined mainly by the frequency of terms they have in common [21]. However some words may have different forms that require some form of processing to change it to its root word. The existence of these variants for a word resulted from its different morphology. Hence morphological variance of a word is the main problem in indexing and retrieving systems [3].

Words can have different forms (variants) in a sentence for the sake of making agreement between subjects and verbs in gender, number, tense, person, mood or voice. Moreover, valid alternative spellings (e.g. “realize” and ”realise”), antonyms (“agreement” and “disagreement”) and different spelling use for the same word in different areas(often and color (America)and offen and coulour (Britain)). Besides abbreviations are also causes to word variants [22]. Generally, morphological variance is the main reason for text word variants [38]. Thus, morphology uses to analyze and describe the forms of words in a language, and usually includes inflectional and derivational.

Derivational morphemes result in new word that have different meaning from the previous one, like “creation” which is formed from “create” by adding suffix “-ion”. On the contrary inflectional morphemes don’t result in meaning difference between variants. An inflectional morpheme varies the form of words in order to express grammatical features, like singular/plural or past/present tense. For instance, student and students are two different forms of same word [42].

Compounding is the process of constructing morphologically complex words from two or more unbound morphemes. That is, compounding is the joining of two linguistic forms, which are functionally independent [43].

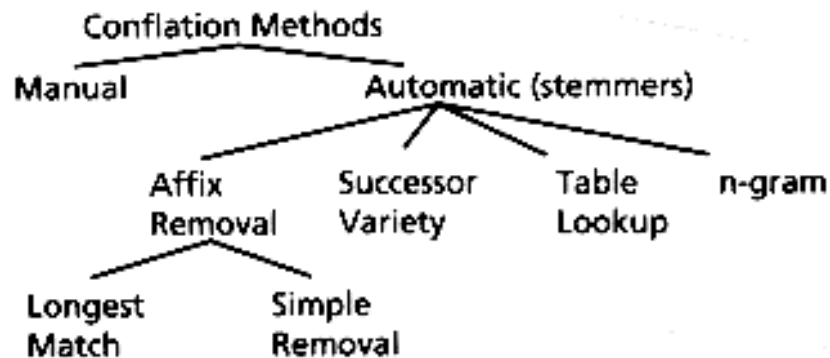
Morphological variants of a given word taken as analogous and, equivalent terms for the sake of information retrieval application. Hence, we have to take in to consideration the morphological variants of words as one of reference of all word forms. Converting these complex forms of words in to their root form is very important in determining the relevance of documents to user's queries in information retrieval. This concept helps to perform stemming in different languages using a number of stemming or conflation methods to get root/stem of morphological variants [21].

Different forms of a word are not recognized equivalently unless there is some kind of natural language processing (NLP). Because word form variations has significant impacts on information retrieval systems effectiveness [21]. Therefore, devising a mechanism to reduce a variation of word to its root form is important. To alleviate the word morphological problems, we use a method that helps to bring together semantically related words to a single form, which is a conflation technique [19].

Conflation is the process of identifying morphological variants of words that have similar concept and represent them by their root words [22]. It can also be defined as the process of matching morphological term variants by mapping term variants to a single form, usually a unique well-formed root for each word root for each word [24]. It performs by bringing similar words that have similar meaning with in common forms. Conflation techniques used to solve problems that are raised due to variations in word forms.

Conflation helps to increase efficiency and effectiveness of the information retrieval systems [40]. It also uses to understand the degree of relationship between different word forms, to overcome vocabulary mismatch problems in addition to better generalizations [20].

Conflation techniques discussed in [22, 19]: the two types of conflation techniques are manual and automatic. Manual conflation is performed at search time with right hand truncation. It is performed on query words but not on documents.



The other technique i.e. Automatic compilation is designed using stemmers. It alleviates the problem of word variants by matching different forms of the same words to their root automatically. Automatic conflation performed both on query and document words. This technique has two classes [23, 19]: stemming algorithm and string similarity algorithms. Stemming algorithms or stemmers are techniques for reducing words to their grammatical root stem form. These types of algorithms are usually designed to handle morphological variants with in language. In contrast, string similarity algorithm is usually language independent and designed to handle all types of variants. This approach involves conflating by calculating and measuring the similarity between an input query and terms and each of the distinct terms in the data base. Terms that have a high similarity to a query terms are then displayed to the user for possible inclusion in the query. The common example to string similarity algorithm is N-gram matching techniques.

In both manual and automatic conflations, over stemming and under stemming problems may occur. In over stemming, the stem is too short (too much of affix is removed from a term) and will result in unrelated words being stemmed to the same stem; while under stemming arises if too short affixes are removed and may result in related terms being described by different stems. Example; the stem of a word ‘instructors’ may be considered as ‘instructor’, rather than ‘instruct’ . These errors are occurred when a word is not stemmed properly. These causes irrelevant documents are retrieved or only few relevant documents are retrieved. The errors that may occur in stemming could be corrected manually using different exceptional lists, or by means of rules in the stemmer [20].

The effectiveness of the retrieval system is usually measured/evaluated using a combination of both recall and precision. Precision is a proportion of the retrieved documents that is usually relevant to the user interests; on the other hand, recall is a ratio of the relevant documents that is actually retrieved from the file. Achieving high recall and retaining a high precision use to improve the salient of information retrieval systems [33]. Hence, to improve the recall level, there are various types of methods such as use of term phrases, use of truncated terms or stems, and use of the synonyms [24].

Another issue that helps to illustrate retrieval effectiveness is stop word lists. Stop words are removed because of the following impacts on information retrieval [20]: they can affect the retrieval effectiveness and weighting process. This is because they mostly have a very high frequency and tend to diminish the impact of frequency differences among less common words. They can also affect efficiency due to their nature and the fact that they carry little meaning, which may result in a large amount of unproductive processing. Therefore, removing stop words results in increase effectiveness of retrieval by reducing storage requirements and increasing the matching of a query with index terms of a document.

2.2 STEMMING ALGORITHMS

Words in the query are matched with the words in the text database to retrieve relevant documents [1]. A query can be Boolean, where Boolean operators used to connect words in the query and documents must satisfy the given condition to be retrieved; or ranked, where the documents with the highest statistical similarity to the query are retrieved as answers. There must be a mechanism for deciding whether a given query term corresponds to a given word in a document for both types of queries. The simplest is to allow exact matching only; for example, "bicycles" would match itself only and a document that contained "bicycles" but "Bicycles" would not be recognized as a match. Another simple correspondence method is case-insensitive exact match.

Stemming is a more general form of correspondence, in which query terms and document terms are reduced to a common root word by removing affixes. Terms match if they have a common root. For example, "trains" might be stemmed to "train", which would match "train", "training", and of course "trains". Each word thus has a set of connotations, that is, words that have the same

root. Several stemming algorithms have been proposed [2, 3, 4, 5, 6, 7, 8], based on different principles; each produces rather different sets of confluations.

Stemming has usually been measured (evaluated) by its impact on querying: since stemming changes the documents that are retrieved in response to a query, it has the potential to change the quality of the set of answers. Despite the fact that stemming and choice of stemmer can significantly change the number of documents that match a query term, the results of information retrieval experiments have been mixed, ranging from negligible difference in average performance to unambiguous but small improvements [9, 10, and 3]. However, measuring average performance can mask significant changes in individual queries, and the metrics used for assessing performance can allow the effects of stemming to seem small, even if the documents sets retrieved are changed significantly.

“Instead of using variation in precision and recall as an indirect measure of stemming accuracy, it is possible to evaluate stemming algorithms by the accuracy of the confluations they produce. This was the approach used by Paice [11].

Although the character of stemming algorithms may vary significantly depending on whether a stem dictionary is being used, whether a suffix list is being used, and of course on the purpose for which the stemmer is designed, most of them are based on certain principles and procedures [5,12]. These procedures involve either the removal of the single longest matching suffix or the iterative removal of several suffixes. The rationale behind iterative approach is the fact that suffixes are attached to stems one after the other. In the iterative approach, suffixes are removed from the stem in the order of their derivational rules. In this approach the stripping starts from the end of the word and working towards the beginning. For instance, if the word FRUTEFULLNESS is considered, the suffix –NESS will be removed in the first iteration and re-considering FRUTEFULL, the suffix –FULL will be removed leaving FRUTE as a final stem. Stripping is done based on the class order that is defined by the programmer. Other approaches such as proposed by [11] may also exist, which are iterative but do not have their endings classified. Using the iterative approach has also its own problem [13] the need of examining large number of endings, which is a difficult activity and compilation of list of ordered class.

Iterative approach stripping procedure is used by Porter. It operates in five stages, using five different classes of suffixes (a total of about 60) to simulate the inflectional and derivational process of words. Some of the rule does not actually delete the suffix but they transform endings of the stem to new endings. For instance, if we remove the suffixes “-ing” from “absorbing” and “- ion” from “absorption”, it gives us “absorb” and “absorpt” which are different stems. But using the recoding rule that transforms endings from “-rpt” to “-rb,” the two words are conflated to the same stem that is “absorb” [6].

In the case of longest mach approach, it removes the longest suffix possible. If the same word FRUTEFULLNESS is considered the suffixes in the word are: -NESS, -FUL, and FULLNESS therefore the algorithm removes -FULNESS from the word. The problem of using longest match approach compared to iterative method is: need for generating all possible combinations of affixes and processing and storage space required and the change of affix during concatenation.

There are context-free or context-sensitive stemming algorithms. Context refers to any property that is attached to the remaining stem and usage of the suffix [14]. In context-free algorithms, no restriction is applied on the stem and hence no additional operations are required to check restrictions. But most of stemming works [5, 6,15, 16] indicated that better result could be obtained by adding constraints to stripping operations, that is, using context-sensitive, which are primarily language dependent. Context sensitive rules specify particular conditions in which each suffix may be stripped from an input word. Savoy [15] illustrates three general types of constraints: quantitative, qualitative and recording rule.

“In quantitative constraints minimum length for the remaining stem is set when a suffix is removed. This helps not to remove ending from a stem in which the ending is in the suffix list but actually not a suffix for that stem. For example for the word ABILITY and suffix ABILITY, as the remaining stem must not be zero, the suffix -ABILITY will not be stripped from the word.”

Qualitative constraints define conditions to be satisfied by ending of the remaining stem. For example in English, a stem does not end with that is, stripping “ize” from “seize” is not allowed.

Spelling corrections and adjustment rules must be used to conflate the words to exact stem. A recording rule is a transformation of the form $AxC \rightarrow AyC$, where A and C specify the context transformation, X is input word (string), and y is the transformation string.

There are stemming ranging from a weak stemmer, which eliminate only plural markers to a intricate one that removes suffixes and prefixes. From the available literature, stemming algorithms for local languages and abroad discussed well.

Most stemming algorithms follow the principles and approaches employed by the two most common English stemmers, Lovis and Porter. These two algorithms that form the bases of most algorithms employ suffix stripping.

Lovins [5] algorithm is based on longest match principle, and uses a list of 260 endings sorted in decreasing order of length. For instance, a word COMPUTATIONALITY would be stemmed to COMPUT if the suffix –ATIONALITY were included in the list. After a suffix is removed, the stem is compared with one of the rules (34 recoding rules) to consider spelling exceptions. The rules define, for instance, the treatment of a suffix preceded by double consonants (such as JUMPPING to JUMP), or minimal stem size must be retained (such as the removal of ING from PLAYING but not from SING) etc. the problem of using longest match approach is the need to have a largest affix list: basic and possible combinations. This requires storage costs and compiling and processing the list.

The problem of stemming has been approached with a wide variety of different methods. According to the Frakes' classification (1992), all stemming algorithms can be roughly classified as Dictionary-Based Techniques, successor variety, Affix removal and Statistical (N-gram).

2.2.1 DICTIONARY-BASED TECHNIQUES

A dictionary technique depends mainly on creating a very large dictionary which stores words found in natural texts with their corresponding morphological parts. Such parts include: stems, roots, and affixation. Each word uses a unique entry in a look up table. For example terms such as engineering process is performed manually, wherein the stem stems are defined for each word and stored in some kind of structured form [60].

This techniques list all words that exist in a specific language which their corresponding variation of words created by different attachments to the root word. To stem a word, the table is queried to find a matching variation of a word. If a matching variation of a word is found, the association root form is returned. Terms from queries and indexes could then be stemmed via table look up. Table entries are listed in alphabetical order. A hash table or binary search list can be used to optimize the search. Even though such approaches are accurate, the problems associated with them include the dictionary maintenance, to keep up with an ever-changing language, and this is actually quite a problem. It is not only that a dictionary created to assist stemming nowadays will probably require major updating in a few years time, but also that a dictionary in use for this purpose today may already be several years out of date as well as storage overhead and retrieval time needed for such data.

The advantage of this approach is that it generates perfect stems. However, the approach is limited to retrieving only those words that have previously being stored. The space occupied for storage tends to grow as the corpus expands, which can make the search process inefficient. In Korvetz's dictionary experiments [74], this direct method, seems effective but in adequate to deal with the "unlimited" words and their formation, especially in inflected languages with elevated morphological structure. This was the main reason that led him to evaluate algorithmic stemmers and conclude that "despite the errors they can be seen to make, they still give good practical results". It is not only that a dictionary created to assist stemming nowadays will probably require major updating in a few years time, but also that a dictionary in use for this purpose today may already be several years out of date".

2.2.2. SUCCESSOR VARIETY

The successor variety approach has been introduced by Hafer and Weiss. It is corpus-based and observes the number of distinct letters following a particular prefix: the successor variety. Successor variety stemmers [71] are based on work in structural linguistics which attempted to determine word and morpheme boundaries based on the distribution of phonemes in a large body of words. The stemming method based on this work uses letters in place of phonemes, and a body of text in place of phonemically transcribed words. It scans the word to be stemmed and finds the cut point where the successor variety increases sharply. Several variations are possible, including: cut-off, peak and plateau, complete word and entropy.

2.2.3. AFFIX REMOVAL ALGORITHMS

This approach is dependent on the morphology of the target language. Affix removal algorithms remove suffixes and/or prefixes from terms leaving a stem. These algorithms sometimes transform the resultant stem. There are many varieties of affix removal algorithms. Most stemmers currently in use are iterative longest match stemmers, a kind of affix removal stemmer first developed by [77]. An iterative longest match stemmer removes the longest possible string of characters from a word according to a set of rules. This process is repeated until no more characters can be removed. Even after all characters have been removed, stems may not be correctly conflated. In addition to [77], iterative longest match stemmers have been reported by Dawson [1974], Porter [81], and Paice & Husk [79].

Affix removal algorithms are often called Rule-Based Technique and is a widely applied stemming technique and the algorithm is introduced by Porter [81]. With specific rules for the English language, this algorithm iteratively removes suffixes from a given word, reducing it to its stem. Even if the algorithm has its limitations, it is the most commonly accepted for its high precision and recall. Lovin's stemmer [77] follows the same rule-based technique but it does not apply its rules iteratively and it is more conservative than Porter's algorithm.

2.2.4. STATISTICAL APPROACH

The stemming process involves statistical methods whereby, through a process of inference and based on a corpus, rules are formulated regarding word formation. Some of the methodologies adopted are N-gram [78] and Hidden Markov Models (HMM) [84]. Most stemmers are language-specific. Single N-gram stemming suggested by [78] tries to bypass this limitation. The idea is to analyze distribution of all N-grams in a document. Since the morphological invariants (unique word roots) will occur less frequently than variant parts (common prefixes and suffixes), a typical statistics like inverse document frequency (IDF) is used to identify them. The authors successfully tested the algorithm on some European languages. The N-gram stemming underperformed stems and required significant amount of memory and storage for index, but its ability to work with an arbitrary language makes it useful for many applications.

Another statistical approach to stemmers design used by [85] is to build a stemming algorithm based on Hidden Markov Models (HMM). It doesn't need a prior linguistic knowledge or a

manually created training set. Instead it uses unsupervised training which can be performed at indexing time HMMs are finite-state automata with transitions defined by probability functions. Each character comprising a word is considered as a state. The authors divided all possible states into two groups (roots and suffixes) and two categories: initial (which can be roots only) and final (roots or suffixes). Transitions between states define word building process. For any given word, the most probable path from initial to final states will produce the split point (a transition from roots to suffixes). That the sequence of characters before this point can as a stem. Using Porter's algorithm as a baseline, they found that HMM had a tendency over stem the words.

In their research, Jinxi Xu and W. Bruce Croft suggested an approach, which allows correcting "rude" stemming results based on the statistical properties of a corpus used [83]. The basic idea is to generate equivalence classes for words with a classical stemmer and then "separate back" some conflated words based on their co-occurrence in the corpora. This approach does not require any linguistic knowledge whatsoever, being totally independent of the morphological structure of the target language.

2.3 STEMMING ALGORITHMS FOR ETHIOPIAN LANGUAGES

Unlike English and other western languages, Ethiopian languages are less researched languages in the areas of information retrieval and natural languages processing applications. Recently there are some researches done in the areas of IR and NLP for Ethiopian languages. Some of the attempts done are presented as follows.

2.3.1 TIGRIGNA STEMMERS

The first attempt to develop Tigrigna stemmer was made by Girma Berhe [70]. The algorithm uses iterative approach but when it finds two affixes that match with the word, it removes the longest one. As Tigrigna is morphologically complex language a context-sensitive stemmer was considered appropriate in this algorithm. Each affix is accompanied by a context-sensitive rule. The techniques for describing the rules are adopted from Porter (1980). The rule for removing an affix is given in the form: Condition $A \Rightarrow SP|S$. if a word has affix A and the condition that accompanies it is true, the word will be changed to a stem with pattern SP or the affix is replaced by the string S which could be null(removed) or more.

This stemmer used five step rules for the purpose of removing affixes. The first step takes the word to be stemmed as an input and removes double letter reduplication. The second step removes prefix-suffix pair. This step takes the output of the first step as input and checks if the words contains match with any prefix-suffix pair. If the word contains a match and the remaining string has a length greater than the minimum length, then the prefix and suffix are removed from a word. The third step removes prefixes and takes the output of prefix-suffix stripping. In removing a prefix, checking for match in the prefix list and counting length of the remaining string is done. The fourth step removes suffixes by accepting the output from the previous stem and checks if the word contains any match from the list of suffixes. If the word has a match and the remaining string is greater than minimum length the suffix is removed from the word. In the last step the algorithm stems reduplication of single letter. This algorithm has recording rule that is applied after each step is applied for checking some spelling exceptions and making readjustment.

2.3.2 AMHARIC STEMMERS

Amharic and Tigrigna are both Semitic languages and have similar morphological system. The inflection and derivation of words follows similar pattern. So, reviewing researches of stemming algorithms done for Amharic language helps for developing Tigrigna.

Amharic is a language with very rich morphology and the main previous contribution in the area of stemming Amharic is the work by [57] which investigated the effect of stemming for information retrieval [86] have developed a stemmer for information retrieval purposes. The stemmer has been developed in a manner analogous to that used previously in a stemmer for Slovene [80]. Specifically the algorithm first identifies a set of stop-words and then a set of affixes associated with the remaining content-bearing words. They have used the characteristics of the resulting affixes were used to guide the development of the stemmer. The stemmer removes affixes by iterative procedures that employ a minimum stem length, recoding and context sensitive rules, with prefixes being removed before suffixes. Once the stem of the word is obtained, the root is obtained by stripping all the remaining vowels.

Those studies also indicated a positive effect from using the stemmed forms in information retrieval: Alemayehu and Willett compared performance of word-based, stem-based, and root

based retrieval, and showed better recall levels for stem- and root-based retrieval over word based [57].

The other Amharic stemmer was developed by Atelach Alemu and Lars Asker [58]. The stemmer finds all possible segmentations of a given word according to the morphological rules of the language and then selects the most likely prefix and suffix for the word based on corpus statistics. It strips off the prefix and suffix and then tries to look up the remaining stem (or alternatively, some morphologically motivated variants of it) in a dictionary to verify that it is a possible stem of the word. The frequency and distribution of prefixes and suffixes over Amharic words is based on a statistical analysis of a 3.5 million word Amharic news corpus. The stemmer had an accuracy of 85% when evaluated on a limited text consisting of 50 sentences (805 words). The stemmer first creates a list consisting of all possible segmentations of the word that is to be stemmed. In a second step, each such segmentation is then verified by matching each candidate stem against the machine readable dictionary. If no stem matches the dictionary, the stemmer will modify the stem and redo the matching. If more than one stem matches, the most likely stem will be selected after disambiguating between the candidate stems based on statistical and other properties of the stems. In the cases when exactly one stem matches the dictionary then that segmentation will be presented as the output from the stemmer.

2.3.3 OROMO STEMMERS

One of the stemmers for Oromigna language is the one developed by Wakshum [13]. The stemmer uses suffix table in combination with rules that strips off suffix from a given word by looking up the longest match suffix in the suffix list. 342 suffixes are compiled automatically by counting and sorting the most frequent endings. Other linguistically valid suffixes are also included manually. The stemmer finds the longest suffixes that matches the end of a given word and remove.

Debela Tesfaye and Ermias Abebe [59] have developed Oromigna stemmer by considering some problems from the stemmer developed by Wakshum M. and come out with new stemmer. This stemmer adopted some concepts from Porter stemmers. Specifically, Concepts about measure, arranging the rules in clusters, analyzing word formation based on the nature of their endings are taken from porter algorithm. This Afan Oromo stemmer is based on a series of steps

that each removes a certain type of affix by way of substitution rules. These rules only apply when certain conditions hold, e.g. the resulting stem must have a certain minimal length. Most rules have a condition based on the measure. The measure is the number of vowel-consonant sequences which are present in the resulting stem. This condition must prevent that letters which look like a suffix but are just part of the stem will be removed.

2.3.4 WOLAYTA STEMMER

Lemma L. (Lemma, 2003) has developed a stemmer for Wolaytta language based on the morphological nature of the language. As stated in his paper Wolaytta is a language dependent on suffixation to form different forms of a given word. Concatenation of suffixes is common in Wolaytta. As a result, two or more suffixes may be concatenated together and attached to a word. In the language, possible list of combination can be very large making difficult to have complete list of combination (concatenations). Besides, concatenation in the language makes suffixes long ones attaching one suffix to another. Hence, iteratively removing each base suffix one by one is considered as the best choice in this algorithm. As a result of the characteristics of the language, the algorithm adopted iterative approach to develop the stemmer for Wolaytta text. The stemmer developed for stemming Wolaytta text is context sensitive.

In his study, Lemma has employed a semi-automatic means to compile the possible suffix list. In the process of compiling the suffix dictionary, the words in the sample text were first written in reverse order. The reversed list of words was then sorted and frequencies of matching substrings identified, finally the sub-strings which occur more than once are selected as suffixes. First, a word is inputted to the stemmer. The algorithm checks whether a suffix from the list is attached to the inputted word. The suffixes are iteratively stripped from the word and after application of necessary condition, the final word is considered as a stem.

2.4 STEMMING ALGORITHM FOR OTHER LANGUAGES

Stemmers are generally customized for each specific language. Their design requires some linguistic expertise in the language and an understanding of the needs of the specific information retrieval and natural language applications. Stemmers have been developed for a wide range of languages including French Language [82], Slovene Language stemmer by [18], Arabic

Language stemmer by [73] and other many languages. The effectiveness of stemming across languages is varied and influenced by many factors.

2.4.1 ENGLISH LANGUAGE STEMMERS

There are many kinds of stemming algorithm available for English language. These algorithms range from the simplest ones like the table look-up to the complicated ones performing iterative longest matching. Some of the English Language stemmers are developed by [77], [58], [81], [12], [79] and [60].

2.4.1.1 LOVINS STEMMING ALGORITHM

The Lovins Stemmer is a single pass, context-sensitive, longest-match Stemmer developed by Julie Beth Lovins of Massachusetts Institute of Technology in 1968 [77] Lovins' paper was the first ever published description of a stemmer. It defines 294 endings, each linked to one of 29 conditions, plus 35 transformation rules. For a word being stemmed, an ending with a satisfying condition is found and removed. A suitable transformation rule is applied next. The aim of this step is to deal with doubled consonants and irregular plurals.

This early stemmer was aimed at both the IR and Computational Linguistics areas of stemming. This stemmer, though innovative for its time, has the problematic task of trying to address two areas (IR and Linguistics) and cannot do well at either. The approach does not give good results with linguistics, as it is not complex enough to stem many suffixes due to their not being present in the rule list. There are problems regarding the reformation of words. This process uses the receding rules to reform the stems into words to ensure they match stems of other similar meaning words. The main problem with this process is that it has been found to be highly unreliable and frequently fails to form words from the stems, or matches the stems of like meaning words. The Stemmer does not satisfy from the IR viewpoint either, as its large rule set, and its receding stage, affect its speed of execution. The Lovins Stemmer removes a maximum of one suffix from a word, due to its nature as single pass algorithm. It uses a list of about 250 different suffixes, and removes the longest suffix attached to the word, ensuring that the stem after the suffix has been removed is always at least 3 characters long. Then the ending of the stem may be reformed (e.g., by un-doubling a final consonant if applicable), by referring to a list of receding transformations.

2.4.1.2 DAWSON STEMMING ALGORITHM

Dawson stemming algorithm [58] is based on the one developed by Lovins [77] which makes use of iterative longest match approval. It is a complex linguistically targeted Stemmer that is strongly based upon the Lovins Stemmer. Initially, Dawson uses the list that contains 260 English suffixes with associated removal condition codes given by Lovins. But, after he has corrected the list, he brings the total up to about 1200 suffixes. To avoid the problems of storage and processing time taken, the suffixes and their condition code numbers backwards are read, stored and indexes by length and final letter. Dawson does not use recording technique in his algorithm to handle stems and instead used an extension of the partial matching procedure also defined within the Lovins Paper. The basic principle of Dawson's algorithm is if two stem matches up to a certain number of characters and the remaining characters of each stem belong to the same stem ending class, then two stems are of the same form.

2.4.1.3 PORTER STEMMING ALGORITHM

The Porter Stemmer is a conflation Stemmer developed by Martin Porter at the University of Cambridge in 1980 [81]. The Stemmer is based on the idea that the suffixes in the English language (approximately 1200) are mostly made up of a combination of smaller and simpler suffixes. This Stemmer is a linear step Stemmer. Specifically it has five steps applying rules within each step. Within each step, if a suffix rule matched to a word, then the conditions attached to that rule are tested on what would be the resulting stem, if that suffix was removed, in the way defined by the rule. For example such a condition may be, the number of vowel characters, which are followed by a consonant character in the stem (Measure), must be greater than one for the rule to be applied.

Once a Rule passes its conditions and is accepted the rule fires and the suffix is removed and control moves to the next step. If the rule is not accepted then the next rule in the step is tested, until either a rule from that step fires and control passes to the next step or there are no more rules in that step whence control moves to the next step. This process continues for all five steps, the resultant stem being returned by the Stemmer after control has been passed from step five.

The objectives for the development of this stemming algorithm are to improve the performance of information retrieval and the success rate for the suffix discarding will be significantly less

than 100% when the process is evaluated. The Porter algorithm consists of a set of condition rules. The conditions are divided into 3 classes; there are conditions on the stem, condition of the suffix and conditions on the rules. Porter's algorithm uses a dictionary of about 60 suffixes and has only a few context-sensitive and recording rules, and therefore is economical in storage and computing time.

2.4.1.4 PAICE/HUSK STEMMING ALGORITHM

The Paice/Husk Stemmer is a simple iterative Stemmer that removes the endings from a word in an indefinite number of steps [79]. The Stemmer uses a separate rule file which is first read into an array or list. This file is divided into a series of sections, each section corresponding to a letter of the alphabet. The section for a given letter, say "e", contains the rules for all endings ending with "e", the sections being ordered alphabetically. An index can thus be built, leading from the last letter of the word to be stemmed to the first rule for that letter. When a word is to be processed, the stemmer takes its last letter and uses the index to find the first rule for that letter. If a rule is accepted then it is applied to the word. If it is not accepted, the rule index is incremented by one and the next rule is tried. However, if the first letter of the next rule does not match with the last letter of the word this implies that no ending can be removed, and so the process terminates. Once a rule has been found to match, it is not applied at once, but must first be checked to confirm that it would leave an acceptable stem. When a rule is applied to a word, this usually means that the ending of the word is removed or replaced.

2.4.1.5 KROVETZ STEMMING ALGORITHM

The Krovetz Stemmer was developed by Bob Krovetz, at the University of Massachusetts, in 1993 [74]. It is quite a 'light' stemmer, as it makes use of inflectional linguistic morphology. The Krovetz Stemmer effectively and accurately removes inflectional suffixes in three steps, the conversion of a plural to its single form (e.g. '-ies', '-es', '-s'), the conversion of past to present tense (e.g. '-ed'), and the removal of '-ing'. The conversion process firstly removes the suffix, and then though a process of checking in a dictionary for any receding, returns the stem to a word.

2.4.2 ARABIC STEMMING ALGORITHMS

Arabic is a highly inflected language and has a complex morphological structure. Some of the applications of Arabic natural language processing require the basic form of the word (root or stem) to be most effective, therefore stemming process is a necessity. There are several stemming approaches that are applied to Arabic language. The morphology complexity of Arabic makes it particularly difficult to develop natural language processing applications for Arabic information retrieval. In Semitic languages like Arabic, most noun, adjective, and verb stems are derived from a few thousand roots by infixing and creating many variants of words [76]. Arabic is highly productive, both derivationally and inflectionally. Definite articles, conjunctions, particles and other prefixes can attach to the beginning of a word, and large numbers of suffixes can attach to the end. A given headword can be found in huge number of different forms. Distributional analyses of Arabic newspaper text show empirically that there is more lexical variability in Arabic than in the European languages for which most IR and NLP work as been performed. Arabic text has more words occurring only once and more distinct words than English text samples of comparable size. The token to type ratio (mean number of occurrences over all distinct words in the sample) is smaller for Arabic texts than for comparably sized English texts [73]. Arabic orthography also contributes variability that can confuse information retrieval systems. It is not the right to left order of the characters, or the context-dependence of the appearance of the characters that are problematic. These are just rendering issues.

The factors described above make Arabic very difficult to stem. Several stemming algorithms for Arabic have been proposed based on different principles and each produces different sets of stem classifications. Larkey et al. gives a good summary of stemming approaches for the Arabic language [9]. The most common approaches used in Arabic stemming are the light and the root-based stemmers. Root-based Stemming is based on removing all attached prefixes and suffixes in an attempt to extract the root of a given Arabic surface word. Several morphological analyzers have been developed for Arabic, e.g. Khoja and Garside [73].

Light Stemming is used not to produce the linguistic root of a given Arabic surface form, but to remove the most frequent suffixes and prefixes. The most common suffixation includes duals and plurals for masculine and feminine, possessive forms, definite articles, and pronouns.

Several light stemmers have been developed, all based on suffix and prefix removal and normalization. Examples of light stemmers include: Aljlal & Frieder's Stemmer (S.Aljlal) Darwish's Al-Stem [72] , and Larkey et al.'s U Mass Stemmer [75].

All light stemmers adhere to the same steps of normalization and stemming. The main difference among them is the number of prefixes and suffixes removed from each one. During the normalization process, all diacritics, punctuation, and glyphs are removed. The light stemmers had different stop word lists consisting of Arabic pronouns, particles and the like removed after minimal normalization. Based on the test results of the stemmers proved that the light stemmer achieved superior performance over the root-based approach since it reduces sense ambiguity by grouping semantically related words into the same class. Although light stemming can correctly classify many variants of into large stem classes, it can fail to classify other forms that should go together. For example, broken plurals for nouns and adjectives do not get conflated with their singular forms, and past tense verbs do not get conflated with their present tense forms, because they retain some affixes and internal differences.

2.5 EVALUATION METHODS FOR STEMMING ALGORITHMS

There are different methodologies employed to evaluate the performance of stemmers. The manual method, vocabulary reduction and Paice's method are identified as stemmer evaluation methods. In the manual method, a human being, who decides the correct stem for each word, performs the evaluation process. Three evaluation measurements are obtained in this manner: the number of correct results; the number of errors due to over stemming; and the number of errors due to under stemming. One of the purposes of stemmers is to reduce the size of the vocabulary for indexing purposes. The vocabulary reduction is obtained by dividing the number of words in the corpus by the number of stems generated, excluding repetitions.

Stemmers are evaluated using Paice's method based on error counting [87]. According to his method, measures of under stemming and over stemming determine how good a stemmer is outside the retrieval context. In Paice's method, three measurements are implemented in order to make a qualitative comparison between different stemmers: the over stemming index (OI); the under stemming index (UI); and the stemming weight (SW). This method requires a word sampling, with no repetitions, separated into conceptual groups in which the words are

semantically and morphologically related. The SW is given by the ratio OI/UI . Paice has compared different English stemming algorithms in isolation from the context of an IR system and he did not use the traditional precision/recall parameters, an ideal stemmer should stem all words in a group to the same stem. If a stemmed group contains more than one unique stem, the stemmer has made under stemming errors. In an IR system this corresponds to a negative effect on recall. If a stem of a certain group also occurs in other stemmed groups, the stemmer has made over stemming errors, which reduce precision. A good stemmer should therefore produce as few under and over stemming errors as possible.

CHAPTER THREE

MORPHOLOGY OF SILT'E LANGUAGE

3.1 INTRODUCTION

Stemming is a fundamental step in processing textual data preceding the tasks of information retrieval, text mining, and natural language processing. The common goal of stemming is to standardize words by reducing a word to its base (root) [51].

Like the other Semitic languages (Amharic, Tigrigna, Ge'ez etc) Silt'e language has the feature of transmitting different messages with a single word alone. Since a single word can have different forms by adding affixes or changing alphabetical orders (patterns), this situation makes the language morphologically rich[52].

The main purpose of stemmer is automatically conflating different words that have similar meaning to their root by removing their affixes and/or rearranging the vowel patterns. Morphology plays an important role to guide the development of the stemmer, and describe the nature and characteristics of affixes in Silt'e words. Therefore developing (designing) a stemmer for a language requires analyzing, understanding and also modeling the language trend in terms of word formation. The above concept shows the importance of studying the Silt'e morphology for the purpose of modeling it and developing automatic procedures for the language. Thus, this chapter deals with the concepts of inflectional and derivational morphologies of Silt'e language.

3.2 THE WRITING SYSTEM OF SILT'E LANGUAGE

Since the 1980s, Silt'e has been written in the Ge'ez, or Ethiopic, writing system, originally developed for the now-extinct Ge'ez language and most known today in its use for Amharic and Tigrigna [52]. However, Silt'e vowels differ considerably from the typical set of seven vowels in languages such as Amharic, Tigrigna (see Table 3.2).

Ethiopic Fidel has the phonetic structure (order) of seven columns (see Table3.1).

See lists of constant, numbers, punctuations and vowels of Ethiopic (Ge'ez language) in appendix 1

First order (ግእዝ)	Second order (ካዕብ)	Third order (ሳልስ)	Forth order (ራብዕ)	Fifth order (ሃምስ)	Sixth order (ሳድስ)	Seventh order (ሳብእ)
(አ läl	አ•lul	አ. lil,	አ lal,	አ lel	አ ə	አ lol
(ሀ lhäl	ሀ•lhul	ሂ lhil,	ሃ lhal,	ሄ lhel	ሀ lhə	ሀ lhol

Table 3.1: Table shows the seven alphabetical order of Ethiopic (Ge'ez) alphabet.

3.2.1 VOWELS AND CONSONANTS OF SILT'E LANGUAGE

The Ge'ez script makes distinctions among only seven vowels, so some of the short-long distinctions in Silt'e are not marked. In practice this probably does not interfere with comprehension because there are relatively few minimal pairs based on vowel length. In written Silt'e, the seven Ethiopic vowels are mapped onto the ten Silt'e vowels as follows: the five short vowels are አ lal, አ lol, እ lil, አ•lul, አ lel and the five long vowels are አ laal, አአ lool, አ. lil, አ•አ•luul, አ.አ. leel. Only the short and long a and the short and long i have been distinguished in the alphabet. The short a is indicated by the first order (form) of the Ethiopic script and the long a by forth order (form). The short i and consonant alone are indicated by sixth order (form) and the long i is indicated by the third form. The third form is also used word finally when the word ends in i [52]. Table 3.2: illustrated more the vowel mapping by using practical words that show the vowel mapping in Silt'e language more clearly with simple example of order of alphabet.

<u>The order (form) of Silt'e alphabets</u>						
1	2	3	4	5	6	7
ሀ lhal	ሀ•lhu (huu)	ሂ lhii	ሃ lhaal	ሄ lhe(hee)	ሀ lhi	ሀ lho (hoo)
ኀ lnal	ኀ•lnu(nuu)	ኂ lnii	ኃ lnaal	ኄ lne(nee)	ኀ lnii	ኀ lno(noo)
ኘ lnāl	ኘ•lnū(nūu)	ኚ lnīi	ኛ lnāal	ኜ lnē(nee)	ኘ lnīi	ኘ lnō(nōo)

<p>➤ ä → a: አለፈ <i>alafa</i> 'he passed'</p> <p>➤ u → u, uu: ሙት <i>mut</i> 'death', ሙት <i>muut</i> 'thing'</p> <p>➤ i → i</p> <p> • ii: ኢን <i>iin</i> 'eye'</p> <p> • word-final i: ሙሪ <i>mari</i> 'friend'</p> <p> • i ending a noun stem: ሙሪክ <i>marikal</i> 'his friend'</p> <p> • impersonal perfect verb i suffix: ባለ <i>baali</i> 'people said'; በባለም <i>babaalim</i> 'even if people said'</p> <p>➤ a → aa: ጋራሽ <i>gaaraaš</i> 'your (f.) house'</p> <p>➤ e → e, ee: ኤፍ <i>eeff</i> 'he covered'</p> <p>➤ i → ə</p> <p> • i (except as above): ኢንግር <i>ingir</i> 'foot'</p> <p> • consonant not followed by a vowel: አስሮሽት <i>asroost</i> 'twelve'</p> <p>➤ o → o, oo: ቆሮ <i>k'oč'e</i> 'tortoise', ከ'ቆሮ <i>k'ooč'e</i> 'he cut'</p>
--

Table3.2: The mapping of seven Ethiopic vowels to ten of Silt'e and order of Silt'e alphabet

Silt'e has a fairly typical set of consonants (26 consonants including **ጥ ገገ**) for an Ethiopian Semitic language. There are the usual ejective consonants, alongside plain voiceless and voiced consonants and all of the consonants, except /h/ and /ʔ/, can be geminated, that is, lengthened. The symbols /p/ and /ʔ/ (glottal stop) play only a marginal role in the system, **ገገ**, because it appears in only a few words and **ገገ**, because (as in Amharic), it is often omitted. The consonants of Silt'e shown in Table 3.3:

<p>ሀ h , ለ l , መ m , ረ r , ሰ s , ሸ š , ቀ q , በ b , ተ t , ቸ č , ነ n , ኘ ñ , አ a , ከ k , ወ w , ዘ z , ገ ʒ , የ y , ደ d , ጅ j , ገ g , ጠ ɬ , ጨ č' , ጰ p , ፈ f , ጥ p </p>
--

Table 3.3: consonants of silt'e writing system

3.3 MORPHOLOGY

Morphology is a branch of linguistic that studies and describes how words are formed in language[55]. Similarly, [58] defines morphology as the study of the structure of words. There are two kinds of morphology: inflectional and derivational. Inflectional morphology is concerned

with the inflectional changes in words where word stems are combined with grammatical markers for things like person, gender, number, tense, case and mood. Inflectional changes do not result in changes of parts of speech. Derivational morphology deals with those changes that result in changing classes of words (changes in the part of speech). For instance, noun or an adjective may be derived from a verb.

3.3.1 SILT'E MORPHEME

A morpheme is a minimal linguistic unit of a language that carries meaning and that can not be further decomposed in to meaningful units [54]. A morpheme in Silt'e can be free or bound, where a free morpheme can stand as a word on its own where as a bound morpheme can not occur on its own as a word, for example these are free morphemes of Silt'e:- ሰብ|sabl 'person', ጋር|gaarl 'house', ቃቀ|kaawaal 'coffee', they can stand by themselves. In contrast to these bound morphemes of silte can't stand by them selves. For example we can add suffix ቸ|lcal to the word ሰብ|sabl to show paucal number; therefore ቸ| is bound morpheme of Silt'e [56].

3.3.2 WORD FORMATION IN SILT'E

Affixing (adding suffix, prefixes and infixes), compounding, duplicating (reduplicating) and changing vowel patterns are used to give different word forms in Silt'e language [53].

The addition of suffix, prefix and infix are the usual ways of word formation in Silt'e. Although it is not common, the addition of prefix and suffix at the same time is also used to form word variant. The possibility of prefixing and suffixing more than three prefix or/and suffix result in long word formation. Hence, the complex morphological structure of a single Silt'e word can give very large number of variants.

The application of compounding is another means to contribute in the process of word variant formation in Silt'e language. A compound is made up of two or more simple words (roots) and has a meaning that is different from that of its parts: አባዴ|abbaadde/abaadel, 'parents' from አባ|abbo/abol 'father' and አዴ|adde/adel 'mother'. ወዘንቁብሌ|wazana-ጃubel 'faint-hearted', ወዘን|wazanal 'heart' and ቁብሌ|ጃublal 'deficient'.

Duplication is also the other means which is used to form word variation: there are at least two distinct kinds of reduplicative formations in Silt'e language: type one involves the reduplication of the first letter (that is ላ llaal) ላላላላላላላ 'send many times' formed from ላላ llaahal 'send', type two by reduplicating the second letter by using its fourth form (that is መመመመ): ጅግመረ ljiimaamaral 'starts' to little which is formed from ጅመረ ljiamaral 'start'.

3.4 INFLECTIONAL AND DERIVATIONAL MORPHOLOGY OF SILT'E

Inflectional affixes describe word as stems are combined with grammatical markers for things like person, gender, number, tense, and case. There are five kinds of speech in Silt'e: adjectives, nouns, verbs, adverbs, and propositions [53]. Prepositions and conjunctions are totally unproductive for IR purpose or other natural language processing. The same thing happens for adverbs and is few in number. Therefore, the discussion of derivational and inflectional morphology concentrates on the remaining three parts of speech, namely verbs, nouns and adjectives.

To study the morphological nature of the language, the researcher used different sources such as [53], [52], [57] and [56]. The translation of the Latin words of the available document that talks about Silt'e morphology to Ethiopic-script is performed and prepared by the researcher. The following morphological analysis part adopted (taken) from [53].

3.4.1 NOUN'S INFLECTION AND DERIVATION

3.4.1.1 NOUN'S INFLECTION

Nouns in Silt'e language inflected for definiteness, gender, number, and case. Morphologically nouns have the following general structure in Silt'e language:

<prep><distriib> stem <acc> <part> <poss> <art> <voc> <cop>

3.4.1.1.1 DEFINITENESS, NUMBER AND GENDER

Definiteness, number and gender are interrelated in Silt'e language. Silt'e has two definite articles, the suffixes -ኢ፡፡-iiil (-ይ፡፡-yi after vowel) for masculine, and the -ቴ፡፡-tel for feminine. However, these two articles can fulfill further functions; the masculine article -ኢ፡፡-iiil can signal that one is dealing with a large number or collective of items, regardless of their gender, and the

feminine article can indicate that one is dealing with one item (singularize) and /or with a small specimen. For example, in addition to denoting definiteness, the expression **ጋርቴ|gaartel** ‘the (f) house’ indicates that one is talking about a single house, and also that the house is considered to be small. By contrast, the expression **ጋርአ|gaariil** ‘the (m) house’ with the masculine definite marker is neutral with regard to size, and it can refer to a collective or large number of houses rather than one.

Gender is mostly determined naturally, not grammatically. Thus humans and animals are given the gender according to their sex. A comparatively small number of words referring to inanimate objects are specified for gender; thus **ወር|waril** ‘moon’ is usually given masculine gender and **አዳር|ayri** ‘sun’ feminine, hence **ወሪ|warii** |‘the moon’ **አዳሪቴ|ayritel** ‘the sun’. Similarly trees are usually treated as **feminine**. The collective use overrides natural gender: **እንዳቺ መጠ** **|Indaaccii mata|** ‘the woman have come’. The gender markers for possessive pronouns shown in table 3.4:

1 st p		2 nd p				3 rd p			
singula r	plural	Singular		plural		singular		Plural	
		m	F	M	f	m	f	m	f
ኤ eel	ነ nal	አ/አሀ aa/aahal	አሽ aas	አሙ ammul		ከ lal	ሽ sal	ኒሙ niimmul	

Table 3.4: Possessive suffix for noun inflection

Table 3.5 shows samples illustrate combinations of several of these suffixes:

አደኑውከ				
<u>Aaddeenawka</u> ‘as to his mother <acc.>’				
አደ ኑ ው ከ				
<u>Aadde</u>	<u>-na</u>	<u>-w</u>		<u>-ka</u>
Mother	acc.	Pragm. Marker	-w	his
ልክኔ				
<u>Likkinee</u> ‘it is my measure/size’				
ልክ ን ኤ				
<u>Likk</u>	<u>-n</u>			<u>-ee</u>
Measure	copula			1 st sg-poss

Table 3.5: example of combination of suffixes

Semantically Silt’e has three forms for number singular, paucal, and plural. However, morphologically only the paucal is marked – the other two categories are unmarked. Thus while **ኡንቸ** |uuncal refers to a small number of stones <paucal>, **ኡን** |uunl refer either to one or large number of stones. Context is used to disambiguate such unmarked expression. Table 3.6 shows Noun paucal formation with suffixation of **-ቸ** |čal:

Noun	Gloss	Suffixes	Paucal
ቻፍ Caaf	‘leaf’	-ቸ čal	ቻፍቸ caaf čal
ሰብ Sab	‘person’	-ቸ čal	ሰብቸ sabčal
በርት Bart	‘stick, club’	-ቸ čal	በርት čalbartiččal
ቦላሌ bolaale	‘long trousers’	-ቸ čal	ቦላሌቸ bolaalčal
ጦላቢ toollaabil	‘beggar’	-ቸ čal	ጦላቢቸ Toollaabčal

Table3.6: Noun paucal formation with suffixation of **-ቸ** |čal

There are two basic morphological types of paucal formation for nouns; suffixation of **-ቸ**– **čal** <-**cč**a after consonant cluster of geminates> and reduplication of the last consonant.

As you can see from table 2, word final **short vowels** are omitted before the paucal suffix.

Many nouns – mostly ending in a short vowel – form the paucal by reduplicating the last consonant, inserting the vowel **አ**laal and replacing the original final vowel by **አ**lol. When the last consonant is marked, only the degeminated, single consonant is repeated.

Most words that form the paucal by reduplication allow the alternative of reduplication plus suffixation as well. Table 3.7 shows paucal noun formation with reduplication and reduplication plus suffixation:

Noun	Gloss	Paucal with reduplication	Paucal Reduplication+suffix
ቦሶ bosol	‘Young enset plant’	ቦሳሶ bosaasol	ቦሳስቸ bosaas čal
አለገ alagal	‘Stranger’	አለጋጎ alagaagol	አለጋጎቸ alagaag čal
አሞሌ amolel	‘salt bar’	አሞላሎ amoolaalol	አሞላለቸ amolaalčal
ቡሬ burrel	‘big tin, big can’	ቡሬሮ burraarol	ቡሬረቸ buraarčal
በለ balal	‘Calamity’	በላሎ balaalol	በላልቸ balaalča

Table 3.7: paucal noun formation with reduplication and reduplication plus suffixation

3.4.1.1.2 CASE

Syntactic case relations are marked either by prefixes (some of which may also be considered prepositions – see below) or suffixes: the accusative case is marked by **ነ**nal, **አ**lal (**አ**l-al after consonants); the genitive by **የ**lya-l; the dative/benefactive relation by **ላ**lla-l; the instrumental or determinative relation by **በ**ba-l, the ablative by **ቸ**lta-l, and the vocative by **አ** l-ol.

3.4.1.1.3 MORPHOLOGICAL PROPERTIES AND PROCESSES WITH NOUNS

Nouns can end in any consonant or semivowels <except the glottal stop>, single or geminate, subject to the general phonological constraints of the language, or in one of the short vowels **አ** al, **ኤ** lel, **እ** lil, or **ኦ** lol (no noun ending in **ኦ** lul has been found.) The following morphophonemic changes are common with prefixes: **Cአ+ወ**→**Cኦ**/Ca+wa→Coo (**here C is to represent constant letter**) – **የ**lyal plus **ወሮ**lwarol ‘dragon’ : **የ**ሮlyooro l‘of a dragon’; **Cአ+አ**→**Cኦ** /Ca+i →Cii – **ለ**lla-l plus **እንባብ**linbaabl ‘snake’: **ለ**ንባብ llinbaabl ‘ for a snake’: (**here V represent vowel letter**) **Cአ+V**→ CV/ Ca+V → CV: **በ**lba-l plus **አደን**ladanl ‘cat’ **ባደን**baadanl ‘by the cat ‘; **ለ**lla-l plus **አም**lumil ‘uncle <maternal>’ **ለ**ምlluumil. Common morphophonemic changes with suffixes are the following. Word final vowels are dropped directly before paucal suffix (illustrated above in table2 and 4). Word final vowels are lengthened immediately preceding the acc. suffix **-ነ**–nal: **ማጠ**lmaatal, **ማጣነ**lmaataanal; preceding the pragmatic marker **-መ**–ml and **ወ**–wl: **አዴኘ**laddennal; **አዴኛም**laddennaaml; preceding the masculine definite article: **ቡኝ**l bucol, **ቡኝይ**lbucooyl.

3.4.1.2 NOUN’S DERIVATION

Nouns can be derived from verbs. The infinitive has already been introduced above. Nouns of agent can be formed in two ways. Verbs ending in a palatal radical add the suffix **ኢሎ** l–iilol to the perfective stem. Table 3.8 shows some example:

Note that the suffix may be added to the palatalized or non-palatalized form of the verb stem.

Non–palatal verbs add the suffix **እ**l–il to the last root consonant of the perfective stem, palatalizing that consonant where possible; the last, non –thematic vowel is lengthened to **አ**l aal: **ቆመረ**lqoommaral ‘be strong’; **ቆመሪ**lqoommaril ‘strong’; **አመሰለ**lamasalal ‘pretend’
 □ **መሰይ**lmasaayl ‘pretending’. Table 3.8 shows noun derivation from perfective stem by adding the suffix **ኢሎ** l–iilol:

Table 3.8: noun derivation from perfective stem by adding the suffix **ኢሎ** |–iilol

3.4.2 ADJECTIVE INFLECTION AND DERIVATION

Adjectives are inflected for number and derived from nouns. Actually most of the nouns in Silt’e can be use as an adjective.

አሞቁ Amoqel ‘work reluctantly and badly’	አሞቁሎ Amoqiilol ‘relectant worker’
ዛጨ Zaace ‘watch look after’	ዛጨሎ Zaaciilol ‘watch man’
በኛ Bacel ‘weep, cry’	በኛሎ Baciilol or በኪሎ bakiilol ‘some one who craies easily’

3.4.2.1 ADJECTIVES’S INFLECTION

There are two basic morphological types of paucal formation for adjectives like nouns do. Suffixation of **-ቸ**– č a <-cča after consonant cluster of geminates> and reduplication of the last consonant. Tables [3.9 and 3.10] show the two types of paucal formation for adjectives:

Adjective word	Gloss	Suffixes	Paucal
ቻፍ Caaf	‘leaf’	-ቸ čal	ቻፍቸ caaf čal
በርት Bart	‘stick, club’	- ቸ čal	በርት –čalbartiččal
ጦላቢ toollaabil	‘beggar’	- ቸ čal	ጦላቢቸ Toollaabčal

Table 3.9: Adjective paucal formation with suffixation of **-ቸ**|čal

Adjective word	Gloss	Paucal with reduplication	Paucal Reduplication+suffix
አለገ aalagl	‘Stranger’	አለገገ alagaagol	አለገገቸ alagaag čal
ጉመራ gumaral	‘white’	ጉመራሮ gumararol	ጉመራሮቸ gumaraarčal

ቡሬ burrel	‘big tin, big can’	ቡሬጅ burraarol	ቡሬሬቸ buraarčal
በለ balal	‘Calamity’	በላሎ balaalol	በላልቸ balaalča l

Table 3.10: paucal adjective formation with reduplication and reduplication plus suffixation

In Silt’e most of the nouns used as an adjective. There are no special adjectival forms for comparison; if the point of comparison is explicit it occurs in the ablative case with the suffix -ኮ|kol. However, very often comparison is expressed by verbs.

E.g. አ-ሀ ተሂኮ መኒጃም /Uha tahiko manijaam ‘He is prouder than I’

3.4.2.2 ADJECTIVES’S DERIVATION

Simple adjectives have no special structural characteristics: ፈየ|fayyal ‘good’ ቦዝ|boozl ‘bad’ የሮሬ|yaroorel ‘big’ ቀል|qalll ‘small’. **Adjectives can be derived from nouns in a number of different ways.** None of these formations seems to be generally productive, but they are all lexically determined. Table 3.11 shows the derivation of adjectives from nouns:

<p>Denominal adjectives in -አኘ -anna :</p> <p>ጉጉ Guttl ‘middle’ ጉጉኛ guttanna ‘average’</p>
<p>Denominal adjectives in -አተኘ -atanna :</p> <p>ከስብ Kasbl ‘work’ ከስብተኛ kasbatanna ‘industrious’</p> <p>ፈሮ Farol ‘injustice’ ፈሮተኛ farootanna ‘unjust’</p>
<p>Denominal adjectives in -አንቺኦ -aancia :</p> <p>ብተር Bitarl ‘marriage’ ብተራንኝ bitaraanco ‘newly –wed’</p>
<p>Denominal adjectives in -አ -al (-አ -al) -ም -ml (-አ -al):</p> <p>መኒጅ Maniijal ‘pride’ መኒጅም maniijaaml ‘proud’</p>
<p>Denominal adjectives in -አ -al(አ -al)ታም (a) taaml</p> <p>አይብ Aybl ‘milk’ አይባታም aybaataaml ‘giving much milk’</p>

Table 3.11: The derivation of adjectives from nouns

3.4.3 INFLECTIONAL AND DERIVATIONAL MORPHOLOGIES OF VERBS

Following the typical Semitic root-pattern morphology, the lexical meaning of most verbs is specified by the consonants of the root. Underived verbs have two to four root consonants, or radicals. (Since the glottal stop is rarely pronounced, it is not written here. However inflectionally these seemingly vowel -initial roots behave exactly as if they had an initial root consonant.) Table 3.12 shows examples of underived verbs with 2, 3 and 4 radicals.

II – RADICAL VERBS:	III-radical verbs:	IV-radical verbs:
ገበ Gabal ‘enter’	ሀረተ haratal ‘swallow’	ድነበተ Dinabatal ‘be surprised’
በደ badal ‘take’	ረወተ rawatal ‘run’	ስነተለ sinatalal ‘know’

Table 3.12: examples of verbs from 2 to 4 radicals

The final radical of a verb can be of four kinds: ending in a consonant only, ending in consonant followed by አ | al, ending in a non-palatal consonant with following አ |al palatalized to ኦ |el.

The nature of the final radical can best be seen in the imperative masc. singular; palatal consonants are ‘depalatalied’ preceding back vowels, so the underlying non-palatal consonant appears, for example, preceding the infinitival ending –አቸ-ootl. Table 3.13 shows the nature of final radical:

(a) CONSONANT ONLY:			
ዋበ waaba	ዋብ waabl	‘give!’	
(b) Consonant plus አ lal:			
ቁረ qeeral	ቁረ qiiral	‘watch’	
(c) Non-palatal consonant plus ኤlel:			
ኖዜ nozəl	ኑዝ nuzl	‘be anger’	
(d) Palatal consonant plus ኤlel:			
ሰቼ sace	ሰኝ sicl	ስኮት sikoot	‘to drink’

Table 3.13: the nature of final radical of verbs

The nasal element appears in two and three radical verbs: አንዜlanzel ‘see’, ኤንዘleenzal ‘hold’, አንቱ loontel ‘close’, አንደረ landaral ‘spent the night’, ኤንቀፈleenqafal ‘embrace’ ; most of these verbs begin with a vowel (underlying preceded by a glottal stop); however , at least one verb has been found that begins with a consonant: ሰንቸ lsooncel ‘smell’. Verbs with a nasal element are special in that the nasal does not count as a radical consonant: thus አንዜlanze l‘see’ inflects like a II- radical, and አንደረ landaral ‘spend the night’ like a III-radical verb, though their roots have three and four consonants respectively.

3.4.3.1 ASPECT

Each verb has three stem forms: perfective stem, imperfective stem and a third stem not marked for aspect, hence referred to here as the non-aspectual stem. The perfective stem is used to form simple past, present perfect and past perfect tenses, the imperfective stem for the formation of the present/future and past-imperfective tenses; the non-aspectual stem serves for infinitive and imperative. Table 3.14 shows examples of the three forms of the verb:

Perfective stem መሰክ masak- 		
(a) Simple past	መሰክ <u>masakal</u>	‘he guided’
(b) Present perfect	መሰካን <u>masakaanl</u>	‘he has guided’
(c) Pluperfect	መሰክ ናር <u>masaka</u> <u>naarl</u>	‘he had guided’
Imperfective stem መስክ mask- 		
(d) Present /future	እ+መስካን i+maskaanl	‘he guides/will guide’
(e) Continuous past	እ+መስክ ናር i+mask naar	‘he was guiding’
Non –aspectual stem ምሰክ m (i)sak- 		
(f) Imper. /jussive	ምሰክ misakl	‘guide!’ (2 m sg)’
	የምሰክ Yamsakl	‘let him guide!’
(g) Infinitive	ምሰኩት misakootl	‘to guide’

Table 3.14: the perfective, imperfective and non-aspectual forms of verb

The stem formation makes use of two distinct morphological means: changes in consonant-vowel pattern of the root and rising of the thematic vowel. Thus Silt’e verbs can be divided in to two broad classes:

Class 1: All verbs with thematic **አ**al; taking the perfective stem as basic, the other two are derived by omitting the first and the second vowel respectively. Table 3.15 and table 3.16 show the pattern of consonant and vowels in class 1 verbs and examples of class 1 verbs of Silt’e:

Perfective	imperfective	non – aspectual
<C>CVC<VC>	<C>CVC<C	CC<VC>
		CCC<C>
(C represents one or more consonants of a word; V represents one or more of vowel sequence)		

Table 3.15: the pattern of consonant and vowel in class one verbs

ሰቼ Sac-el ‘drink’	እሰቻን i-sac-aanl	ሰች sicl
ብተክ Batakal ‘pull out’	እብትካን i-batk-aanl	ብተክ bitakl

Table 3.16: examples for class 1 verbs of Silt’e

Case 2: class two verbs are all verbs with a thematic vowel other than **አ** | a|. Thematic **ኤ** | e| is raised to **እ** | i|, **ኦ** | o|, to **ኡ** | u| in the non- aspectual stem: Table 3.18 shows examples of class 2 verbs:

ጄጄ Jeejel ‘reaches’	እጄጃን i-jeej- aanl	ጄጅ jijj
ኖዜ Nozel ‘be angry,	እኖዛን i-noze-aan	ኑዝ nuzl

Table 3.18: examples of class 2 verbs of Silt’e

Raising verbs with thematic **ኣ** | aal are exceptional in that the **ኣ** | aal is raised to **ኡ** | i| in the imperfective aspect, see these examples:

ጮመ Caamal ‘test well’	እጮማን i-ciim-aanl	ጮም caaml
ላሀ Laahal ‘send’	እላሃን li-liih-aanl	ላሀ laahl

(Note that not all verbs with thematic **ኣ** | aal undergo rising; **ራጅ** | raaje | is an example: **እራጃን** | i-raaj-aanl.)

Class 2 verbs with more than two radicals also omit the vowel of the last radical in both the imperfective and non- aspectual stem forms, see these examples:

ሼበለ Seebalal ‘dance’	እሼብላን i-seebl-aan	ሼብል siibl
ቶቀሰ Tooqasa ‘beg’	እቶቅላን i- tooks –aanl	ቶቅስ tuuksl

3.4.3.1.1. VERB INFLECTION

The inflectional structure of Silt’e verbs is as follows:

<neg ><pers> STEM <pers> <ben/detr> <cps> <aux>

The double occurrence of the person marker indicates that Silt’e marks person by both prefix- and suffixes in forms based on the imperfective and non-aspectual stems.

Like in other Ethio-Semitic languages, sentence constituents other than the subject can optionally be marked in the morphology of the verb, as is indicated by the benefactive (ben), detrimental/instrumental (detr) and complement person suffix (cps) categories above. The complement person suffixes can occur without either the benefactive or detrimental/instrumental markers, but these last two categories require the presence of the person markers.

Table 3.19 shows the suffixes that used to form the present time are as follows:

1 st p				2 nd p				3 rd p			
Singular		Plural		singular		plural		singular		plural	
m	F	m	F	m	f	m	f	m	f	m	f
አሁ-	አው-	አን-		አሀ-	አሽ-	አሙ-		አን-	አት-	አን-	
aahul	aaw	aanl		aahal	aasl	aammul		aanl	aatl	aanl	

Table 3.19: suffixes used to form the verb in present tense

The Impersonal አን |–aanl, the past auxiliary is ኖር |naar|– (past tense form of verb of existence).

3.4.3.1.2. PERFECTIVE

In Silt'e verb, the perfect nature of verbs uses to create the past and complete actions. These types of verbs are the bases to other forms that may inflect to their base. The suffixes in table 3.20 are added to the perfective stem to mark the subject in the verb:

1 sg	ከ- kul after C(consonant)
	ሀ- hul, ወ- wl elsewhere
1p	ኃ- nal
2f	ሻ- sl, word finally
	ሻ- sil elsewhere
2m	ከ- kal after C(consonants)
	ሀ- hal, ኣ laa lelsewhere
2p	ከ-ሙ- kumul after C
	ሙ- mmul elsewhere
3m	ኣ- el in palatal verbs
	ኣ- a elsewhere
3f	ተ- tal preceding object suffix
	ተ- tl elsewhere
3p	ኣ- ul
Impers.	ኣ- il

Table 3.20: suffixes added to the perfective stem to mark the subject in the verb

The category ‘impersonal’ in the last row of the table is used to indicate what some unspecified group of individuals or people in general do. Hence it is glossed here ‘people’.

3.4.3.1.3. SIMPLE PAST

The suffixes in table 3.20 added to the perfective stem to form the simple past tense of the verb. Table 3.21 shows the simple past form of a verb:

Simple past

Simple past of C – verb	Suffix	Gloss
መሰከ masakkul	-ከ -kul	‘I guided’
መሰከ masakkal	-ከ -kal	‘You (2m) guided’
መሰከሽ masaks	-ሽ -s	‘You (2f) guided’
መሰከ masakal	-ከ -kal	‘He guided’
መሰከት masakt	-ት -t	‘She guided’
መሰከን masaknal	-ን -nal	‘We guided’
መሰከ-ሙ masakkumul	-ከ-ሙ -kumul	‘You (pl) guided’
መሰከ masakul	-ከ -kul	‘They guided’
መሰከ masakil	እ -il	‘People guided’

Table 3.21 Simple past form of verbs

Simple past of palatal verbs	Gloss
ከሼሁ Kaseehu/w	‘I wanted’
ከሼህ Kaseehal	‘You (m) wanted’
ከሼሽ Kasees	‘You (f) wanted’
ከሼ kasel	‘He wanted’
ከሼት Kaseet	‘She wanted’
ከሼን Kaseenal	‘We wanted’
ከሼሙ Kaseemmul	‘You (p) wanted’
ከሱ kasul	‘They wanted’
ከሱ kasil	‘People wanted’

Table 3.22: Simple past form of palatal verbs

Simple past of Ca-verbs	Gloss
በላሁ/በላው Balaahu/w	‘I ate’
በላሀ Balaaha	‘You (m)ate’
በላሽ Balaas	‘You (f)ate’
በለ Balal	‘He ate’
በላት Balaatl	‘She ate’
በላን Balaanal	‘We ate’
በላሙ Balaammul	‘You (pl)ate’
በሉ Balul	‘They ate’
በሊ Balils	‘People ate’

Table 3.23: Simple past form of Ca-verbs (constant followed by vowel a)

3.4.3.1.4. PRESENT PERFECT

Two further tenses can be formed from the perfective stem: the present perfect and the pluperfect by adding the present and the past tense auxiliary respectively. Added to the simple past tense forms given above tables, they yield the set of present perfect forms in Table 3.24:

Present perfect form	Gloss
መሰኩ masakkool	‘I have guided’
መሰካ masakkaal	‘You have (2m) guided’
መሰክሼሽ masakseesl	‘You have (2f) guided’
መሰካን masakaanl	‘He has guided’
መሰክታት masaktaatl	‘She has guided’
መሰክናን masaknaanl	‘We have guided’
መሰኩዋሙ masakkumoommul	‘You have (pl) guided’
መሰኮን masakoonl	‘They have guided’
መሰኬን masakeenl	‘People guided’

Table 3.24: present perfect form of verbs

These forms involve the following vowel contraction processes: $አ\text{-}lul + አ\text{-}laal \rightarrow አ\text{-}አ\text{-}lool$, $አ\text{-}lal + አ\text{-}laal \rightarrow አ\text{-}laal$, $እ\text{-}lil + አ\text{-}laal \rightarrow ኤ\text{-}ኤ\text{-}leel$. For the first person singular there is the following additional contraction involved: $መሰክ + ኩ + አሁ | masak + ku + aahul \rightarrow መሰኮሁ | masakkoohul \rightarrow መሰኮ | masakkool$. The resulting ending $-ኮ\text{-}kool$ is plausible considering that the first person ending $ሁ\text{-}hul$ has a common voiced variant $ዉ\text{-}wl$. Thus the ending $ኮ\text{-}kool$ could easily be explained if we assume the following further contraction rule: $አ\text{-}አ\text{-}lool + ዉ\text{-}wl \rightarrow አ\text{-}አ\text{-}lool$.

For palatal and Ca –verbs the present perfect is formed analogously, except that the following additional contraction rule is required for the 3rd person masculine of palatal verbs. $ኤ\text{-}lel + አ\text{-}laal \rightarrow አ\text{-}laal$: ከሻን|kasaanl ‘he has wanted’.

3.4.3.1.5. PLUPERFECT

The pluperfect is formed by the simple past tense forms followed by the auxiliary **ናር|naar|** which is uninflected in all persons except the first person singular. Table 3.25 shows the pluperfect of verbs:

Pluperfect	Gloss
መሰከ ናርኩ Masakku naar	‘I had guided’
መሰከ ናር Masakka naar	‘You (m) had guided’
መሰከኛ ናር Masaks naar	‘You (f) had guided’
መሰከ ናር Masaka naar	‘He had guided etc’

Table 3.25: the pluperfect of verbs

3.4.3.1.6. NEGATIVE

Negative forms based on the perfective stem are formed with the prefix **አል|al-|**. In this way negatives can be obtained from all three tenses based on the perfective stem: **አልመሰከ**|almasaka| ‘he did not guided’, **አልመሰከን**|almasakaan| ‘he has not guided’, **አልመሰከ**|almasaka naar| ‘he had not guided’.

3.4.3.2 IMPERFECTIVE (NON-PAST)

The imperfect verb in Silt’e language uses to describe non-past action. It uses prefixes and/or suffixes. As is typical for Semitic languages, some of the person markers used with imperfective stem involve both prefixes and suffixes. The affixes are shown in Table 3.26:

ይ (y)/ እ li-l	/--C	1sg, 3m, 1pl, 3pl, impers.
ይ ly-l	/--V (for vowel ending verbs)	
ት ti-l	/C (for consonant ending verbs)	2m, 2f, 2pl,3f
ት it-l	/--V	
The suffixes are እ-l-l 2f, ት-na 1pl, ኡ-l-ul 2pl and 3pl, l-l impersonal		

Table 3.26: the prefixes and suffixes of imperfective verb

. The second person female suffix exerts a palatalizing influence on the last consonant of the stem, according to the following rules: **d, g →j; t, k →c; s→s, z →z, n →n, l →y.**

3.4.3.2.1. CONTINUOUS PAST

The presence of these affixes in table 3.26 is best seen in the continuous past tense, which consists of the conjugated imperfective stem followed by the past tense auxiliary.

Continuous past	Gloss
እመስክ ናር Imask naarkul	‘I was guiding’
ትመስክ ናር Timask naarl	‘You(m) were guiding’
ትመስኪ ናር Timaski naarl	‘You (f) were guiding’
እመስክ ናር Imask naarl	‘He was guiding’
ትመስክ ናር Timask naarl	‘She was guiding’
እመስክነ ናር Imaskina naarl	‘We were guiding’
ትመስኩ ናር Timasku naarl	‘You (pl) were guiding’
እመስኩ ናር Imasku naarl	‘They were guiding’
እመስኪ ናር Imaski naarl	‘One was guiding’

Table 3.27: the affixes of imperfective verb in continuous past tense

3.4.3.2 .2. PRESENT /FUTURE

The present/future tense is formed by adding the present tense auxiliary (አሙ-|aammul, አን-|aanl, አት-|aatl, አን-|aanl, አሙ-|aammul, አን-|aanl, አት-|aatl, አን-|aanl) to the imperfective with its person affixes(-ኩ-|kul,-ከ-|kal,-ሽ-|s, -ከ-|kal,-ት-|t,-ን-|nal,-ኩሙ-|kumul,-ኩ-|kul,እ-|il).

Table 3.28 shows this form:

Present/future	'Gloss'
እመስካሠ limaskaahul	'I guide'
ትመስካሠ timaskaahal	'You (m) guide'
ትመስኬሽ timaskeesl	'You (f) guide'
እመስካን limaskaanal	'He guides'
ትመስካት timaskaatl	'She guides'
እመስክን limaskinnanal	'We guide'
ትመስኩሙ Timaskoommul	'You (pl) guide'
እመስኮን limaskoonl	'They (pl) guide'
እመስኬን limaskeenl	'People'

Table 3.28: the present/future form of Silt'e verb

The vowel contraction processes involved in these formations: (The second person female suffix exerts a palatalizing influence on the last consonant of the stem, according to the following rules:

d, g →j; t, k →c; s→s, z →z, n →n, l →y.)

3.4.3.2.3. NEGATIVE

In the negative the imperfective stem takes the prefixes in table 3.29 rather than the person prefixes listed above (ይ (lyl)/እli-l, □ly-l, ትlti-l, □lt-l).

3.4.3.3 NON – ASPECTUAL

3.4.3.3.1. INFINITIVE

The infinitive is formed by suffixing **አት** |oot| to the non- aspectual stem: **ምሰከት** |misakoot| 'to guide'; **ክሶት** |kisoot| 'to want'.

3.4.3.3.2. IMPERATIVE AND JUSSIVE

The jussive forms also involve a combination of prefixes and suffixes; the imperative uses only suffixes. For both imperative and jussive the suffixes are the same as those used with the imperfective stem. The prefixes are shown in Table 3.31:

ለ L (አ a) -	1sg, 1pl (የ ya - is also used for 1pl)
ት T (አ a) -	3f
ይ Y (አ a) -	3m, 3pl, impersonal

Table 3.31: the prefixes of imperative and jussive verb

The parentheses around the **a** above table indicate that this vowel is omitted preceding another vowel. Table 3.32 shows the imperative and jussive forms:

ለምሰክ Lamsakl	let me guide	ምሰክ Misakl	guide! (m)
ምሰኪ Misakil	‘guide (f)’	የምሰክ Yamsakl	‘let him guide’
ተምሰክ Tamsakl	‘let her guide’	ለምሰክነ Lamsaknal	የምሰክነ /yamsakna
	‘let us guide’	ምሰኩ Misakul	‘guide! (pl)’
የምሰኩ Yamsakul	‘let them guide’	የምሰኪ Yamsakil	‘let one guide’
ዩ.መ Yiimal	‘let him swear’	ተንት Tuuntl	‘let her close’

Table 3.32: Shows the imperative and jussive forms of a verb

3.4.3.3.3. NEGATIVE

The negative infinitive is formed from the affirmative by prefixing **አል**|al-|: **አልምሰኩት** |almisakootl ‘to not guide’.

The suffix set for the negative imperative/jussive is the same as for the affirmative, but the prefixes are different. The prefixes are shown in Table 3.33:

አል - al-	1sg, 1pl
አት at-	2m, 2f, 2pl, 3f
አይ lay-	3m, 3pl, (1pl), impersonal

Table 3.33: the prefix for negative imperative/jussive verb formation

Examples : **አልምሰክ**|almisakl ‘let me not guide’; **አትምሰክ**|atmisakl ‘don’t guide (m)’; **አይምሰክ** laymisakl ‘let him not guide’ **አይምሰኪ**laymisaki |‘let people not guide’.

3.4.3.3.4. MOODS

So far most verb forms have been given in the active mood. The passive mood is formed by prefixing **ተ**|ta-| to perfective stem forms (**ት**|t-| after other prefix), and **ት**|t-| to imperfective and non-aspectual forms: **ተመሰከ**|tamasakal ‘he was guided’: **የትመሰከ**|yatmasakal ‘he who was guided’; **እትመሰካን**|itmakaanl ‘he will be guided’.

3.4.3.3.5. MORPHOPHONOEMIC CHANGES WITH PREFIXES

Contraction processes take place with vowel and glide – initial verbs. With h- initial verbs, the **u/h** is elided (omitted), a following short vowel is lengthened, and the prefix is realized by the allomorph it has preceding vowels. Table 3.34 shows examples of the changes with prefix:

ዩዳን Yeedaanl	እ + ሄዳን li+heedaanl	‘he will go’
ዮኖን Yoonaanl	እ + ሆኖን li+hoonaanl	‘he will be’
ያርታን Yaartaanl	እ + ሀርታን li+hartaanl	‘he will swallow’

Table 3.34: the morphophonemic changes with prefixes

This rule does not apply after the affirmative jussive prefixes **ለ** |la-|, **የ**|ya-|, and **ተ**|ta :**ለሂድ**|lahiidl ‘let me go!’ **የሂድ**|yahiidl ‘let him go !’ ; **ተሂድ**|tahiidl ‘let her go!’

With **ወ**|w-| initial verbs, regular morphophonemic changes occur only with verb forms beginning in **ወ**|wa-| and **ወ**|w-| followed by consonant: **pf +wa → oo**; **pf +wC → pf + uC**:

(Here the **C** represents the constant terms)

ዮርዳን Yoordaanl	ዮርዳን+ y+wardaanl	‘he goes down’
እሎርድ lloordl	እለ+ወርድ lila+wardl	‘he does not go down’
አሉተ Aloota l	አል+ ወተ lal+watal	‘he did not go out’
አዩቀ Ayuqal	አይ+ወቀ lay+wqal	‘let him not hit’
አቱቀ Atuqal	አት ወቀ lat+wqal	‘don’t hit’
አሉቆት Aluqootl	አል ወቆት lal+wqootl	‘not to hit’

Table 3.35: Morphophonemic changes for verbs which begin with **ወ**|wa-| and **ወ**·|w-

When prefixes are added to vowel – initial verbs, the vowel, if short, is lengthened. Table 3.36 shows these phenomena:

ያርሳን Yaarsaanl	ይ+ አርሳን y+arsaanl	‘he will plough’
ታድጋት Taadgaatl	ት+ አድጋት t+adgaatl	‘she will leave (tr)’
ቲኛት Tiinaatl	ት+ እኛት t+inaatl	‘she will sleep’

Table 3.36: vowel initial verbs lengthened after they get prefixed

3.4.3.4 DERIVATION OF VERBS

It seems best to deal with derived stems from the morphological point of view, since the semantic characteristics tend to overlap and are unpredictable at times. The following derivational processes have been found:

- Prefixation of **አት**|at-| plus change of thematic **አ**|a| to **ኤ**|ee|; this has a causative meaning. Thematic vowels other than **አ**|a| get lengthened. Verbs with long thematic vowels only add the prefix. Table 3.37 shows some examples:

ወቀ waqal	‘hit’;	አትቁቀ atweeqa	‘cause to hit’
ጀመረ Jammara	‘begin’;	አጀመረ ajjeemmaral	‘cause to begin’
ደነበተ Dinabatal	‘be frightened ‘;	አደነበተ addineebatal	‘make frightened ‘
ፎጋ Fogel	‘inflate’;	አትፎጋ atfooge	‘cause to inflate’
ዛጨ Zaace	‘herd’;	አዛጨ azzaace	‘caused to hared’
ካፈረ Eeffel	‘cover’;	አቴፈረ ateeffel	‘cause to cover’

Table 3.37: derivation of causative verb

This causative formation is the only derivation process that can be applied to virtually any verb, including other derived verbs.

- Prifixation of **አ**la- to intransitive verbs; tends to result in a transitive meaning, but in many cases the resulting **meaning is not predictable**. Table 3.38 shows some examples of this type of derivation:

ራጅ Raajel	‘be old’;	አራጅ araajel	‘make old’
ወከበ Wakabal	‘buy’;	አወከበ awakabal	‘sell’
ገበ Gabal	‘enter’;	አገበ agabal	‘place in side; marry’

Table 38: derivation of transitive verbs by prifixing **አ**la-

- Preffixation of **ተ**ta-; in most cases this passivizes underived transitive verbs. Table 3.39 shows examples of derivation of passive verbs:

ጨኘ Ceenel	‘give birth’;	ተጨኘ taceenel	‘be born’
ዋበ Waabal	‘give’;	ታበ taabal	‘be given’
ኤወደ Eewadal	‘tell’;	ቴወደ teewadal	‘be told’

Table 3.39: derivation of passive verbs by prefixing ተ|ta-

- Lengthening of thematic **አ** |al to **አ** |aal; this usually accompanied by the prffixation of ተ| ta-l; semantically it often indicates a reciprocal, iterative or very intensive process. The Prefix may be replaced by አት|at-l, resulting in a causative meaning: **ፈጅ** |fajel ‘finish’; **ተፋጅ** |tafaajel ‘destroy each other ‘; **አትፋጅ** |atfaajel ‘cause people to destroy each other’. **ረወተ** |Rawwatal ‘run’; **ተራወተ** |taraawatal ‘to run here and there’; **አትራወተ** |atraawatal ‘cause to run here and there’.
- Reduplication: there are at least two distinct kinds of reduplicative formations: type 1 involves the reduplication of the first letter, type 2 that of the second letter by **preciding** with it’s forth order. Semantically reduplication often indicates multiple actions – either doing something repeatedly or doing something to many objects, people etc. For a number of verbs it indicates that the activity is done ‘a little’, not completely or properly. Table 3.40 and Table 3.41 shows examples of type1 and type 2 reduplication respectively:

Type 1 reduplication <C1VC2 (V) →C1aaC1VC2 (v) (here C and V shows consonant and vowel respectively)

ለሀ Laahal	‘send’;	ለለሀ laalaaha	‘send many people or something to many places’
ቁረ Qeeral	‘wait’;	ቃቁረ qaaqeeral	‘watch repeatedly, many times’
ጨኘ Ceenel	‘give birth’;	ጫጨኘ caaceenel	‘ give birth many times’

Table 3.40: Type 1 reduplication

Type two reduplication < (C1) VC2 (V) →(C1) iC2aaC2 (V)>

(Here C and V shows consonant and vowel respectively)

ፈጅ Fajel ‘finish’;	ፍጃጅ fijaajel ‘finish many things’
አጅ Ajel ‘hit’;	አጃጅ lijaaje ‘hit many times, in many places, many people’
ኤመ Eemal ‘slander’;	ኤማመ imaamal ‘slander many times’
ኤንዘ Eenza ‘hold’;	ቲንዘዘ tiinzaazal ‘hold each other’
ቀተለ Qatalal ‘kill’	ቅታተለ qitaatalal ‘kill many animals, people’
ጀመረ Jammaraal ‘begin’	ጅማመረ jimaammaraal ‘to start a little bit’

Table 3.41: Type two reduplication

Some verbs employ both types of reduplication, with different meanings: **ፍጌ** | fogel ‘inflate’; type one reduplication: **አትፋፍጌ** | latfaafoogel ‘cause to be swollen in many places’; type 2 reduplication: **ፉጋጌ** | fugaagel , **ፍጋጌ** | fogaagel ‘inflate many times’

With both types the original thematic vowel may also be lengthened: **ክነበለ** | kinabalal ‘return(vt)’; **ክናናበለ** | kinaanaabalal ‘turn again and again, turn over many times’ ;**ፍጌ** | fogel ‘inflate’ **አትፋፍጌ** | latfaafoogel ‘cause to be swollen in many places’

The reduplicated verb stems can have one or more of the derivations that a basic verb can have: causative in **አት** | -lat-, passive in **ተ** | -ta-, and though much more rarely, transitive with **አ** | al.

Table 3.42 shows examples of the three forms (reduplication, transitive and causative) together:

ሶንቹ Sooncel ‘smell’		
REPLICATED	TRANSITIVE	CAUSATIVE
ሳሶንቹ	አሳሶንቹ	አትሳሶንቹ
<u>Saasoonce</u>	<u>asaasoonce</u>	<u>atssaasoonce</u>
Many things	to smell different things	cause each other to smell something
Give a smell, something		
Gives a smell many times		
ጅመረ Jammara ‘begin’		
REPLICATED	TRANSITIVE	CAUSATIVE
ጅማመረ	አጅመመረ	ተጅማመረ
<u>Jimaammara</u>	<u>ajjimammara</u>	<u>tajmaammara</u>
To start a little bit	cause to be started a little bit	to be started a little bit

Table 3.42: reduplication, transitive and causative derivation

SUMMARY

This chapter presents the Silt'e language morphology, both inflectional and derivational. The inflectional and the derivational morphologies in Silt'e language involve prefixation, infixation, prefix-suffix pair and suffixation. Hence, Silt'e words are very complex morphologically.

The language also shows some morphophonemic changes in prefixes. Hence, this also can increase the varieties of the words as well.

The complexity of the language is one of the main reasons for a language to desire a stemmer. A stemmer helps to conflate variants of words as well as to access documents in natural language text.

In Silt'e gender is determined naturally, therefore nouns do not inflected for gender but definiteness and possessive marker affixes are used according to the gender of the noun(that is naturally given) to mark it in the verb.

Verbs, nouns and adjectives are very reach morphologically to agree with person, number, and gender (not for noun) to agree with the subject in Silt'e language, the discussions are mainly focused on noun, adjective and verb morphologies.

Silt'e language has two types of reduplication and discussed in this chapter as type1 reduplication and type 2 reduplication. Type1 with first consonant reduplication and the second (type 2) indicated by the second consonant reduplication.

The forms of Silt'e verb such as perfect (past and completed actions), imperfect, jussive, and infinitive are discussed. Its derivational characteristics: causative, replicated transitive causative and passivating of verbs also discussed.

The next chapter presents the development of a stemmer to conflate variants of Silt'e words with respect to concepts in this chapter.

CHAPTER FOUR

A STEMMING ALGORITHM DESIGN AND IMPLEMENTATION (PROTO TYPE) FOR SILT'E LANGUAGE TEXT

4.1 INTRODUCTION

A morphology part of the Silt'e language has been presented (reviewed) in chapter three. It has been shown that the main word formation process in Silt'e is done through affixation. As described before the parts of affix are: - prefix, suffix, prefix-suffix pair and reduplication (type1 and type2 in Silt'e language). These affixes are used for inflecting and deriving words. Nouns are inflected for definiteness, number, and person. Verbs are inflected for gender, number, person, aspect and tense. Like other Semitic languages, Silt'e experiences complex morphological structure that resulted in very large variants of a word. In designing retrieval systems and other natural language processing systems for the language, reducing these variants into one form improves the performance of the systems. This can be achieved by a conflation technique, which is usually stemming. The main purpose of a stemmer is reducing different variants of words in to their standard form (root). Thus, the main reason for this chapter is designing (developing) a stemming algorithm for the language. Therefore, the following part presents development of a stemming algorithm for the language. The compilation of stop words and affixes and evaluation of the stemmer has been presented as well.

4.2 THE CORPUS

The researcher has utilized different sources of sample text for the development of the stemmer and the experiments. Since there is no document collection for any of Ethiopian languages available likes TREC and CLEF collections for the English language, the researcher took documents from different sources and organized them into one collection. The corpus consists of 10,922 tokens or 5,288 distinct words. The document has been used for stop word list collection; affix collection and testing the algorithm. The test data was compiled from the document randomly to check the stemmer from different angles of varieties of words.

4.3 WORD DISTRIBUTION OF SILT'E

Word distribution in sample text documents of a language helps to study language's behavior, and this distribution can be shown using word-ratio (number of distinct words to total number of words). This helps to show how many words are morphologically distributed within a document [9]. **Table 4.1** shows the word ratio of sample document for Silt'e language.

Name	Description	Word tokens	Word types (Distinct)	Word ratio (distinct to total word)
SampleDocument	Different books of Silt'e in academic area and other: <ul style="list-style-type: none">• Text book from grade 1 to 8• Text book from grade 9 to 10• Fictions• poems	10,922	5,288	48.42

Table 4.1: Number of words and their distributions to Silt'e sample text data

Sample document in **Table 4.1** used for compilation of stop words, affixes (prefix and suffix) list and for testing the developed algorithm.

To compare the sample sets of Silt'e word distribution with other languages such as Ge'ez [44], Tigrigna [40], and Amharic [42], **Table 4.2:** show sample text of (Text1) with 1,518 numbers of words was taken randomly from the SampleDocument in **Table 4.1**. The word ratio for sample text of Geez, Tigrigna, Arabic and English languages has been adopted from [44].

language	Text	Total number of words	Distinct words	Word ratio (distinct to total word)
Silt'e	Text 1	1518	823	54.22 %
Ge'ez	Lukas	1,866	1,064	57.02 %
Tigrigna	Text1	1,632	918	56.25 %
Arabic	Text1	1,600	902	56.38 %
English	Text1	1,600	621	38.81%

Table 4.2: comparison of word distribution ratio of the sample data set among different languages

The higher value of the word ratio corresponds to more distinct words in a text and vice versa [88]. The word ratios obtained for Silt'e sample text Text1 is somewhat similar to Ge'ez, Tigrigna, Amharic and Arabic texts as shown in Table 4.2. However, it is absolutely different from English text. This similarity among Silt'e, Ge'ez, Amharic, Tigrigna and Arabic might be due to the fact that each is in Semitic language groups. As Table 4.2 shows numbers of unique words found in Silt'e document is similar with Semitic languages and it is larger compared with the English language distinct words. Hence, Silt'e language has comparable distinct words with morphologically very complex languages when its word distribution is compared with others based on their sample data sets. And when it is compared with English it is morphologically very complex.

4.4. DOCUMENT (MORPHOLOGICAL) PREPROCESSING

The stream of characters in natural language text must be broken up in to distinct meaningful units before any language processing can be performed. Preprocessing is the most important part of all text processing. Preprocessing must ensure that the source text be presented to NLP in a form usable for it. For example, NLP programs usually need their input to be tokenized, i.e. text elements usually word forms or sentences are identified and placed on separate lines of the input [50]. In the preprocessing stage this study addresses tokenization, normalization and stop word removal.

4.4.1. NORMALIZATION

All punctuation marks, control characters, numbers and special characters are removed from the text before the data is processed.

In contrast to Amharic and Tigrigna, Silt'e uses distinct 26 Ethiopic alphabets (in Silt'e language writing system different letters that has similar sound, like ሐ, ኀ, ኸ, ሠ, ዐ are not used, instead only ሀ, አ, ሰ are used). But, in some document of the real life, people experienced writing of these similar sound characters which are used in Amharic and Tigrigna writing systems. Even though these symbols sound the same, in Silt'e language they are not part of the alphabet. But in case people use them they must be changed to alphabets that are used in Silt'e, that have similar sound form. Therefore, for the sake of flexibility the system will handle this kind of situations as they appear in the documents.

4.4.2. TOKENIZATION

In this study, words are taken as tokens. Since all punctuation marks have been converted to spaces, space is used as a word demarcation. Hence, if a sequence of characters is followed by space, that sequence is identified as a word. A consecutive sequence of valid characters was recognized as a word in tokenization processes.

4.5. COMPILATION OF THE STOP WORD LIST

The stopword list was compiled from the Sample document in [Table 4.1](#) by collecting the most frequently occurring words. Frequency of words in the document was generated by using python program. [Table 4.3](#) lists the top 54 frequent items from the document. As can be seen from the [Table 4.3](#), the stopword list consists of prepositions (ከ-ሰጥ|in), demonstrative adjectives (ሂተይ|this, the masculine singular), articles (እታይ|the, masculine singular, እተይ| the, feminine singular) and conjunctions (such as |ገንመ| but, |ዎ|and). As we can see in [Table 4.3](#) there are words with high frequency but are not stop words (such as አሰነት, ቅሬታቶ, ገረድነት, የአዳ, አደ, አህላቅ, ጋር, ባሎት, የገረድ, ወክት). As indicated in chapter three, function words such as prepositions, conjunctions and articles in Silt'e exist affixed to words. For this reason the

frequency of function words in Silt'e is low. Moreover the stopwords in Silt'e can exist in different forms by concatenating affixes. Therefore this also decreases their frequencies. **Table 4.4** shows sample list of words that have low frequencies but are stopwords. See the complete list in appendix IX. Because of this reason, the use of frequency alone did not help to generate all the stopwords and the researcher was used research paper that talks about the morphology of Silt'e language to collect stopwords from the SampleDocument. Many function words were also added manually to the stopword list by consulting dictionary [53] and research paper [52]. **Table 4.5** shows the list of sample stop words that are used in the development of the algorithm. Generally, 184 stopwords are used to develop the stemmer.

No	Word	frequency	No	Word	Frequency
1	ደር	73	41	አደ	17
2	ሰብ	71	42	ሀድ	16
3	አነግነ	71	43	ጀዋብ	16
4	ኤት	62	44	ኮሎ	16
5	ወክት	57	45	ጉት	16
6	ደረሳሶ	57	46	ሚሽ	15
7	ወልድ	44	47	ያሎ	15
8	አሰጥ	43	48	ዮነኮ	15
9	ቅሬታቶ	39	49	ገረድነት	15
10	ባድ	39			
11	ሃለት	38			
12	ሱልቸ	36			

13	የገረድ	35			
14	አሰነት	35			
15	ሉላሉሌ	34			
16	ያለይ	29			
17	ያለይሙ	29			
18	የአዳ	29			
19	አድ	29			
20	ገረድ	28			
21	አሀላቅ	27			
22	ዋ	26			
23	ወግ	25			
24	ግዝ	24			
25	ጋር	24			
26	የሰብ	21			
27	ባለ	21			
28	የውን	21			
29	ስንቅ	20			
30	ገንገናም	20			
31	ገዳ	20			
32	ባሎት	20			
33	ጃንገ	19			
34	ድራሙ	19			
35	ያቀትሎን	18			
36	ያቴራን	18			

37	የጎ	18			
38	ቡር	17			
39	በቅራቲ	17			
40	ሬር	17			

Table 4.3: The most frequent words from the SampleDocument

Words with lower frequencies but are stopwords

Word frequency	Word frequency	Word frequency	Word frequency
ገናኒሚ 2	ልቶን 3	ኢነ 2	ጊዘቸ 2
ጃንጎ 3	ቦንት 2	ኢንኩ 2	ሂንኩ 2
የሰሀደደይ 2	ያፍቴቴ 2	ኢነኮ 2	ዩሃን 2
ለገገኒሞ 2	አቴታን 1	ባታይ 3	በገገሽ 2
አደደሞ 2	አነይ 2	ገነግነ 2	
ሂተይ 2	አታይንገ 2	ሎኔት 2	ለሆነ 2
ጊንመ 2	ሊሊ 2	ሂንኩምከ 2	ሂንኩም 5
አደኒሙ 2	ዮናነይ 2	ታት 2	ዮኑ 2

Table 4.4: words with lower frequencies but are stop words

አዩ	አልተቴ	እሲ	እነይ
እታይ	እትታይ	አሲ	በልዳሌ
አተቴ	እትተቴ	እነይ	አልቀሬ
አታይ	አትታይ	አነይ	ደር
አተቴ	አትተቴ	እቲ	ኤት
ዩታይ	እብታይ	አቲ	ፈሬ
ዩተቴ	አብተቴ	እቢ	ፎኖ
ያታይ	እብታይ	አቢ	ግነ
ያተቴ	እብተቴ	ሱር	ቀደ
እልታይ	አብታይ	አድ	
እልተቴ	አብተቴ	አድአድ	
አልተይ			

Table 4.5: sample list of stop words

4.6. COMPILATION OF AFFIXES

In Silt'e affixes consist of three different types, which are the prefix, suffix, and infix. Unlike English stemmers which work quiet well just by removing suffixes and few prefixes to obtain the stems, an effective and powerful Silt'e stemmer not only must be able to remove suffixes, but also prefixes, and infixes as well. Without removing all this affixes, the stemmer cannot be effectively used to stem Silt'e documents.

4.6.1. COMPILATION OF PREFIXES

A set of prefixes that are used to develop the algorithm is collected from the SampleDocument by using Evan Gut's research paper "Concise grammar of Silt'e" [53] as a guide and from the paper

itself. Generally 36 prefixes are used to develop the stemmer. Table 4.6 shows some of the prefixes collected for the development of the algorithm

እለው	ሰ
እለው	ቃ
በሰ	ጫ
ከይ	ከ
እት	ለ
እት	የ
እል	ይ
እለ	እ
እለ	ተ
ከል	ት
ከት	

Table 4.6: sample list of prefixes

4.6.2. COMPILATION OF SUFFIXES

To collect a list of suffixes from the SampleDocument the researcher uses Evan Gut's research paper "Concise grammar of Silt'ee" [53] as a guide and collected from the paper itself. A total of 68 suffixes are used for the purpose of developing the stemmer. Table 4.7 shows samples from a list of suffixes used in the algorithm.

አሙ	ከሞሙ	ኧ
ከሙ	አን	ይ
ኒሙ	ኤን	ሚ
አሀ	አኘ	ዋ
አወ	አተኘ	ሶ
አሀ	አንቾ	ዬን
አሽ	አም	ኒሙ
አነ	አታም	
አት	አታም	
አሙ	ኧሎ	

Table 4.7: sample suffixes

4.7. THE PROPOSED STEMMER

In this research the data have been represented in Unicode and the stemmer accepts Unicode data directly without transliterating it to Latin form. Therefore, the length is used to represent the size of the word while developing the algorithm. The minimum length for a meaning full Silt'e word is two. Most words in Silt'e language have length of two [53]. Thus, the minimum length of word size has been used in the algorithms based on this concept.

Techniques developed for English and Semitic languages such as Amharic, Tigrigna and Geez have been studied. Some of the techniques used in these algorithms were incorporated in developing the Silt'e stemmer. The algorithm uses iterative approach but it removes the longest affix first.

As indicated in chapter three, the process of inflection and derivation in Silt'e is done through one or the combination of the following processes.

- Pure external affix (without modifying the stem)
- External affix and with modification of the letter(s) of the stem at the beginning, end or other positions
- Pattern change (insertion or deletion of consonant and vowels)
- Reduplication

Accordingly, these characteristics were taken in to consideration in developing the stemmer.

Of the two approaches (context free and context sensitive), a context sensitive approach was considered appropriate, as Silt'e is morphologically complex language (See Chapter three).

In the stemmer developed, three steps were used for the purpose of removing affixes. The purpose of each step is given as follows:

The first step removes prefixes. In removing a prefix, checking a match in stop word list, prefix list and counting length of the word and checking contexts sensitive conditions is done. If the word gets match in the prefix list and satisfies the conditions for context sensitive, it will be returned without removing the prefix, otherwise the prefix stripping process takes place.

In the second step the removal of suffix is done. First the word is checked against stop word list. If there is no match in the stop word list then the length of the word will be checked. If it is greater than two then the suffix list file will be opened and checked if there is a match of suffix with word. If the word has a match and doesn't satisfy the context sensitive conditions the suffix stripping process will be performed. Otherwise no suffix will be removed.

The last or the third step is used to stem reduplication of letters from the word. This step deal with type 1 and type 2 reduplication in Silt'e language. First the length of the word is checked. If length greater than two then the program checks existence of reduplication in a word. If the word holds reduplication of the first letter then the first reduplicated letter will be removed from the word. Otherwise the reduplication of the second letter will be cheked if exist then get stripped. If the word doesn't have any kind of reduplication then it will returned as it is.

Except the removal of reduplication, after each step is applied, the word is checked against some context sensitive conditions. This is because to avoied the removal of non genuen affixes.

The proposed stemmer gets the word and then checks if it is in stop word. If a match is found the stemmer starts again and takes the next word if end of file not reached else the word will be given for the next process and the length of the word will be counted. If length of the word is less than three it will not be stemmed and get written to the stemmed file otherwise the word will be given the affix removal process to strip the appropriate affixes. The algorithm checks each time the word against the stop word after prefix and suffix stripping process. The recoding and context sensitive checked and based on the conditions the appropriate action will be taken. Flow chart of the stemmer algorithm shown in figure 1.

4.8. RECODING AND CONTEXT SENSITIVE RULES

A context-free stemmer is one that removes strings with out consideration of the remaining stem, and can thus remove strings that are similar to, but that actually are not, genuine affixes. English examples of such behavior include removing “re-“from “regular” or “-al” from “metal” [86]. And equally poor results are obtained if context free is adoped for Silt’e. For example አባት|abotl ‘father’, አረሰ|arasal ‘plough’, have a leading stringአ lal which is in the list of prefix removing this leading string result in ባት|botl, ረሰ|rasal doesn’t have any relation to the correct stem. There are many examples in Silt’e and it was hence decided that a context-sensitive approach should be adopted, with the affix removal is being controlled by two action cods and three conditions. The two types of context sensitive actions are:

Action 1: don’t perform any affix removal

Action 2: remove affix

4.8.1. RECODING RULS

Recoding1: The word ends in the letter of six order form: this is to recode the change in order form at the end.

Recoding 2: The word begins with affix: If character ‘ፆ’ are part of identified stem word change it to ፆፀ

4.8.2. CONTEXT SENSITIVE CONDITIONS

Condition 1: The affix followed by: order of letters. This is to avoid the removal of non-genuine affixes.

Condition 2: The letter to be removed is reduplicated and in the first order: take action 1, this is to avoid removing non-genuine type 1 reduplication like **ጠጠፈ** |tataral

Condition 3: The length of word to be considered: if the length of the word is less than two, take action 1 this is to maintain the minimum stem length of the word in the language

Certain amount of terminal character recoding can occur in Silt'e. Examples of characters that can be recoded (converted) are ሻ, ሽ, ቸ, አድ, ኩ. . A recoding routine was thus included to handle such occurrences. Table 4.8 shows some recoding examples:

Ending	Recode	Example
ሻ	ት	አሰሻ to አሰት
ት	ደ	ሄት to ሄደ
ሽ	ድ	ባሽ to ባድ
ቸ	ድ	ገራቸ to ገራድ

Table 4.8: examples of recoding

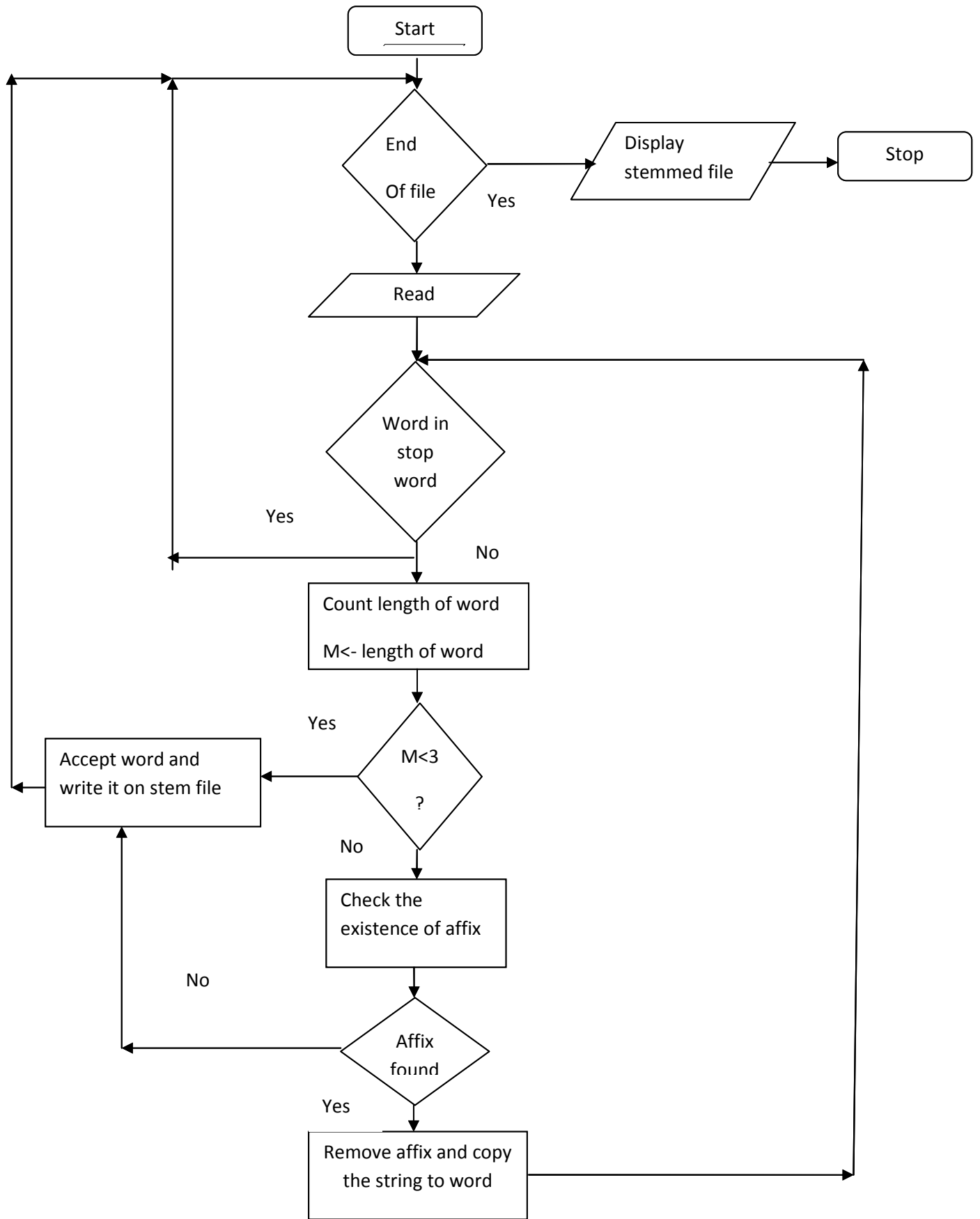


Figure 4.1: Flow chart for the general affix removal algorithm

4.9. COMPONENTS OF THE STEMMER

The stemmer processes prefix stripping and suffix stripping based on the affix lists stored in the affix list files. As the affix found, the word will be checked with the context sensitive conditions if the conditions are not fulfilled the affix removal takes place. There are also affixes that are stripped using iterative approach with rules of the language, for instance, lists of characters such as **ŋ**lbal, **ʔ**lyal, **ʔ**lta or **ŋ**lsal can occur by concatenating each others as prefix of words. In removing of affixes, the sequences that are done in this work are : <checking word length>, <prefix checking>,[prefix removing], <suffix checking>,[suffix removing],<checking letter reduplication type1 and 2>,[removing the duplication type1 or two]. In each activity there are stemmed word length checkers, stop word checkers and the recoding and context sensitive rules to result in good stemmed words as possible.

4.9.1 THE AFFIX- REMOVAL TECHNIQUES

Before this technique is applied, some common stop words from sample text are removed by matching with stop word lists. In the process of affix removals, the length of each word is also checked.

The stemming process is done whenever length of word is greater than two. This is because the minimum length of words in Silt'e is two [53]. The stemmed word, after affixes are removed is also checked with list of stop words and if it is part of stop words, it will not go to the affix removal process or will not be part of the stemmed file.

Basically, two actions are taken in the stemming processes.

Action 1: do not remove any affix

Action 2: remove the concerned affix

To take any one of the above actions, there are rules (conditions) that are used to check whether the word is context sensitive or not and apply action 1 or action 2. The conditions are:

Condition 1: after getting the assumed prefix or suffix, if length of word is not more than two or fulfill context sensitive rule take action 1.

Condition 2: if part of the assumed affix found and the length of a word is greater than two and if the word is not fulfill context sensitive rule take action2.

In the stemming process, based on the above conditions the appropriate action is taken.

4.9.1.2 PREFIX STRIPING

The procedure takes a word and checks the existence of a true prefix, and removes it whenever the condition is fulfilled. For example a verb which begins with **አ**lal, and if the word remaining is more than two it will strip the prefix after checking the context sensitive rules. If there is no condition to be fulfilled the prefix is stripped otherwise the word is returned. For instance the word **አቦት** labotl will not go under prefixing process because, **አ**lal is not a true prefix for **አቦት** labotl as the remaining string **ቦት** lbotl has no meaning in Silt'e language. Whereas **አ**lal is a true prefix for **አበለ** labalal. Hence it undergoes the stripping process and **አ**lal will be removed and the remaining stem **በለ** lbalal, 'eat' will be returned. Table 9 shows the prefix stripping stemmer algorithm:

<p>1 Get word</p> <p>2 Read stop word file and check if there is match with word</p> <p style="padding-left: 40px;">If word exists in the stop word list then</p> <p style="padding-left: 80px;">Go to step 5</p> <p style="padding-left: 40px;">Else</p> <p style="padding-left: 80px;">Find the length of the word</p> <p>3 If if length of word < 3 then return word and Go to step 5</p> <p style="padding-left: 40px;">Else :</p> <p style="padding-left: 80px;">If length of (word) <= (Maximum Prefix length + 2) then</p> <p style="padding-left: 120px;">assign the first length (word)-2 characters to a temporary</p> <p style="padding-left: 120px;">prefix variable</p> <p style="padding-left: 40px;">Else</p> <p style="padding-left: 80px;">assign the first Maximum Prefix length to a temporary prefix variable</p> <p>4 Read the prefix list file and check if match found with a temporary prefix variable:</p> <p style="padding-left: 40px;">If match found then check conditions:</p> <p style="padding-left: 80px;">If context sensitive rule fulfilled:</p> <p style="padding-left: 120px;">Return word and Go to step 5</p>
--

<p>Elese :</p> <p>Remove prefix and Go to step 2</p> <p>Else</p> <p>If length(a temporary prefix variable) >1 then:</p> <p>assign length(a temporary prefix variable)-1 characters to a temporary prefix variable</p> <p>And go to step 4</p> <p>Else</p> <p>Return word and Go to step 5</p> <p>5 If end of file not reached Go to step 1</p> <p>Else:</p> <p>Stop processing</p>
--

Table 4.9: prefix striping stemming algorithm

```

1 Get word
2 Read stop word file and check if there is match with word
    If word exist in the stop word list then
        Go to step 5
    Else
        Find the length of the word
3 If if length of word < 3 then return word and Go to step 5
Else :
    If length of (word)<=(Maximum suffix length+2) then
        assign the last length (word)-2 characters to a temporary suffix
    Else
        assign the first Maximum suffix length characters to a temporary suffix
4 Read the suffix list file and check if match found with temporary suffix :
    If match found then ckeck conditions:
        If context sensitive rule or recoding rule fulfilled:
            Take the appropriate recoding
            Return word and Go to step 5
        Elsele :
            Remve prefix and Go to step 2
    Else
        If length( temporary suffix ) >1 then:
            assign length( temporary suffix )-1 characters to temporary suffix
            And go to step 4
        Else
            Return word and Go to step 5
5 If end of file not reached Go to step 1
Else:
    Stop processing

```

Table 4.10: Suffix Stripping Stemming Algorithm

4.9.2 SUFFIX STRIPPING

The procedure takes a word and checks the existence of suffix match in the suffix list. If match found then will checked against the context sensitive conditions and based on the conditions the appropriate action will be taken. The process is almost the same with the prefix striping except the stripping and the affixes are on the right side of the word. Table 4.10 shows the suffix striping stemmer algorithm:

4.10. LETTER REDUPLICATION STRIPING

4.10.1. TYPE 1 REDUPLICATION

This procedure deals with reduplication of the first letter of the word, i.e. **ላላሀ**llaalaahal. The procedure takes as an input the output of the previous procedure. It checks the existence of the reduplication and removes the first letter of the word if it is reduplicated and the letter is not in the first order of the alphabet. Otherwise the procedure returns the word as it is. For instance the word **ላላሀ**llaalaahal contains reduplication and the first letter will be removed and resulted in word **ላሀ**llaahal, 'He send'. The procedure checks the alphabetical order of the letter to avoed the removal of the reduplicated letter if it is like Amharic word **ጠጠረ**ltataral.

Table 4.11 shows type 1 reduplication stemming algorithm:

1 Get Word

2 Find the length of the word

3 If length of word < 3 then:

return word

Go to step 4

Else:

If first letter of word is the same to second letter of word then:

If the first and the second letters in first order form then:

Return word

Go to step 4

Else :

Remove the first letter

return the stripped word

Go to step 4

Else:

If the first and second letters are the forth and the fifth orders of the same alphabet:

Remove the first letter

return the stripped word

Go to step 4

Else: return word and Go to step 4

4 If end of file not reached Go to step 1

Else :

Stop processing

Table 4.11: type1 reduplication stripping stemmer algorithm

4.10.3. TYPE 2 REDUPLICATION

This procedure deals with the reduplication of the second letter (infact by reduplicating it self to its fourth form), i.e. **ጅማመረ** ljimaamaral. It checks the existence of the reduplication and removes the second letter of the word if it is reduplicated and returns the stripped word, otherwise the procedure returns the word as a final output. For example, the word **ጅማመረ** ljimaamaral contain reduplication and the second letter will be removed and result in word **ጅመረ**ljamaral. The procedure will not deal with words like **አበበ** lababal, in Amharic because the reduplication in the word **አበበ** lababal is not in the fourth and first order of the alphabet, i.e letter **በ**lal and **በ**lal are both in first order of the alphabet. Table 4.12 shows type 2 reduplication stemming algorithm.

<p>1Get Word</p> <p>2 Find the length of the word</p> <p>3 If length of word < 3 then:</p> <p style="padding-left: 40px;">return word</p> <p style="padding-left: 40px;">Go to step4</p> <p>Else:</p> <p style="padding-left: 40px;">If second and third letter of word are the fourth and the first orders of the same alphabet then :</p> <p style="padding-left: 80px;">Remove the second letter</p> <p style="padding-left: 80px;">return the stripped word</p> <p style="padding-left: 80px;">Go to step 4</p> <p>Else:</p> <p style="padding-left: 40px;">Return word</p> <p style="padding-left: 40px;">Go to step 4</p> <p>4 If end of file not reached Go to step 1</p> <p>Else :</p>
--

Stop processing

Table 4.12: type 2 reduplication stripping stemmer algorithm

4.11. IMPLEMENTATION OF THE STEMMER

The Silt'e stemmer implemented as a sequential program using python programming language. The algorithm implemented is the longest match first and the lists of affixes are checked against the word. It removes affixes iteratively until the entire affixes are removed. Table [4.13, 4.14, 4.15, and 4.16] show sample words prefix stripping; suffix stripping, letter reduplication (type 1) and letter reduplication (type 2) results of the Silt'e stemming algorithm:

No.	Prefixes to be removed	The word to be stemmed	Word after stemming
1	በ	በደውሰይ	ደውሰይ
2	የ	የሚሸትክ	ሚሸትክ
3	የ	የስልጤ	ስልጤ
4	በ	በቡርደክ	ቡርደክ
5	ተ	ተዌሴ	ዌሴ

Table 4.13: Sample of prefix striping

No.	word before stemming	suffixes to be removed	Word after stemming
	ደውሰይ	ይ	ደውሰ
	ሚሸትክ	ክ	ሚሸት
	ቡርደክ	ክ	ቡርደ
	በቶጵያ	በ	ቶጵያ
	ሙትቸ	ቸ	ሙት
	ሰሙኒሙ	ኒሙ	ሰሙ

Table 4.14: Sample of suffix striping

Reduplication to be removed	The word to be stemmed	Word after stemming
ና	ናናገ	ናገ

ሳ	ሳሳሀ	ሳሀ
ሳ	ሳሳሰ	ሳሰ
ሳ	ሳሳደ	ሳደ

Table 4.15: Sample of type1 the reduplication striping

Reduplication to be removed	The word to be stemmed	Word after stemming
ባ	ድባበለ	ደበለ
ቃ	ንቃቀለ	ነቀለ
ቻ	ንቻቹ	ነቹ

Table 4.16: Sample of type2 the reduplication striping

4.12. EXPERMENTS AND DISCUSSIONS

4.12.1. EVALUATION OF THE STEMMER

There are several criteria for judging stemmers: correctness, retrieval effectiveness, and compression performance. There are two ways in which stemming can be incorrect: over stemming and under stemming. When a term is over stemmed, too much of it is removed. Over stemming can cause unrelated terms to be conflated. Under stemming is the removal of too little of a term. Under stemming will prevent related terms from being conflated. To evaluate the performance of the stemmer, manual counting technique was used. This helps to compare number of errors that are not conflated correctly with the correct one. The stemmer was tested the sample text1 that is taken randomly from the sample text document. Sample text1 contains 1,486 words. During the stemming process two types of errors are observed. These are under stemming and over stemming. The stemed result of the sample text is 966. Out of these words 4.24% (41) words were under stemed, 10.04 % (97) words were overstemed. Totally the stemmer generates 14.28 % (138 words) stemming error. As a result, the accuracy of the stemmer becomes 85.71 %. See some errors of the stemmer in Table 4.17, Table 4.18 and table 4.19, shows the summarized result of the evaluation.

Words	Resulting stem	Expected stem	Error type
አይነትቸ	ነት	አይነት	Over stemming
አመት	አመ	አመት	Over stemming
	ቡርዳመ	ቡርዳ	Under stemming
□□□□□	ደም	አበድ	Over stemmingr

□□□	ወቅ	ወቅት	Over stemming
ብላትቸ	ብላት	ብላ	Under stemming
ሰብት	ሰብ	ሰብ	Under stemming

Table 4.17: examples of stemming errors

Table 4.17 shows sample errors observed during the stemming process of silte stemmer.

4.12.2. THE RESULTS

Table 4.18 and table 4.19, shows the summarized result of the evaluation.

In table 4.18 column name Data set contains the sample text used in the evaluation process, column words show the number of words in the sample text1, the figure 966 resulted after the sample text1 experienced the stemming process, out of 966 words correctly stemmed are 828 hence resulted in the accuracy of 85.71%.

Data set	words	stems	Correct stems	Percentage (%)
Sample text1	1486	966	828	85.71

Table 4.18: correct stem

Table 4.19: illustrates the distribution of errors for the incorrectly stemmed words in Sampletext 1.

Out of the stemmed words 14.28 % (138) words were incorrectly stemmed. Over stemming account for 10.04 % (97) words and under stemming 4.24 % (41) words.

Data set	Errors (%)	Over stemming (%)	Under stemming (%)
Sample text1	14.28	10.04	4.24

Table 4.19: distribution of errors

The errors that observed from the stemmer could be due to the following reasons:

1. Because of the complexity of the language, it was difficult to come up with complete list of affixes at a time.

2. More context sensitive rules are required based on the detailed study of the morphology of the language

WORD COMPRESSION RATIO

Sample text=1,486 words. The number words of the test data after stemming also counted and are 966.

Stemed=966

The stemmer is also evaluated with percentage of compression. For calculating the compression rate(C), the expression used by [9] was used. The expression is shown as follow:

$$C = \frac{100 * (W - S)}{W}$$

Where W number of total words

S number of stemed words from W

Sizeof the data=1,486

Number of stems=966

Hence, the percentage of compression for Silt'e text based on the training text for

Stemmer becomes $100 * (1,486 - 966 / 1,486) = 34.99\%$

The compration ratio is 34.99% for the stemmed words. It reduces the sample text by 34.99%.

SUMMARY

When the stemming process is done, over stemming and under stemming problems is observed.

For the experiment, 1486 words, which collected randomly from the available text document sources, have been used. From the data the performance of the stemmer is evaluated and the result shows 85.71 % accuracy. From this data set of words, 34.99% of compression of stemmed words found. The total errors observed from the experiment were 14.28%.

The stemmer conflates derivational and inflectional affixes. It does not conflate irregular and compound words.

The next chapter presents conclusions of findings and recommendations for future research.

CHAPTER FIVE

CONCLUSION AND RECOMMENDATION

5.1 CONCLUSIONS

Silt'e is one of the Semitic languages. These languages have common grammatical system based on the root pattern structure. The main word formation process in Silt'e is done through affixation. Silt'e uses prefix, infix and suffix. It also uses reduplication of part of a word, which is first constant reduplication and the second constant reduplication of a word.

The process of adding one affixe to another can result in relatively long words, which often contain an amount of Semantic information equivalent to a whole English phrase, clause or sentence. A Silt'e word can give rise to a very large number of variants, with a consequent need for effective conflation procedures if high recall is to be achieved in searches of Silt'e text data bases. In this research, the possibility of developing stemming algorithm for the language was investigated.

The quantitative analysis done on the sample text showed that words are distributed through the texts in their morphological variants.

Stemmers developed for other languages could not be applied for this language because of the morphological complexity and difference in features of the language as discussed in chapter three. However, commonly used methods of stemmers such as using affix dictionary, stop word list and context sensitive rules are employed. Also some techniques are adopted from Porter and Nega in developing the stemmer. Stripping suffix is not enough to conflate variant words of Silt'e to one form. Hence the stemmer also includes procedures to remove prefix and reduplication of letters.

To evaluate the stemmer developed in this study, test data of size 1486 words were selected randomly from the sample texts. The experiment showed that the stemmer performs at accuracy of 85.71 % and reduce dictionary size by is 34.99% for stems. This shows that using a stemmer for Silt'e brings a significant reduction in dictionary size as a result of conflating variant words

to the same stem. The results obtained from the experiment are promising and using the stemmer in IR system of the language could improve the performance of the system.

The stemmer conflates only inflectional and derivational affixes. It does not conflate compounding and irregular forms. Infix and prefix-suffix pair also not covered in this research. The reason is that it is possible to put prefix part to prefix list and suffix part to suffix list to strip the prefix-suffix pair of the word. As far as the infix part of the language concerned it is difficult to handle it within short period of time. There are different ordering possibilities of applying the procedures. In this experiment the stripping order was prefix, suffix and letter reduplication. Other possibilities could not be tested due to limitation in time.

5.2 RECOMMENDATIONS

Although encouraging result has been obtained in this research work, the following recommendations are identified for further work in order to make the result useful in operational retrieval environment:

This research showed the possibility of developing a stemmer to conflate word variants of Silt'e language, which has complex morphology and where a word can have very large variants. However the study was based on a limited size of sample texts and not tested in IR environment because of the time constraint. The operation of the algorithm depends on the order of the stripping procedure, despite the fact that it is not clear in what order the procedure should be applied to an input word to obtain correct stem. The ordering could be based on arbitrary criteria or on linguistic analysis. Due to limitation in time the merits of the different possible ordering could not be studied in this experiment. Studying the effect of ordering the stripping procedure on the performance of the stemmer and selecting the best possible ordering.

Evaluating the stemmer on text collection of large size collected from different sources. This is because large size sample can represent the characteristics of the language more than small size sample. Therefore, the accuracy of the stemmer can better be checked this way.

This stemming algorithm developed based on Unicode data, it is recommended to use transliterated data to make it more efficient.

Evaluating the stemmer in IR environment (if developed in the future) to measure its performance in actual retrieval session.

One can add more context-sensitive and recoding rules in order to increase the accuracy of this stemmer.

Moreover, by doing additional research on the language, this work can help to develop application tools such as spell checker, parser, thesaurus and dictionary.

REFERENCES

- [1] Gerard Salton. Automatic Text Processing: The Transformation, Analysis, and Retrieval of Information by Computer, Addison-Wesley, Reading, MA, 1989.
- [2] Christian Jacquemin, Guessing morphology from terms and corpora, In Belkin et al [1], pages 156-165.
- [3] Wessel Kraaij. Viewing stemming as recall enhancement, In Hans-Peter Frei, Donna Harman, Peter Schauble and Ross Wilkinson(editors), Proceedings of the 19th Annual International ACM-SIGIR Conference on Re-search and Development in Information Retrieval, pages 40-48, Zurich, Switzerland, 18-22 August 1996. ACM
- [4] Robert Krovetz. Viewing morphology as an inference process, In Robert Korfhage, Edie Rasmussen and Peter Willett (editors), Proceedings of the 16th Annual International ACM-SIGIR Conference on Research and Development in Information Retrieval, pages 191-202, Pittsburg, U.S.A., June 27 - July 11 1993. ACM
- [5] J. B. Lovins. Development of a stemming algorithm, Mechanical Translation and Computation, Volume 11, Number 1-2, pages 22-31, 1968.
- [6] M. Porter, An algorithm for suffix stripping Program, Volume 14, Number 3, pages 130-137, 1980.
- [7] Evelyne Tzoukermann, Judith L. Klavans and Christian Jacquemin, Effective use of natural language processing techniques for automatic conation of multi-word terms: the role of derivational morphology, part of speech tagging, and shallow parsing. In Belkin et al [1], pages 148-155.
- [8] Jinxi Xu and W.Bruce Croft, Corpus-based stemming using co-occurrence of word variants. ACM Transactions on Information Systems, Volume 16, Number 1, pages 61-81, January 1998.
- [9] Donna Harman, How effective is suffixing?, Journal of the American Society for Information Science, Volume 42, Number 1, pages 7-15, 1991.

- [10] David Hull, Stemming algorithms: a case study for detailed evaluation. Journal of the American Society for Information Science, Volume 47, Number 1, pages 70-84, 1996.
- [11] Chris D. Paice, An evaluation method for stemming algorithms, In Croft and van Rijsbergen, pages 42:50.
- [12] Lenneon, M. Peirce, D. Tarry, B. and Willett, P. "An evaluation of conflation algorithms for information retrieval." Journal of information science, 3,177-183, 1981 .
- [13] Wakshum Mekonen, "development of Stemming algorithm for Afaan Oromo Text", M.Sc. Theses, Addis abeba University, 2000.
- [14] Tesfi Tewedrios, "Word formation in Tigrinya" M.S.c Theses, Addis Abeba Unversity, 1993.
- [15] Savoy J. "Stemming of French words Based on grammatical Categories "Journal of American society for information science, 44(1) 1-9, 1993.
- [16] Ahmed F., et al, Experiments with stemming algorithm for Malay words, Journal of the American Society for Information Science, 47(12), 909-918, 1996.
- [17] Al-Kharashi, I.A. and Evens, M.W, "Comparing Words, Stems and Roots as Index Terms in an Arabic Information Retrieval Systems "Journal of the American Society for information science, 45(8), 546-560, 1994).
- [18] Popovic M. and Willett P. , "The effectiveness of Stemming for Natural Language Access to Slovin Textual Data" Journal of American society for Information sciences, 43(5), 391-395, 1992.
- [19] Ekmekcioglu C.F., Lynch M.F and Willett P. , "Stemming and N-gram Matching for term conflation in Turkish Texts", 1996.
- [20] B.L. Narayan and Sankark.Pal, Distribution Based Stemmer refinement Machine intelligence unit, Indian, Statstical, Institute, Calcutta-700108, India, Pp.1-6.
- [21] David A. Hull, Stemming Algorithms-A case study for Detailed Evaluation, Pp.1-27, 1995

- [22] S.Srinivasan and ZP.Thambidurai, Stans Algorithm for Root Word Stemming, Information Technology Jornal 5(4): 685-688, India, 2006.
- [23] Abdule Basit M., Gowende R., Husen .A, and Abduselam M., “A Hibrid Method of Stemming for Arabic Text”, Libya, Pp.1-7, 2008.
- [24] H. Harmanani, W. Keirous and S. Raheel, The International of Arab Jornal of information Technology, Vol.3,No.3, A Rule-Based Extensible Stemmer For Information Retrieval Application to Arabic, Pp.265-271, 2006.
- [25] F.W. Lancaster, Online Information System, Encyclopedia of Library and Information Science, Vol. 20, Marcel Dekker, New York, 1977.
- [26] Girma Berhe, Stemming Algorithm Development for Tigrigna Language Text Document, MSc Thesis, Addis Ababa University, 2001.
- [27] Birungi, P., Improved Strategies for Employment and Human Resources Utilization in the Information and Documentation Sector, Strategies for Human Resources Development for Information Management in Africa, Ed. 49-57, Addis Ababa: UNECA, PADIS, 1995.
- [28] Nega Alemayehu and Petter Willett, The Effectiveness of Stemming for Information Retrieval in Amharic for Information Retrieval: Electronic Library and Information Systems, Vol. 37, Num. 4, PP. 254-259, 2003.
- [29] Lemma Lessa, Development of Stemming Algorithm to Wolytta Text, M.sc Thesis, Addis Ababa University, 2000.
- [30] Wessel Kraajj, Viewing Stemming as Recall enhancement, In Hans-Peter Frei, Donna Harman, Peter Schauble and Ross Wilkinson (editors), Proceedings of the 19th Annual International ACM-SIGIR Conference on Re-search and Development in Information Retrieval, Pages 40-48, Zurich, Switzerland, 18-22 August 1996, ACM.
- [32] Amit Singhal, Modern Information Retrieval: A Brief Overview, 2001.
- [33] F.W. Lancaster, Online information system, encyclopedia of library and information science, Vol.20, Marcel Dekker, New York, 1977.

[34] Baeza-Yeates and Riberio-Neto B., Modern Information Retrieval, England, Addison Wesley Longman Limited, 1999.

[35] <http://www.omniglot.com/writing/Silt'e.htm>

[36] Ricardo Bauza-yates and Bertnier Ribeiro-Neto, 1999.

[37] Jelitan A.Hugh,E.Williams & S.M.M Tahaghoghi, Stemming Indonissian language, 28th Australian Computer Science Conference (ACSC), Conferences in Research and Practice in Information Technology, Vol.38.pp.1-8, 2005.

[38] Alemayehu, N and WillettP., Stemming of Amharic words for Amharic words for information retrieval, Literary and Linguistic computing, Vol.17, No.1, pp.1-17, 2002.

[39] Basen O.Alijla, Stemming in Natural Language, Faculty seminar, Departement of Information Technology System , Faculty of Information Technology, The Islamic Unversity of Gaza, 2009.

[40] Girma Berhe, Stemming Algorithm Development for Tigrigna Language Text Document, MSc Thesis, Addis Abeba University, 2001.

[41] Birungi P., Improved Strategies for Employment and Human Ressources Utilization in the Information and Documetation Sector, Strategies for Human Resources Development for information management in Africa, Ed. 49-57, Addis Abeba:UNECA, PADIS, 1995.

[42] Nega Alemayehu and Petter Willett, The Effectiveness of Stemming for Information Retrieval in Amharic for Information Retrieval, Electronic Library and Information Systems, Vol.37, num.4, pp.254-259, 2003.

[43] Lemma Lessa, Development of Stemming Algorithm to Wolytta Text, M.sc Thesis, Addis Abeba Univerity, 2003.

[44] Abebe Belay, Designing A Stemmer For Ge'ez Text Using Rule Based Approach, M.sc Thesis, Addis Abeba Unversity, 2010.

[45] Salton G. & McGill N.J, Introduction to Modern Information Retrieval, New York: McGraw-Hill, 1983.

- [46] Hetzron R., "The classification of Ethiopian Semitic languages", University of California, California, 1969.
- [47] Hiava M.K, Cross language Retrieval–English/Russian/French at URL <http://www.accessinn.com/aaai.htm>, 1997.
- [48] ELSE II project, Report on the feasibility of adding multilingual functionality to the search and browsing facilities of the ELSE system at URL http://nile.dmu.ac.uk/elise/el2_dels/d427d.htm, 1999.
- [49] Peters Carol, "Multilingual Information Retrieval", Presented at Cross Language Evaluation Forum, 2000.
- [50] Kaplan, a Stemming Algorithm for Latin Text Databases, 2000
- [51] Eiman Tamah Al-Shammari, Kuwait University, George Mason University, ISBN: 978-972-8924-63-8 © 2008 IADIS.
- [52] Gutt E.H.M. & Hussein Mohammed, Silt'e–Amharic–English dictionary (with a concise grammar by E-A Gutt), Addis Ababa: Addis Ababa University Press, 1997.
- [53] Gutt E.-A, "Concise grammar of Silt'e", in: Gutt, E.H.M. 1997, pp. 895–960, 1997.
- [54] Dillmann, the Ethiopic language book, vo.2, n.3, p.107-112, 1958.
- [55] Hull, Stemming Algorithms: A Study for Detailed Evaluation, Journal of the American Society for Information Science, 47(1), 70-84, 1996.
- [56] Alemayehu N., Willett P., Stemming of Amharic words for Information Retrieval Literary and Linguistic Computing 17(1), 1-17, 2002.
- [57] Alemayehu N., Willett P., The effectiveness of stemming for information retrieval in Amharic, Emerald Research Register 37(4), 254-259, 2003.
- [58] Argaw A.A., Asker L., An Amharic stemmer: Reducing words to their citation forms, In proceedings of the 45th annual meeting of the association for computational Linguistics, PP. 104-

110, ACI, Prague, Czech Republic Workshop on Computational Approaches to Semitic Languages, 2007 b.

[58] Dawson J., Suffix removal for word conflation, In Bulletin of the Association for Literary and Linguistics computing, Vol. 2(3), pp. 33-46, 1974.

[59] Debela Tesfaye, Ermias Abebe, Designing a Rule Based stemmer for Afaan Oromo Text, International journal of computational linguistics (IJCL), Volume (1): Issue (2) October 2010.

[60] Frakes W., R. Baeza-Yates, Information Retrieval: Data Structures and Algorithms Englewood Cliffs, NJ: Prentice-hall, 1992.

[70] Girma Berhe, "A stemming algorithm development for Tigrigna language text documents", Master Thesis, Faculty of Informatics, Department of Information Science, 2001.

[71] Hafer M., and S. Weiss, "Word Segmentation by Letter Successor Varieties", Information Storage and Retrieval, 10, 371-85, 1974.

[72] K. Darwish and D. Orad, CLIR experiments at Maryland for TREC 2002: Evidence combination for Arabic-English retrieval, In proceeding of TREC 2002, Gaithersburg, Maryland, 2002.

[73] Khoja S. and Garside R., Stemming Arabic text, Computing department, Lancaster University, Lancaster, 1999.

[74] Krovetz R., Viewing Morphology as an inference process, In proceedings of the 16th annual international ACM SIGIR conference on research and development in information retrieval, pp. 191-202, ACM New York, 1993.

[75] L. S. Larkey and M. E. Connell, Arabic information retrieval at UMass, In proceedings of TREC 2001, Gaithersburg: NIST, 2001.

[76] L. S. Larkey Ballesteros and M. E. Connell, Improving stemming for Arabic information retrieval: light stemming and co-occurrence analysis, In SIGIR 2002, Tampere Finland: ACM, 2002, 2002.

- [77] Lovins J., Development of a Stemming Algorithm, Mechanical translation and Computational linguistics, 11 pp. 22-31, 1968.
- [78] Mayfield J. and McNamee P., Single N-gram stemming, In proceedings of the 26th annual international ACM SIGIR conference on research and development in information retrieval, pages 415-416, Toronto, Canada, ACM Press, 2003.
- [79] Paice C., Another stemmer, In proc. Of SIGIR Forum, Vol. 24(3), pp. 56-61, 1990.
- [80] Popovic M & Willett P., Processing of documents and queries in a Slovene language free text retrieval system, Literature and linguistics computing, 5, 182-190, 1990.
- [81] Porter M, An algorithm for suffix stripping Program, 143, Pp. 130-137, 1980.
- [82] Savoy J., Stemming of French words based on grammatical categories journal of the American society for information science, Vol. 44, No 1, pp. 1-10, 1993.
- [83] Xu J., Croft B., Corpus Base stemming using concurrence of word variants, In ACM transactions on information systems, Vol. 16, No 1, Pp. 61-81, 1998.
- [84] Melucci M., Introduction to Modern Information Retrieval, New York, 2003.
- [85] Massimo Melucci and Nicola Orio, the Measurement of Term Importance in Automatic Indexing 2003.
- [86] Alemayehu, Nega and Willet P., Stemming of Amharic words for Information Retrieval, In Literary and Linguistic Computing, 17 (1): 1-17., 2002.
- [87] Paice D.C., Methods of Evaluation of stemming Algorithm Based on Error Counting, Journal of the American Society for Information System, 47(8): 632-649, 1996.
- [88] Hmeidi I., Kanaan G. and Evens M., Desighn and Implimentation of Authomatic Indexing for Information retreaval with Arabic Documents, Jornal of the American society for information Systems, 48(10):867-881, 1997.

APENDECES

APPENDIX I: List of stop words compiled for the stemmer

አብተቴ እብታይ እብተቴ አብታይ አብተቴ እሊ አሊ እነይ አነይ

እቲ አቲ እቢ ዋ አቢ ሱር አድ አድአድ እነይ በልዳሌ አልቀሬ

ደር ኤት ፈሬ ፎኖ ግነ ቀደ ስር እንኮ ኡንኮ

ማ ምን አይኔ አይነኮ ላይሽ አይታይ ወይ ታሌ ናር

ናርት ናሩ ብቾ ለገነ ገነ ሁለምክ ለሳድባድክ ለሰባድሽ

ለሰባድኑም ለሰገግሽ ለሰገግኑም ለሰገግክ ለሰገጋ ሀዳድኑም

ድባዩ አብሌ ናርት ናር ናሩ ቱፍትሉፍት ተሬር ታቼን

ቀለ በለ ቂጦ ለደር አበይ ጉት ኤጋህ ኤጋሽ ኤጋሙ

ታዮኑ ገናገናይ ዮናነይ ቢትላይም አሎነ ሂነሚ በሁነትነሙ

ሁኖተኒመዋ አታይ ሀነይ ሁኖ ገና ዮልስ ሉላሉሌ

በሉሌ አለቢ ኢንኩምንገ ቅጩ ብቾ ገገ ኡስጥ ግዝ

ሆነምታሌ በሁኖት ለላሉሌ በልዳሌ ወክት ደር አይነኮ

ሉሌ አቢሌ ተደር በውኖትምክ የገግ ገገኑ ቲታሚ አናኔ

ገጌ ሀነገነ እነይ ሱርዋ ገናሚ ሂንኩምንገ ኡሀ

ያሽ ዩዩ በዩ ዋ ለገነገናም ለሂው በሰቼ አተም ግዝቸ

ሬራቀደን ሬራ ቀደን አናግና ኢንኪምንገ ሉሉሌ ዮለ

ሱር ሂንኩምንገ ሉሀ ለሂ አዪ እታይ እተቴ አታይ አተቴ

ትዮኑ ዩታይ ዩተቴ ያታይ ያተቴ

APPENDIX II: List of Silt'e prefixes compiled for the stemmer

ኢለው እለው በል ኢለ በሰ አይ

እተ እት እል እለ አል አት አይ

አት ተይ በ ሳ ቃ ጫ አ ሰ

ለ የ ይ ት ከ ቲ ተ ሻ ና ያ እ

APPENDIX III: List of Silt'e suffixes compiled for the stemmer

ኩሞሙ ቢያኔ አታም አተኘ አሙ ሼታ

ኩሙ ኒሙ ሼሽ ታት ኤን አኘ ተኘ ንቾ

እሎ ዩን ኔት ታም ተኛ ኩ ኢ ይ ቴ ኤ ነ ሽ

ኤ ኔ ከ ሽ ሽ ቾ ይ ኮ ኩ ሁ ው ነ ን ሽ ሺ

ኤ አ ተ ት ኡ እ እ አ ኦ ው ም እ ን ኤ ኢ

ቴ ቾ ሁ ው ሀ ሽ ነ ት ሙ ሁ ው ከ ሀ

አ ሽ ሺ ኤ አ ት ነ ሙ ኡ እ ኮ

ካ ን ም ሙ እ ይ ሚ ዋ ሶ ኒ ኝ

ሎ ያ ከ ኛ ቶ ቻ

APPENDIX IV: Silt'e alphabet

ሀ	ha	ሁ	hu	ሂ	hii	ሃ	haa	ሄ	he	ሀ	hi	ሆ	ho
ለ	la	ሉ	lu	ሊ	lii	ላ	laa	ሌ	le	ለ	li	ሎ	lo
መ	ma	ሙ	mu	ሚ	mii	ማ	maa	ሜ	me	ም	mi	ሞ	mo
ረ	ra	ሩ	ru	ሪ	rii	ራ	raa	ሪ	re	ር	ri	ሮ	ro
ሰ	sa	ሱ	su	ሲ	sii	ሳ	saa	ሴ	se	ሰ	si	ሶ	so
ሸ	ša	ሹ	šu	ሺ	šii	ሻ	šaa	ሼ	še	ሸ	ši	ሻ	šo
ቀ	qa	ቁ	qu	ቂ	qii	ቃ	qaa	ቄ	qe	ቀ	qi	ቆ	qo
በ	ba	ቡ	bu	ቢ	bii	ባ	baa	ቤ	be	ብ	bi	ቦ	bo
ተ	ta	ቱ	tu	ቲ	tii	ታ	taa	ቲ	te	ት	ti	ቶ	to
ቸ	ca	ቹ	cu	ቺ	cii	ቻ	caa	ቼ	ce	ቸ	ci	ቻ	co
ነ	na	ኑ	nu	ኒ	nii	ና	naa	ኔ	ne	ን	ni	ኖ	no
ኘ	ña	ኙ	ñu	ኚ	ñii	ኛ	ña	ኜ	ñe	ኘ	ñi	ኞ	ño
አ	a	ኡ	u	ኢ	ii	ኣ	aa	ኤ	ee	አ	i	ኦ	o
ከ	ka	ኩ	ku	ኪ	kii	ካ	kaa	ኬ	ke	ከ	ki	ኮ	ko
ወ	wa	ዉ	wu	ዊ	wii	ዋ	waa	ዌ	we	ወ	wi	ዎ	wo
ዘ	za	ዙ	zu	ዚ	zii	ዛ	zaa	ዜ	ze	ዘ	zi	ዞ	zo
ዠ	ža	ዡ	žu	ዢ	žii	ዣ	žaa	ዤ	že	ዠ	ži	ዡ	žo
የ	ya	ዩ	yu	ዪ	yii	ያ	yaa	ዬ	ye	የ	yi	ዮ	yo
ደ	da	ዱ	du	ዲ	dii	ዳ	daa	ዴ	de	ደ	di	ዶ	do
ጀ	ja	ጁ	ju	ጂ	jii	ጃ	jaa	ጄ	je	ጀ	ji	ጆ	jo
ገ	ga	ገ	gu	ጊ	gii	ጋ	gaa	ጌ	ge	ገ	gi	ጎ	go
ጠ	ta	ጡ	tu	ጢ	tii	ጣ	taa	ጤ	te	ጠ	ti	ጡ	to
ጨ	ça	ጩ	çu	ጪ	çii	ጫ	çaa	ጬ	çe	ጨ	çi	ጰ	ço
አ	Pa	ኦ	Pu	ኢ	Pii	ኣ	Paa	ኤ	Pe	አ	Pi	ኦ	Po
ፈ	fa	ፉ	fu	ፊ	fii	ፋ	faa	ፌ	fe	ፈ	fi	ፎ	fo
ፐ	pa	ፑ	pu	ፒ	pii	ፓ	paa	ፔ	pe	ፐ	pi	ፑ	po

APPENDIX V: Ge'ez (Ethiopic) alphabet

ሀ	he	ሁ	hu	ሂ	hi	ሃ	ha	ሄ	hE	ህ	h	ሆ	ho
ለ	le	ሉ	lu	ሊ	li	ላ	la	ሌ	lE	ለ	l	ሎ	lo
ሐ	he	ሐ	Hu	ሐ	Hi	ሐ	Ha	ሐ	HE	ሐ	H	ሐ	Ho
መ	me	መ	mu	ሚ	mi	ማ	ma	ሚ	mE	ም	m	ሞ	mo
ሠ	se	ሠ	`su	ሚ	`si	ማ	`sa	ሚ	`sE	ሥ	`s	ሥ	`so
ረ	re	ሩ	ru	ሪ	ri	ራ	ra	ሪ	rE	ር	r	ሮ	ro
ሰ	se	ሰ	su	ሲ	si	ሳ	sa	ሲ	sE	ስ	s	ሶ	so
ሸ	xe	ሸ	xu	ሺ	xi	ሻ	xa	ሺ	xE	ሸ	x	ሻ	xo
ቀ	qe	ቁ	qu	ቂ	qi	ቃ	qa	ቁ	qE	ቅ	q	ቆ	qo
በ	be	ቡ	bu	ቢ	bi	ባ	ba	ቢ	bE	ብ	b	ቦ	bo
ቨ	ve	ቨ	vu	ቪ	vi	ቫ	va	ቪ	vE	ቭ	v	ቮ	vo
ተ	te	ቱ	tu	ቲ	ti	ታ	ta	ቱ	tE	ት	t	ቶ	to
ቸ	ce	ቸ	cu	ቹ	ci	ቻ	ca	ቸ	cE	ች	c	ቼ	co
ኀ	`he	ኁ	`hu	ኂ	`hi	ኃ	`ha	ኄ	`hE	ኅ	`h	ኆ	`ho
ነ	ne	ኑ	nu	ኒ	ni	ና	na	ኑ	nE	ነ	n	ኖ	no
ኘ	Ne	ኘ	Nu	ኚ	Ni	ኛ	Na	ኚ	NE	ኘ	N	ኞ	No
አ	a	አ	u	አ	i	አ	a	አ	E	አ	l	አ	o
ከ	ke	ከ	ku	ከ	ki	ካ	ka	ከ	kE	ክ	k	ኮ	ko
ኸ	`ke	ኸ	`ku	ኺ	`ki	ኻ	`ka	ኺ	`kE	ኸ	`k	ኻ	`ko
ወ	we	ወ	wu	ወ	wi	ወ	wa	ወ	wE	ው	w	ወ	wo
ዐ	`e	ዐ	`u	ዐ	`i	ዐ	`a	ዐ	`E	ዐ	`l	ዐ	`o
ዘ	ze	ዘ	zu	ዘ	zi	ዘ	za	ዘ	zE	ዘ	z	ዘ	zo
ዠ	Ze	ዠ	Zu	ዠ	Zi	ዠ	Za	ዠ	ZE	ዠ	Z	ዠ	Zo
የ	ye	የ	yu	የ	yi	የ	ya	የ	yE	የ	y	የ	yo
ደ	de	ደ	du	ደ	di	ደ	da	ደ	dE	ደ	d	ደ	do
ጀ	je	ጀ	ju	ጀ	ji	ጀ	ja	ጀ	jE	ጀ	j	ጀ	jo
ገ	ge	ገ	gu	ገ	gi	ገ	ga	ገ	gE	ገ	g	ገ	go
ጠ	Te	ጠ	Tu	ጠ	Ti	ጠ	Ta	ጠ	TE	ጠ	T	ጠ	To
ጨ	Ce	ጨ	Cu	ጨ	Ci	ጨ	Ca	ጨ	CE	ጨ	C	ጨ	Co
አ	Pe	አ	Pu	አ	Pi	አ	Pa	አ	PE	አ	P	አ	Po
አ	Se	አ	Su	አ	Si	አ	Sa	አ	SE	አ	S	አ	So
ፀ	`Se	ፀ		ፀ	`Si	ፀ	`Sa	ፀ	`SE	ፀ	`S	ፀ	`So
ፈ	fe		`S	ፈ	fi	ፈ	fa	ፈ	fE	ፈ	f	ፈ	fo
ፐ	pe	ሁ	fu	ፐ	pi	ፐ	pa	ፐ	pE	ፐ	p	ፐ	po

ሰ	lWa	ገ	gWu	::
ሐ	HWa	ቀ	qWi	::
ጣ	mWa	ኀ	hWi	፤
ሢ	`sWa	ኁ	kWi	፤
ረ	rWa	ገ	gWi	፥
ሰ	sWa	ቋ	qWa	፦
ሻ	xWa	ኂ	hWa	፥
ቄ	qWe	ኃ	kWa	፥
ቧ	bWa	ገ	gWa	
ደ	dWa	ቋ	qWE	
ጿ	jWa	ኄ	hWE	
ገ	gWe	ኅ	kWE	
ጣ	TWa	ገ	gWE	
ጫ	CWa	ኆ	ea	
ጸ	PWa			
ጺ	SWa			
ፋ	fWa			
ፐ	pWa			
ቀ	qWu			
ኀ	hWu			
ኁ	kWu			



Figure 1. The basic Feedel (Ethiopic base characters and their six forms)

APPENDIX VI: Comparison of Silt'e words before stem and after stem

Silt'e words before stem

ቀደምጅታሽ ሸኚ ሺፕት ቀደምጅታሽ አያሜ አያሜ አዮለሽ በቡርደሽ ያለትወደ ያለትኮ ማ ኤወደ ለጦራሚ ኡርበ ለሳንኩራይ ኢለውነቅ የሀደኡመት ኪሼክ በሉሉሌ ስንፌ ዩውዲበያን የቀውልን ሱብሱብ ለፋፍነት ለገረድነት በሚሽ ሚሽ ገፈራኔ በሚሽሚሽ መዬ በሬደ የሙሃበ አይጦፋን ሼጣኔ ሰመይ ወዘና ወልዴ አያም አተረ በሀቅል መርከተኛ አልበለይ ኢሼማን በጨልመ የጠለሀ አስለመጠ መረጨክ ኪዝበኛን አራሺን የቀበጠ አዋለኩ ወክቡ ተቀቡ ጩጠ ይልብ አዋልክ ሀድ ሰብ በቀውል አናጊን በኩትብ ለያቲለልፍ ኢክሸነይ ሉክተ ፈያክ አሻኔ ሉሃ የትሬ ቁጨ ለገና ሰብ ኢትሪናኮዋ ኢቻለኮ ያሸቢያንን ዋልክ አይነት ስነ ዱትብ ጡሪ በለይ አዋለክ በቀውል አናፊን በኩትብ ኢቀርባን ግዝ አዋልክ ሊግኖት ለሰብ ወልድ ቢሃል ቀይማን ዋ ቢያቀር በይማን ፍደልቸ ኡስቤክ አናፊና ሉክታክ ለቲላል ፎት በኩትብ ኢትጎበለን የቀውልን ሱብሱብት ኩትብ በስንፌ አናጊን በቡትን ሊቀርብ ያቀትለን ኩትብ ወገሺን ተገደገደ ቲለልፋን አሊመ አበ አዴነሽ ኪቲሰል ወራቤ በሂዶት ደር ታለት ቴድበያቲ መኪና ታብለሻኔ ጉራይ ኢንግሽ የትሜለ የተረሻት መሀመድ ሞተ ሚሽተክ ከዲጀም ሞት ቢሉ መሰው ሆሽት ሰብቻ ለሚነ ሞቱ ለሚነ ሊሞቱ አቀተሉ ኢለንሱል ሊነቀ ያቀትለን ኩትብ ለቀረበይ ሱል ኪምባዩ ለዩባናክ ኢለጦፍ ሎናምክ መሀመድ መኪነ ታግለበጠቢያኔ በምቶትክ ሚሽትክ ከዲፋ የሚሽ ሙቶት ሰማታኔ ቱሃ ተለዬው ኢለውነብር ባታኔ ኢቀትለን ደዌ ሰቼት ሞት ኢታይ ኩትብ መንቁዋ ውጣት ሌንዘክ ለቀሪሎይ ያጦፋን ቢቢ መሳወ መሀመድ የሞተቢመንቁክ የከሚናን አምሰ የሚሽትክ መውት መንቁክ ደዌዋ የመሀመድ ተመውት የስልጤ አፈርክ ኢቤዣን ለላቶ ዘማን በቴለቀ አሽር በጤማይ ባድ አርክ በታለቀ በደውሰይ በድ ኮሌጅ መኒ ባሉነ ወራቤ ቢሌድ አጋኔ አይቺሌን ግዝ ኢነካሌ አማድ ኡመተው በድሁን አለሃን አሀምድ ለታንጨ ለኖታም ጫፍ ሲር ቲያበሊ አለ ሚነ ኢላቴር ሰበው በርደባሊ ሃኪመ ጋር ለርከቦት ራቁ ቲትጉላሊ ተውኖት ፈመረ ቦራቤ ሆስፒታሊ አብሺሪስ የኘቴ ወራቤ ከተማ ሰሙይ በቁበበይ ኢዝነሽ ተሳሴማ ሙሶ ቀኚቢሻን ሱማሽ ሉልም ጫም የሲልጤን ወልድ ሆናን ወራቤን ያለንገፍ አደም ጊዝ ታያቀር አለሃው ዋበሽ አሮት ማልተ ታይቢል በጃደ ወልደሽ፣ ቶቢያ ባዴያ ጌሰም ቲሌቁሽ ኪመ አበቃሎ ዮባነ አፈረሻይ በሂም ሉለሉሌ ኤንዘን ገዎሻይ ከለፌ አለፈዱይ በለንኒ እመሻይ ያፍሪቀቴን ወሪ ቶቢ በድናይ ለዚጠኚ ወሪ የለበስኩይ ጂስማሽ ሆሽተ ዛን ዱነክይ የጠወረኝ ጨነሽ አሮሮ ቲትቢይኝ የራመኚ አንጀሽ አለሃን ኢጦቅሳያው ለያቴሪኝ ቡቻሽ ዱማነሽ ሹሩበ አቅናኔ ታቱሻን

ቢሴቹ በቅባት ፈያክ ሚሰሱሻን የሰመይ አመረ ሙግ በለ ያንገርሻን ኢታይ ታክማመቶት ላሽው ሚነ ያኒንሽን ሲልጢነታሽ መጢክ ኢዝኒክ ኢከታሻን አዝጋግ በሎት በቡርደነ ሀንጊነ ባቢሌ ይትረከቦን ሀለትቸ ዩንዜ ሀደ ግዝ ሆኖኔ ሰብ ሀንጊነ ሩሀ ያለይ ገነ ጊዝ ሀንጊነ ያአመኝ ሰብ ሩክቦ ያሺያነን ጊዝ ዩልቀተኝ አዝጋግቸ የደች ቃምቸ ለባይትከ ቆቶ ወንደበን ሚዴ ገነገናም ዩንዜ በሂምቸዋ አበቃላሎነ ቡንዱላሌ ዩንዜን ሰበኝ አዝጋግቸንገ ያቀኑይሙ አዝጋግቸነ ቢሾ ተዮን ያመኝ አዝጋግቸነ ለባይትከ የአሀላቅ ሀለትቸነ ሲያሰኝዋ ኢኮኖመኝ ሀለትቸዋ የአደኝ አዝጋግቸነ ለባይትከ ሰብአኝሎ አቆምስቸ አርት ዲን ሰብጊሾ ሳይንስኝ ቴክኖሎጂኝዋ ላቁምሶት ዩድገሎን ጊዝቸ ዩንዛን ያዝጋግ ሳይንስ ባድለኝ ሳይክለ በፈትለኝ አሽር መቃም ያደ አሽረኮ ለደረሳሶ ቲቀርብ ለደር የጩቀሙይሙይ ጊዝቸነ ባደኝ ቡንቁፌ ዩንዛነኮ አቻበሩያን ለሆነም ኤት አሽረይ ለላሉሌ አሽርቸነ አቻበራኔ ኤንዛን ለባይትከ የጅስመ ሳይንስ የሃያት ሳይንስ ያድኖት ንባረት ሳይንስ ያመኝ ሳይንስ አፊየ ሂርሰት የንባረት ብላት አሽርዋ ዩንጅ ብልቸነ አቻበራኔ ዩንዛን ያዝጋግ ሳይንስ በጢሽት ያቀራነይዋ ጉተኝ ሂንገጮከ በሰብያዮ አዝጋግ ሃለትቸ ደር ይትሰበታነይ ሰበዬይ ቲሌቅ ቲዩድ ተቡርደክ ነቀላኔ ቲፌት ዩዳነይ አዝጋገከ ፊልፊሎት ይጀሚራን ሂታይ ሊቆት አዝጋገከ የቻሎት ፊልፊሎት ሃለት ሄደ ሄዳኔ ቱሀ አቢሌ ታለይ ፈቲይ አዱኒያ ተራከቦትከ ለሰቢዬይ መሊቅዋ ተጀረቦት በሉላሉሌ ኡንገ ይድገለያን ሂ ሰቢዬይ ተቡርደክ ጀመራኔ ለቻሎትዋ ገዝ ለትጀረቦት ያሼያነይ ጅሀድ ሰቢዬይ በቡርደክ ለቻሎት ያለይ ቆመፌቸ አቁሚሮት ደበላኝ ገዝ ተፈቲይ አመኝ አዝጋገከ ለቻብሮት ያድገለያን ሰቢያዮ ቢታይ ሃለት በሰብቸ ጉት ያለይ ሩክቦዋ ሃድነት ለትጀረቦት ተድግሎትምከ ደበላኔ በባድቴ ያሉይ ለላሉሌ አደቸዋ ሰብቸነ የንበሮት ሃለትቸ ለትጀረቦትዋ አደኒሙዋ የንበሮት ሃለትኒሙ ላሂብደት መሰ ዮናን ሂነሚ ያትሬጊጣን በድባያም የአዝጋግ ሳይንስ አሽር በሰወደም የሁን ባመኝ አዝጋገነ ያክማምቶትዋ የትቃቂሮት ወሻይብነት ይነበረናነኮ ያሻን ያዝጋግ ሳይንስ ሰብያዮነ ገዝኒሙ ለትጀረቦት ይድገላን ኢልምዋ በሰበ ጉት ፊየ ዮነ ርክቦ ያሶነኮ ይግዘይማን ትዮን ፊየ አሽላቅም ሊነበረይሙ መሳን አሽረይ የሰወደም የሁን ያመኝ ፊየ አሀላቅ ባትሪክቦት ሰብያዮነ ሳይንሰኝ የሱል ብላትቸዋ ዩምርዋ የሃያት የቻሎት ኢልምቸ ለርከቦት አሰነት በውኖት ቆማሬ ዮነ ፈይሰለ ባቦት

ያውጄም የሁን የከሸፈ ወቅተኒሙ በሱተ ይድገሎነኮ መሰ በውኖት ፈይደ ዮብይማን አሽረይሚ ሰቢያዮነ ለባድኒሙ ይድጋሎን ጠረኝ ሰብ ያሺይማን

ያዝጋግ ሳይንስ ያሽር ንድፌ ጎልጌ ጃንጌ ቲትሲናድ ባጂሲይ ያሽር ፎሊሲ ለቶጲ ያሽር ጋርቸ ቢቤሻርን ሃለትዋ ለቀደ በትረጅይ ሁንዱሉሌ ዞፎፎ አሰነትን የንድፌ አሽረይ በቂጦ ኡምር ይትረከቦን ሰቢያዮዋ በቂጦ ጎልጌ ላሉ ደረሳሰነ ኡስቤ ባሶትን ኢትቁራነይ የቻሎት መቃምዋ ኢንገይይ ጎልጊቸ ቢከተይማን የቻሎት መቃም ሃለትን ተስናጃን ሂርሰት ቴክኖሎጂዋ ኤዲሰ በትረጅይ ሁንዱሉሌ ዞፎፎ አሰነት ቢቤዣንዋ ቢድገሎን ቂጨ በፊቲክ ሃደኘ ይቻበሮነኮ ባሽሪይ ውስጥ አገባይማን ንድፌ አሽሪ ለደረሳሰይ ለአቅልኒሙ ቢያብያይናን ሃለት ያለ ግዝ ባውጦት ለደረሳሰይ ቢትፋሃመኒማን ሃለት አጣጣይን

ሀጂሲይ ያሽር ከሪኩለም ለቶጲያ ያሽር ጋርቸ ቲትሲናድ ሰአብት ሲጀ ወቅቲቸነ ባድ ሳምት ኡስጥ ለአዝጋግ ሳይንስ አሽር ሳደያን ያድ ያሽር ዘማን ሳምቲቸ ሲጀ ወቅትቸ ወሰዳኔ ለትሲናጂ ኤት አሽርጌታቶይ ፈት ወቅት ይነብረይማነኮ ባይትም ዲባያዮ ሳምትቸ አትረከቢኒይማን ቢጴታሚ ወቅት አሽርጌታቶይ ፈት ኩፍተ ሲጀ ወቅትቸ ሊረከበንኤት ሚካት ያለቢይሙዋ ሉሌ ኡግሻር ያትኬጅይማን ደረሳሰነ ሊገዙት ያቀትሎነኮ መሰ ዮነይማን ኩፍት ወቅት ኒበሮትከ ደረሳሰይ ሊገዙትዋ ገነ ዲባያ የአሽረ ቢል ላሶት መሰ ተሁኖትምከ ደበላኔ የአሽረይ አይዶ በሱት በድጋለሎት የትጎበለይ ቀስድነ በሱተ ለከሚሎት ይድገላን

ለሆነም ኤት አሽርጌታቶይ ኩፍተይ ወቅት አሽረይ ቀደሞኔም ቢፊዱ በማን ማነምግዝ ታቲልፎት ለደረስቻይ ኡግዣር ቢያኛን ሃለት ተድጋለሎት ይነብረይማን ሃጂሲይ ፎሊሲ የአሽረይ ሃለት ቲዬውድ የቅሮት አቅሮት ሃለቲይ ዩለጊነ ዞፎፎነ ባዴኘ ባቻበረ ሃለት ኢንዞት አለቢይን ይላን ሂነሚ ለከውኖት የቀውል የክተቦትዋ የከውን ሃለት ተድጋለሎት መጥን ለሆነም ኤት ሂ የአሽር ንድፌ አሽርጌታቸ ዩለጊነ ዞፎፎነ በጎልጌ አሽር ቲያቀሩ ይትዲገለሎነኮ ኡስቤ ያሻን ቲዮን ኡሀነም የጎልጌ ቢል ባፍዋ በኩትብ ጃድ የጋረ ቢል ባሶት ዋ ያቢሌ የከውነ ብል ባሶትዋ ዞፎፎ ባሶት በሰመጢክ ደረሳይ ቢል ያሻነኮ ባኮትዋ ባዙፎት ገነገናም ከውኖትን

ቢታይ ንድፌ አሽር በሱራይ በሰወዶ መትፈጄከ ኤት ላዙፎት ይድገሎን ቀውልቸግ ፈሀሚቸ ያሉ ቲዮን አሁንም አሽርጌታቶነ ሙላይ ያሽረይ ቀስድ በሱተ የትከመለዋ ያልቲከመለኮ ላጢሮት ይቤሻርይማን ቢያዣህሮን ወቅት ኢከታን የቻሎት መቃም ቂጨ” ይላነይ ኡስቤ ገነ ኤት ውሰዶት ኤለቢይ ለበሎት የከሲይዋ ተጀረቦት ያለቢነይ ሱተኘይ ቢልዋ በሱተይ መቃም ሀድ ደረሰ ተጀረቦት ያለቢይዋ ለኢታይ ላለቢይ ጎልጌ ቂጨ ዮነ ከሚሎት ያለቢይ ኢልምን ለበሎትን ቲታይ ቂጨ ከሎም የሁን ደረ አዘረክ ያሉ ደረሳሰነ አሽርጌታቶይ ልቤ በሎት አለቢይሙ በደርክ ያሉይ ያቆምሮነኮ አንሻሽጦት ያለቢነ ቲዮን በኮሎ አዘርክ ያሉይ ጋኔቶ ሚካተኒሙ ላሌኔ

ኡግዛር አሶት ያትኬሻን ሂይ ቀስድ በሱተ ለከሚሎት አሽርጌታቶ ደረሳሶ ያሽርጋረይ አሚርቸዋ የአመኛ አዝጋግ አበሮስ ሁሉምከ በቂጠ ተቻበሮኔ ጃድ አሶት ያትኬሻን ሂይ ያሽር ንድፌ ጎልጌ ጃንጎ በሉሌ የአዝጋግ ሳይንሰ ያሽር ንድፌ የትሴናዱይ ያጥቃቀሎይ የሙለይ በቶጵያ ያሽር ሚኒስቴር ያሽር ንድፌ አስናዶትዋ አቆምሮት ጎልጌ ቲዮን በፈረንጅይ ሂልቅ ያሰናዱይ ለደር የቀረቡይ ባይትቻነ አይነኮ በነቶ ለይቴንዚ ያግዘነኮ ባበይኖት ወባጃይ ጀምሮት ለሰአደዲይ ነታቶ ደረሳሶይ ሚነ ተረሶት ያለቢይዋ ተቀጠብቻነ ዩሰቦነኮ አሶት የሳንበ ነቶ በአየር ይትብራጫነይ፣ አይነኮ ይትቃጠብነያን? አየርንገ አግቦት አውጦትነ አቃኖት ኢላቀትሊ፣ ጊንመ የጋረነ አየር መሰከተከ በክፈቶት ጠሊል አሶት ያቀትሌን የደል ሙርጠጤ ሹርት ይሰራጫነይ ተዌሴ ሹማትን አይነከ ይትቅራቀርነያን በሸማኒ ዞፍ በሱተ እንደ ተራጦት አለቢነ በኒጥር ኢንጅ ሲንቅ ያቀረበ ሰብ ስንቀከ ብሎት ኤለቢና ሪጀ ይትሲራጫነይ በሪጀ ጩንጩ ቲነከሲን አይነኮ ይትቅራቀርነያን በሪጀ አጎበርዋ የቃነነ መይ ባፊሶትዋ በድሪቆት ኤድሰ በደምን ይትሲራጫነይ አይነኮ ይትቅራቀርነያን ደመነ ተገነሰብ ደመኒካኮት ኤለቢነ ተዚነዋ ሱጠ ዮኑ ሙትቸነ ደመ በነኩዋ በቂጠ ተድጋለሎት ቢደጎት የጅስመ ቆመፌቻቾ ሚነ ዮኑኮዋ ሱመኒሙ ጨቅሞት ደረሳሶቻይ የሰአደዲይ የጅስመ ቆመፌቻቾነ ይጨቅሞኔ ያቴሮነኮ ይጠሩይማነኮቢለኒይሙ ያበይኖነኮ አሶት ደረሳሶይ በገገኒሙ ተሮሻት በነቀ ይሸሉይማን ዮኑ የኑርቡቃቆነ ይጠሮነኮ አሶት ሉላሉሌ ዮኑ ለባይትከ ምስለ በትዲጋለሎት ሃቀኛ ሙትቸ ሻመ ካሌ የኮረንቴ ኑር ሙቃድ ይነዳን ኢንጩ አይርዋ ወሪ ባቲሮት ደረሳሶይ ሉላሉሌ ኑርቸን አቲሮት ሎነኮ ተሳሎት የነታቶይ ሂልቀኒሙ በድሎት ደረሳሶይ አይታይ ሱር በቡርደኒሙ ሁለጊነ ይትረከቦንወ በለ አይትረከቦን በብሎት ይለሎነኮ አሶት ደረሳሶቻይ ኤድሰነ ቲገናይ በጠቀስ ፈየኮ ለይችሎም ያቀትሎን ነቶይ ቀታይዋ ደዊ ዩለይ ሆኖተኮ ኢውዶት ለኢምኮ በቡርደኒሙ ተገናይ ቲያትዋሩይ በብሻ የልቲረከበይ በበሎት ለደረሳሶይ ለሰብ ወልድ ሲንቀ ላሊቆት ለቁሚሮትዋ ቁወ ላቦት ያድገለያነኮ አበይኖት ደረሳሶቻይ ይትረዘቆነይ ቶኑይ ሃዳላን ሴንቂቸ ኡስቀ ሱመኒሙ ይጠሮነኮ አሶት ዮዱያነይ ሲንቀ ሱመም ይጠሮነኮ አሶት

Silt'e words after stem (stop word exclu

ሸኚ ሺፕ ቀደ ያሜ ዮላ ቡርደ ማ ኤወደ ኡርበ ውነቅ ሀደ ኪሼ ለሉሌ ስንፌ ዩውዲበ ቀውል ሱብሱብ ሚሽ ገፈራ መዩ ሬደ ሙሃበ ጠፋ ሹጣ None ወዘና ወልዴ ያም ተረ ሀቅል መር ለይ ኢሼማ ጨል ጠለ ስለመጠ መረጫ ኪዝበ ራሺ ቀበጠ ዋለ ወክቡ ጩጠ ዋል ሀድ ሰብ ጊን ኩትብ ፈያ ሻኔ ኡሃ ትሬ ቂጨ ኢትሪና ኢቻለ ሸቢ ነት ስነ ዱትብ ጡሪ ፊን ኢቀርባ ግዝ ግኖ ወልድ ቢሃል ቀይማ ዋ ቢያቀር ፍደል ኡስቤ ፊና ኡክታ ፎት ኢትጎበለ ቡት ቀርብ ቀትለ ወገ ሊመ አበ ዴነ ኪቲሰል ወራቤ ሂዶ ደር ታለ ቱድበያቲ መኪና ታብለሻ ኢንግ ሚለ መሀመድ ሞተ ዲጀ ሞት ቢሉ መሰ ሆሽ ሞቱ ቀተሉ ንሱል ነቀ ኩትቤ ሱል ኪምባዩ ጠፍ ሎና መኪ ታግለበጠ ምቶ ዲፋ ሙቶ ሰማታ ቱሃ ውነብር ታኔ ኢቀትለ ደዌ ሰቼ መንቄ ውጣ ሌንዘ ቢቢ መሳወ ሞተቢመንቄ ሚና ምሰ መው ደዊ ስልጤ ፈር ኢቤዣ ዘማ ቱለቀ ሸር አድ ርክ ታለቀ በድ ኮሌጅ መኒ ለነ ቢሌድ ጋኔ ቺሌ ኢነካሌ ማድ ድሁ ሀምድ ጫፍ ሲር በሊ አለ ሚነ ኢላቴር ሰበ ርደባሊ ሃኪ ጋር ራቄ ጉላሊ ፈመረ ቦራቤ ሆስፒታሊ ብሺሪስ ኘቴ ተማ ኢዝ ሙሶ ቀኚቢሻ ሱማ ኡል ጫም ሲልጤ ሆና ደም ጊዝ ታያቀር ዋበ ሮት ታይቢል ጃደ ቶቢ ዴያ ጌሰ ሌቄ ኪመ ቃሎ የወ ሂም ኡለሉሌ ኤንዘ ለፌ ፍሪቀ ወሪ ጂስማ ዛን ጠወረ ጨነ አሮሮ ቢይ ራመኚ ንጀ ኢጠቅሳ ቡቻ ዱማ ሹሩበ ቅና ታቱሻ ቢሴቼ ቅባ ሚሰሱሻ መረ ሙግ በለ ንገርሻ ታክማ ሸው ኘን ሲልጤነታ መጠ ኢክታሻ ዝጋግ ሎት ሀነጊ ቢሌ ሀለ ዩንዜ ሆኖ ሩሀ ገነ መኚ ሩክቦ ሺያ ዩልቀተኚ ደች ምቸ ቆቶ ወንደበ ሚዴ ገነገና ሎነ ቡንዱላሌ ሰበኚ ዝጋግቸንገ ቢቾ ሀላቅ ሲያሰኚ ኢኮኖመኚ ደኚ ሰብአኚ ቆምስ ርት ዲን ሰብጊቾ ሳይንሰኚ ቱክኖሎጂኚ ቁም ዩድገ ዩንዛ ሳይንስ ድለኚ ሳይክለ ፈትለኚ መቃ ያደ ሸረ ቡንቁፌ ቻበሩ ኤት ኡላሉሌ ሸር ቻበራ ኤንዛ ጅስ ሃያ ድኖ ንባረ ፊየ ሂርሰ ብላ ዩንጅ ብል ጢሽ ጉተኚ ሂንገፍ ብያዮ ሃለ ሌቅ ዩድ ነቀላ ፌት ዝጋገ ፈልፌ ቆት ቻሎ ሂደ ሂዳ ቱሀ ዱኒ መሊቅ ለላሉሌ ኡንገ ሂ ጀመራ ገገ ጅሀድ ቆመፌ ቁሚሮ ደበላኚ ድገለ ሰቢያዮ ብቸ ጉት ሃድ ደበላ ባድ ለይ ደቸ ንበሮ ደ ሂብዶ ሂነ ሬጊጣ ድባ ወዶ ሁን ክማ ቂሮ ወሻይብ ሻን ሰብያዮ ኢል በ ፈየ ርክቦ ሶነ ዮን ሸላቅ መሳ ሰወዶ ሪክቦ ሳይንሰኚ ዩምር ሰነ ውኖ ቆማሬ ፈይሰለ አቦንድፌ ሸር ወዶ መትፈጀ ኤት ዙፎ None ቀውልቸግ ፈህ ያሉ ዮን ሁን ሸርጌታ ቀስድ ሱተ መለ ልቲከመለ ጢሮ

ቢያገህሮ ወቅ ኢከታ ቻሎ መቃ ኡስቤ ገነ ውሰዶ ሲይ ቢል ሀድ ደረሰ ጎልጌ ቁጨ የወ ሚሎ ኢል ታይ ሎም ደረ ዘረ ደረሳ ልቤ ሎት ደር ሉይ ቆምሮ ንሻሽጦ ኮሎ ዘር ጋኔ ሌኔ ኡግሃር ሶት ኬሻ ሚር መኚ ዝጋግ ሮስ ሁሉ ቂጦ ጃድ ጃንጎ ሉሌ ሳይንሰ ሸር ቶጵ ሚኒስቴር ስናዶ ሂልቅ ቻነ ነኮ ነቶ ግዘ ጀምሮ ነታ ሚነ ዩሰቦ ሳንቦ የር ርንጎ ግቦ ውጦ ኖት ጊን ጋረ ክፈ ጠሊል ቀትሌ ደል ሙርጠጤ ሹር ሹማ ነክ ሸማ ዞፍ ንደ ኒጥር ኢንጅ ሲንቅ ቀረቦ ሰብ ስንቀ ብሎ ኤለቢና ሪጀ ጩንጩ ነከሲ ጎበር ነነ መይ ፊሶ ድሪቆ ኤድሰ ደም ደመ ኤለቢ ሱጠ ሙት ነኩ ቂጠ ቢደጎ ጅስ ቆመፌቻቾ ወዋ ሱመ ጨቅሞ ቴሮ ገገ ነቀ ኑርቡቃቆ ኡላሉሌ ምስለ ዲጋለ ሃቀኚ ሻመ ካሌ ኮረ ኑር ሙቃድ ኢንጨ ርዋ ወሪ ሮት ሎነ ሂልቀ ድሎ ሱር ቡርደ ሁለጊ በለ ረከቦ ጠቀስ ፊየ ቀት ደዊ ኢውዶ ብሻር ወልድ ሲንቀ ሊቆ ቁወ ቦት ድገለ ሃዳላ ሴንቂ ኡስቀ

Appendix VII words with lower frequencies but are stop words

Word	frequency	Word	frequency	Word	frequency	Word	frequency
ገናኒሚ	2	ልቶን	3	ኢነ	2	ጊዘቸ	2
ጃንጉ	3	ቦንት	2	ኢንኩ	2	ሂንኩ	2
የሰሀደደይ	2	ያፍቴቴ	2	ኢነኩ	2	ዩሃን	2
ለገጊኒሞ	2	እቴታን	1	ባታይ	3	በገጊሽ	2
አደደሞ	2	እነይ	2	ገገገነ	2		
ሂተይ	2	እታይንገ	2	ሎኔት	2	ለሆነ	2
ጊንመ	2	ሊሊ	2	ሂንኩምከ	2	ሂንኩም	5
አደኒሙ	2	ዮናነይ	2	ታት	2	ዩኑ	2
ሊኢታይ	2	ገገሽ	3	ለሄት	2	ዮናን	6
ቴቲይ	2	ቦነንገ	2	ሚነ	3	ለሆነምከ	2
ኩሃም	2	ወናኔ	2	ዩኑኮ	3	ያፍቴ	3
ተደር	2	ብቼ	2	ሂንኩሞ	3	ናር	2
ሬረ	3	የብቸ	1	ታይታይ	4	ተዮን	5
ሆኖት	2	በዮን	2	ዩነኮ	2	ሊታይ	3
በገጋሙ	2	እሊ	2	ለሚን	2	በምን	5
በዮንም	2	የኛነይ	2	ቢቻ	1	ሎንት	4
ሂንካሌይ	2	ባሎትን	9	ሂ	3	ቢታይ	13
ገናኒሚ	2	አዲ	6	ሂነይ	10	ሁኖተኒሙ	3
አነጊነ	7	የገናይ	6	ለደር	9	በገና	3
ዩልይ	2	ሂንኩምንገ	4	ታዮኑ	2	በርሬ	3
ዮነይ	4	የሰአደደነ	2	ሊትኬተሎነይ	9	ገገረሙ	6
ለታት	8	ገገከ	5	ዩነይ	4	ቴኢኮሎ	3
አይኔ	2	ገነገናሚ	2	ኢታይንገ	2	ሁልም	6
እታይ	2	ተገናይ	10	ገናም	3	አዘር	14
የሄተ	2	ኩሁኑም	6	የገገከ	7	ሁሎም	3
ሀኔ	2	የናረይ	8	ለገናሚ	3	በሁኖተኒሙ	3
ለኢኮ	6	ኢነይ	9	ገነገናምኒሙ	1	ታቲ	3
ሆነም	4	ሂሊኮ	2	ገገ	4	በገገኒሙቀ	3
ታሌ	4	ሎነም	9	አይነት	16	የገናም	3
አልዳሌ	3	ኢታይ	12	ግን	3	ቀደ	6
አዮናነይ	2	ጊነ	6	ኩስቤ	26	አቱም	3

አይታኒ	8	በገነ	7	ለሳድባድ	2	ዮነን	4
የገነ	2	ሂታሚ	4	ገና	2	በሁኖት	3
ዞፍ	9	ሱር	9	ያበደም	2	የሁን	5
በጉትኒሙ	6	ጉተ	2	በአበድም	2	ሁሉምግነ	2
በባይት	2	ሩክቦ	22	ለምን	8	ቲኢ	5
ገናይ	9	ግዝቸ	11	ገነ	9	ላፍቲ	3
በአድ	3	ሁለም	5	ለሰሀአደዲ	2	ኮሎ	16
ሂታይ	5	ባፍቲ	2	ምን	12	ዮኑ	13
ለሰሀደዲ	3	ገነይንገ	2	አይነኮ	13	ኡሁንም	2
ቶህ	3	ታቲም	3	ጉት	16	ጉተኝ	2
የቶህ	3	ሁኖተከ	11	ለሰደዲ	2	ሀለት	14
ግነ	13	በልዳሌ	12	ሀደደ	4	በጉደር	2
ሬሬሰ	5	ቢኢታይ	3	በገግ	6	ለባይትምከ	7
አሶት	16	ኢንኩምንገ	7	የገግ	5	ያሌተ	2
ቢኢሸ	2	አድባድ	3	ለገናይ	6	ሁለምግነ	2
የሳድባድ	4	ለሆነምኮ	3	ምነ	13	በጉተኝ	2
ለሰደደኒሙ	2	በውኖት	7	አይታይ	5		

Declaration

I declare that the thesis is my original work and has not been presented for a degree in any other University.

This thesis has been submitted for examination with my approval as university advisor.
