

# Amharic Text Summarization for News Items Posted on Social Media

Abaynew Guadie\*, Debela Tesfaye, Teferi Kebebew

Department of Computing, Jimma Institute of Technology, Jimma University, Jimma, Ethiopia

## ABSTRACT

This paper introduces Amharic Text Summarization for News Items posted on Social Media, to summarize the news items posted Amharic texts a time posted documents from social media on Twitter and Facebook; The main problems of the social media posted texts are that most people would probably read they're posted in Amharic texts with duplicate posted documents. However, to find the information the user is looking for to find summary posted texts and read important portions of posts as Amharic documents to extract desired information on social media. Summarization is dealing with information overload presenting and posting with a text document for the current time representation of the posted documents to summarize. Our proposed approach has three main components: First, calculate the similarity between each posted document within the two pairs of sentences. Second, clustering based on the similarity results of the documents to group them by using Kmeans algorithm. Third, summarizing the clustered posted document individually using TF-IDF algorithms that involve finding statistical ways for the frequent terms to rank the documents. We applied the summarization technique is an extractive summarization approach that is assigned an extract the sentences with the highest-ranked sentences in the posted documents to form the summaries and the size of the summary can be identified by the user. In experiment one, the highest F-measure score is 87.07% for extraction rate at 30%, in the clustered group of protests posts. In the second experiment, the highest F-measure score is 84% for extraction rate at 30%, in droughts post groups. In the third experiment, the highest F-measure score is 91.37% for extraction rate at 30%, in the sports post groups and also the fourth experiments the highest F-measure score is 93.52% for extraction rate at 30% to generate the summary post texts. If the system to generate the size of the summary is increased, the extraction rate also increased in posted texts. For this, the evaluation system has shown that very good results to summarize the posted texts on social media.

**Keywords:** Amharic; Similarity measure; Text summarization; Clustering; Tf-IDF algorithm; Social media; Facebook; Twitter; News posted texts

## INTRODUCTION

Amharic language is one of the main African languages and also it is the working language of the Federal Government of Ethiopia and is widely spoken throughout the country [1,2]. It is an Afro-Asiatic language belonging to the Semitic group of its unique alphabets [3]. Its speaker post Amharic texts on Twitter and Facebook in the social network at this time, which is increasing volumes of posting data are available on the social media, which is observed on the growing online posts, websites, and digital storage in the language. These Amharic text documents are available digitally and the amount is highly increasing every day from a user posting to the social media on Facebook and Twitter. The task of developing for easy retrieval of relevant information is especially challenging on Amharic texts because there are only a few recent and uncoordinated efforts of automation and language processing [2]. Currently, with the

enhancement in the most people to use and posted and reposted many Amharic text documents on social networks, the amount of data one has to deal with has increased rapidly on Twitter and Facebook for the readers. Text summarization aims to investigate the summarization of user posts over time for the documents from the social media within the period of each post as a stream posted documents on the social media. Therefore, documents shall be iterated over in news items order for the news items posted on Facebook and also its aims to create and evaluate the news items of posted texts summarization systems for the news post on the social media. Services such as Twitter and Facebook generate to rapidly access phenomenal volume of content for most real-world posts on a daily source [4].

In this paper, we address summarizing targeted user posts of interest in a human reader by extracting the most representative poets of the irrelevant Amharic tweet stream for the post news that could be

**Correspondence to:** Abaynew Guadie, Department of Computing, Jimma Institute of Technology, Jimma University, Jimma, Ethiopia, E-mail: abaynew.guadie@ju.edu.et

**Received:** September 27, 2021, **Accepted:** October 11, 2021, **Published:** October 18, 2021

**Citation:** Guadie A, Tesfaye D, Kebebew T (2021) Amharic Text Summarization for News Items Posted on Social Media. J Inform Tech Softw Eng. S5: 001.

**Copyright:** © 2021 Guadie A, et al. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

summarized without duplicate posted Amharic text document [5]. News posted Amharic text documents include summarizing posted political protests, natural disasters, bombing, earthquake shooting, storms and social media for accident things to happen to compile information from a large set of documents into single news posted document or generate a summary of user posts from Twitter and Facebook [6]. A good way to get up-to-date, monthly, and yearly information Amharic texts will be to get a stream of sentence length in posted documents about the situation as it develops [7].

### Text summarization for news items tasks

The work of Eidheim proposes the temporal summarization task, which is to inform readers of important novel information about a particular topic [8]. For that, the DUC and TAC Update Summarization tasks were designed that as a single pass collection process, processing all the documents at once, while in that year TREC temporal summarization tracks the task designed requires the generation of continuous and immediate updates. As with the earlier work, sentences are the unit of selection. The author experimented with several different methods of sentence selections based on language models, TF-IDF, and word counting. One of the primary issues the authors wanted to address was the fact that earlier methods were based on the assumption that all the evaluated documents were relevant. The problems covered by this paper were very similar to the temporal summarization task and it shows that great care is necessary to be taken when selecting documents to ensure that the results are good. The author had been done only English news articles were considered when generating updates and the main reason for this is that the chunks and queries are all English, allowing other languages is a potential source of noise, but it is unlikely to contain any interesting information. It is also likely that a user monitoring an update stream would like the stream to be in language, making it a reasonable decision to only look at documents in a specific language. For the sake of keeping the collection of documents small, non-English documents are removed by the preprocessor of the tweets [9].

### Temporal patterns text documents

The work of those Premlatha and Geetha, for extracting the temporal patterns that are from texts requires the expression for handling the types of temporal term categories that were explicitly, implicitly, and vaguely conveyed to the temporal information for the given document [10]. The authors performed by using Finite State Automata (FSA) in the natural language expression are converted into a calendar-based timeline. The authors give the temporal expressions the first the temporal explicitly are examples like on the 18<sup>th</sup> of March 2009, in Nov 2010, 12 Jul 2011, etc. We applied to use this temporal in the news items of posting term, text summarization explicit expression to identify the posted documents by using date, month and year format to group by monthly for the summary posted document in social media posts.

Finally, the vague temporal expressions are like July, after several weeks before, etc. In the work of this article, for the topic-based clustering for the calendar model which goes to enable the temporal intervals posted documents was carried out their work for the papers [11]. Those authors to preprocess the use of the clustering such as stop word removal, stemming were performed

the document represented using tf-idf representation which is used for clustering for the similarity.

### Text summarization for Amharic texts

Tamiru, had described and proposed to this papers for Amharic texts the two generic texts summarization approaches [12]. For the first technique, to put the topics Latent Semantic Analysis (LSA) and the second technique mixes the graph-based ranking among the Latent Semantic Analysis algorithms to identify the topics of a document were used to select the semantically the main sentences for the summary generation. The author used the algorithm was Latent Semantic Analysis and graph-based ranking algorithms to explore this work. This work evaluates to propose the performance of the summarization approaches and prototype Amharic news text summarization system. For this author, to evaluate the summaries systems with manual summaries were generated by six independent human evaluators to be taken the evaluator for the experiments. Author prepared the dataset corpus used for evaluating the summarization system was 50 Amharic news from Ethiopian news reporters the items were in the range of 17-44 sentences. These results are evaluated by comparing the system summaries with their corresponding manual summaries.

The work of Teklewold, had proposed open sources of customizing by Amharic texts automatic summarization using an open text summarizer tool of two ways of executing the two experiments, the first one experiment is done without changing the code of the tool and the second is done for the changing the Porter stemmed tool for the Amharic stemmed algorithm [13]. This work uses the frequency of terms to determine the relative importance of a sentence in a text. This paper's evaluation of the experiments was producing 90 news articles and to test its performance for the summaries for each rate at 10%, 20%, and 30% extraction rates for the results. This author evaluate the system was evaluated using subjective and objective evaluation.

According to the work of Yirdaw, had proposed topic-based Amharic text summarization to investigate the six algorithms to explore as the use of terms by concept matrix to implement for thesis [14]. In this algorithms take two common steps, the first step, to identify the keywords for the documents that were used to select the term of use of the concept matrix to find the document. The second step, for the sentences to find the best keywords contained that were selected for presentation in the summary. For this author to take the experiment with news articles in Amharic texts to explore the algorithms for selecting the first sentence of the document for inclusion in the summary of the Amharic news texts. In this paper author evaluated the paper by using the precision/recall for summaries of 20%, 25%, and 30% extraction rates to evaluate for the news articles. The author of comparing system with the previous methods of developed for other languages based on topic modeling approaches summarization that had been used in this Amharic data set of the papers. Since the work of these authors had investigated for Amharic texts to an automatic single document summarization tasks for the graph-based automatic Amharic text summarizer were proposed. They were to work the generic and domain independent graph based model could successfully make extracts from Amharic texts for their papers [15]. These authors use the two graphs-based ranking algorithms were introduced, their thesis used for PageRank and HITS that we're using the two-sentence centrality measures

and sum the relation between the sentences in a text for a graph that show the results from their experiments. They worked to prepare the data sets for the experiments shown 30 news articles used on economics, politics, society, and sports were conducted for Amharic news articles from collected for Ethiopian reporter news web sites and Addis Admas to test its performance for the summaries for amharic news articles [15].

### Single and multi-document text summarization for posting texts

In multi-document summarization system is developed for the web context. Information may have to be extracted from many different articles and pieced together to form a comprehensive and coherent summary [16]. One major difference between single document summarization and multi-document summarization is the potential redundancy that comes from using many source texts. The solution presented is based on clustering the important sentences picked out from the various source texts and using only a representative sentence from each cluster [16]. Multi-document summary is when a summary of one topic is prepared for many different documents. For these reasons, multi-document summarization is much less developed than its single-document and various methods have been proposed to identify cross-document overlaps from different researchers. Summons (which is a paper issued by informing a person that a complaint has been filed against it), or order, a system that covers most aspects of multi-document tweet summarization, takes an information retrieval approach [17]. Assuming that all input documents are parsed into templates (whose standardization makes comparison easier), summons clusters the templates according to their contents, and then apply rules to extract tweets items of major imports. In contrast, the problem of organizing information on multi-document summarization so that the generated summary is coherent has received relatively little attention [18]. Multi-document summarization poses interesting challenges to single post documents and also the information overloads faced by today's society pose great challenges to researchers that want to find a relevant piece of information [19]. Automatic summarization is a field of computational linguistics that can help humans to deal with this information overload by automatically extracting the idea of documents. An important study shows that even newspaper article type and some simple procedures can provide essentially perfect results [20]. According to the work of these author Inouyes and Kalita their paper described that an algorithm for summarizing microblog documents on Twitter [21]. Firstly, they presented algorithms that produce single-document summaries, but later extend them to produce summaries containing multiple documents and also they evaluated the generated summaries by comparing them to both manually produced summaries and, for the multiple post summaries, to the summary results of some of the leading traditional summarization systems. The author discussed to overflowing with information on Twitter, they taught that just being able to search for text tweets and receive user tweet the most recent posts whose text match the keywords are not enough. Summarization can represent the tweeted document with a short piece of text covering the main topics for the user posts, and help users select through the Internet, the most relevant document, and filter out redundant information on the social media to find the relevant tweets [22-24]. So, the author describes the document summarization has become one of the most important research

topics in the natural language processing and information retrieval communities.

### For the summarize the user posts on Twitter and Facebook for evaluation techniques for previous studies

The most existing evaluations of summarization systems are fundamentally typical; the evaluators create a set of summaries, one for each test text, and then compare the Summarizer's output, measuring the content overlap that often by sentence recall and precision. To simplify evaluating extracts, independently developed an automated method to create extracts corresponding to abstracts for Twitter [25,26]. The two evaluation methods of the text summarization as intrinsic and extrinsic evaluation methods. An intrinsic method is the evaluation of the system for summaries according to have the evaluators take the ratio of some scale of the interval the extraction rate (readability, Informativeness, facility, coverage and others). It was prepared by creating the human or ideal summaries of the given input text document for comparing the summary of the summarizing system and the human summary of the evaluators. The evaluator was used for measuring the average scores found from each evaluation criterion for the given document. The other evaluation method is extrinsic evaluation for measuring the acceptability and also the efficiencies of the automatic summaries of documents to achieve the tasks and easy motivate, this evaluation method is the main difficulty to ensure that are applied to correlate or link with the task of the performance efficiency for the summary documents. Some of the intrinsic evaluation methods to evaluate the summaries are as follows.

1. **Coherence and structured for summary:** This is one of the evaluators to measure the extraction rate-based methods to use the flow of the information to be structured by using the cut and paste processes of the documents on phrases, sentences, or paragraphs that produces the result in a summary extracted for the documents which results from incoherence to get for the documents.
2. **Informativeness for summary:** This is one of the summaries of comparing the system summary generated the documents within the input texts for summarizing the key ideas about manual summary that are included in the automatic summaries and the summary information are presented in the input texts.
3. **Sentence precision and recall:** It is a standard measure for information retrieval in terms of sentences in the given document. So, the precision of measures by using how many of the sentences in the system summary generated and also in the reference human summary are producing the results. On the other hand, recall for the sentences to measure for how many of the sentences in the reference ideal summary are extracted from the system summary of the documents [27]. The Twitter observatory that allows observing, searching, analyzing, and presenting social media is introduced as a part of the research, and illustrative examples of using this proposed pipeline show how the Twitter user interaction with the social media data [28]. According to this author Sharifi, et al. [29]. Described that this algorithms process collections of short posts as specific topics on social media the well-known site called Twitter and create short summaries of those posts as the Twitter. The goal of the research is to produce summaries. In this paper evaluated the summaries produced by the summarizing algorithms, compare them with

human-produced summaries, and obtain excellent results. Since this author P. Meladianos to deal with the task of sub-posts detection in evolving twitter posts using posts collected from the Twitter stream and also by representing a sequence of successive tweets in a short time interval as a weighted graph of words, they were able to identify the key moments sub-posts that combine a post using the concept of graph degeneracy on the Twitter [30]. They have selected a tweet to best describe each sub-post using a simple yet effective heuristic (which is used for experimental).

4. **The co-selection summary evaluations:** This method is that to discuss here are the simplest of all evaluation measures co-selection sentences and this simplicity creates a value score to check the intersection with the system and human summary of sentences [14]. For that, the main problem is the difficulty in accounting for variations in what humans consider ideal summary sentences in the documents. The summary evaluation for using the co-selection measure was to take from information retrieval evaluations techniques and for describing the formula for Recall (R), Precision (P), and F-measure given to calculate the system and the human selected the fraction of sentences that the system has chosen from the total of sentences found in the ideal summary as follows to formula (Equations 1-3).

$$R = \frac{|system\ and\ human\ choice\ overlap|}{|sentences\ chosen\ by\ human|} \quad (1)$$

Precision (P) measures the fraction of system summaries that are correctly chosen.

$$P = \frac{|System\ and\ human\ choice\ overlap|}{|Sentences\ chosen\ by\ system|} \quad (2)$$

F-score (F) is the harmonic mean of recall and precision.

$$F = \frac{2 * P * R}{P + R} \quad (3)$$

Many works were being done on the area of summarization on social media by Twitter in western world language for news items summarization. But most of the work is being done for major technology languages like English, Chinese, German, and French [16,31].

Due to those authors that described the shortness of tweets and also TwitIE makes the assumption that each tweet is written in only one language [32]. The choice of languages used for categorization is specified through the arrangement file, complete as an initialization parameter. The authors take three tweets one English, one German, and one French. TwitIE TextCat was used to allocate automatically the language feature of the tweet text (denoted by the Tweet annotation). Those give a collection of tweets in a new language, it is possible to train TwitIE Text Cat to support that new language as well, and also that is done by using the fingerprint generation (which is patterned to generate) included in the language identification plugin. It builds a corpus of documents and reliable tweet language identification allows them to only process those tweets written in English with the TwitIE English and named entity recognizer [32].

### The use of social media in the summarization posted documents

Social media was defined as a group of Internet-based applications

that build on the ideological and technological foundations of Web 2.0, and that allow the creation and exchange of user-posted generated content [5]. The use of social media has exploded in recent years and the high availability of such information is then used by many researchers. In this author's work, the author described as to design novel features "identifying valuable information on Twitter during natural disasters" that can be used as input to machine learning classifiers to automatically and accurately identify informational tweets from the rest in a timely fashion [33]. This approaches used for machine learning algorithms that would be discussed for URL extraction are used extensively in tweets to link to external sources that could not ordinarily be acceptable to the long-restricted structure of a tweet. However, URLs found in tweets, are shortened to accommodate for the length restriction. Inherently, a few features can be extracted from the URL itself, considering that each shortened URL has a base domain and a randomly generated code appended to that (i.e. https: "t.co/[code]"), which, when clicked on, will redirect to an actual web page [33].

### Sentence extraction by TF-IDF and position weighting

In this article to describe the sentence extraction via using tf-idf, they would be discussed the posted Japanese Newspaper to create a summary and this system is implemented with the sentence extraction approach and weighting strategy to mine from several documents [34]. They created an experimental system for the Japanese Summarization to compute the importance value of each sentence based on Japanese newspaper terms. In this paper author used the important sentences whose sum of characters exceeds the restricted character amount are eliminated and the remaining sentences are then sorted as they appeared in the original document. The author of summarization used single and multi-document summarization to implement the different evaluation between a long summary and its summary is more remarkable. In this article, a frequent term texts summarization algorithm is designed and implemented in Java and a designed algorithm is implemented using open source technologies like Java, Porter stemmer, etc., and verified over the standard text mining corpus. Japanese summarization to compute the importance values for each sentence based on Japanese newspaper terms [35].

### Hybrid TF-IDF documents

TF-IDF stands for term frequency-inverse document frequency, is a numerical statistic that reflects how important a word is to a document in a collection or corpus is occurring, it is the most common weighting method used to describe documents in the Vector Space Model (VSM), particularly on IR problems. In a hybrid TF-IDF algorithm development and the idea of the algorithm is to assign each word, sentence within a document a weight that reflects the words, sentences the most important to the document. The sentences are ordered by their weight from which the top sentences with the most weight are chosen as the summary. To avoid redundancy or duplicate posted words the algorithm selects sentences and tokenize, the next sentences or terms and checks them to make sure that it does not have a similarity with a given threshold with any of the other previously selected because the topmost weighted tweets may be very similar [36]. The tf-idf of term  $t$  in document  $d$  is calculated as (Equation 4):



$$tf - idf_{t,d} = tf_{t,d} \times idf_t \quad (4)$$

TF-IDF Algorithm (tf-idf) is a mathematical statistic that is meant to show how important a word is to a document in a collection of documents. In information retrieval systems, it is used as a weighting factor. As the number of times, a word appears in the document increases, the tf-IDF value increases proportionally. But this tf-idf value is decreased by the frequency of the word with the collection. This helps to take into account the fact that some words appeared more frequently in general. It is a logarithmic obtained value and it is obtained by dividing the total number of documents by the number of documents containing the term, and then taking the logarithm of that quotient and the log of this term is calculated to a value obtained is the IDF (Equation 5-6).

$$idf(t, D) = \log \frac{N}{|\{d \in D : t \in d\}|} \quad (5)$$

Then tf-idf is calculated as (Equation 6):

$$tfidf(t, d, D) = tf(t, d) \cdot idf(t, D) \quad (6)$$

The Inverse Document Frequency (IDF) is a measure of how much information on the documents provides, that is, whether the term is common or rare (some words such as common stop words are frequent that words do not help discriminate between one document over another) across all documents. Term Frequency-Inverse Document Frequency (TF-IDF) is a statistical weighting technique that has been applied to many types of information retrieval problems. For example, it has been used for automatic indexing, query matching of documents, and automated summarization [5,37]. One other important contribution to the Hybrid TF-IDF equation is its normalization of words of a document occurs, which allows it to carefully control the overall target summary length [38]. Had established TF-IDF is very sensitive to English text document length and often overweights terms from longer documents for capital and small characters. In our application of TF-IDF, we observed not the same effect in Amharic language as English words. Without a normalization word, established TF-IDF, given the most weight to the longer documents since the weight of a document is the simple sum of the weights of the composing words. The given below is the research study of the related literature about these approaches and systems for multi-document summarization.

**Similarity measure:** The cosine similarity measure is one of the common techniques used to measure similarity between a pair of sentences vectors. Here sentences are represented as a weighted vector between the sentences in the posted documents on social media from Twitter and Facebook.

**Cluster-based method:** Basically, clustering is to group similar sentences score values of their classes. For clustering of multi documents, these documents refer to the sentences and the cluster that a sentence belongs to be represented by classes of the group.

**Word frequency:** The basic idea of using word frequency is that important words are found many times in the document. Tf and idf are some of the most common measures used to calculate the word frequency. After this, we calculated the combine the tf and idf to score the importance of the word to find in the documents to the summary.

**Feature-based method:** The extractive summarization types include

identifying the relevant sentences from the text and putting them together to create an accurate summary. Some of the features that are considered for the selection of sentences are the significant word location of a sentence, length of the sentence from the input documents, etc.

## METHODOLOGY

### Data sets

We used to collect training dataset on the Twitter and Facebook posted documents to Amharic text input corpus 4951 posted sentences in totals (protests (3943), droughts (667), floods (101), and sports (240)) in different news items. For those in the collected posted documents from Facebook and Twitter to obtain the sources in protests in totals 120,862 posts, in droughts posted is 43,774 posts, in sports posted is 10,299 posts and in floods is 1,209 are the training sets posts are identified by the format of date, month and years. We have the following training data sets and we could be tested the data sets are #protests, #droughts, #sports, and #floods data on social media posted texts that could be found in the news items posted Amharic text documents over a time in social media. The data set consists of 4 news items posted texts to collect data onto the Facebook and Twitter posted text documents and selected to extract from summary the important sentences in the input to the summary (Table 1).

In the above (Table 1), we made for training and testing the posted Amharic texts for social media on Facebook and Twitter posted in Amharic texts to take and analyze the samples to process the similarity and cluster for the similar posted sentences automatically. Based on these data to find the similarity with each posted Amharic text with a pair of sentences on social media and identify the group of the similar clusters, the numbers of input k are three (created on distance values or nearest distance value to group the similar items of sentences to cluster and after cluster, we had been found three clustered documents for each posted text documents for sentences to summarize individually clustered (Figure 1).

### System design for the research study

**Pre-processing of input post document:** This would work on the Amharic posted documents for processing by the rest of the system. The preprocessing involves sentence segmentation in the document, tokenization of the segmented sentences, stop word removal of the list of words, stemming words

**Similarity measure:** Sentence similarity was computed as a linear combination of sentence similarity and word similarity. The cosine similarity measure is one of the common techniques used to measure similarity between a pair of sentences to be calculated the similarity between each sentence and after similarity to group the similar clusters to form.

**Clustering:** For clustering of multi documents, these items refer to the posted documents for sentences and the cluster that a similar sentence belongs to represented by classes of the group in a shorter or nearest distance value calculated by the centroids using the formula of Euclidean distance function. Once we used with post text data, k-means clustering can provide a great way to organize the thousands-to-millions of words, phrases, or sentences for posted documents. Since, the sentences for randomly for each

Table 1: Preparation of data sets pre-processes.

Datasets items	Size (in sentences)	Similarity posts (in sentences)	Clusters (in documents), K=3 groups, sentences		Amharic, posted texts (dd-mm-yyyy)format
Protests	3,943	5,34,994	c1(doc1)	201486	43,236
			c2(doc2)	166797	39,805
			c3(doc3)	166711	37,821
			Total	534994	1,20,862
Droughts	667	224971	c1(doc1)	84,380	16,850
			c2(doc2)	70,429	12,979
			c3(doc3)	70,162	13,945
			Total	2,24,971	43,774
Sports	240	28,824	c1(doc1)	10,804	3,470
			c2(doc2)	9,010	3,631
			c3(doc3)	9,010	3,198
			Total	28,824	10,299
Floods	101	4960	c1(doc1)	1,833	459
			c2(doc2)	1,554	419
			c3(doc3)	1,573	331
Totals	4,951	7,93,748	Total	4,960	1,209

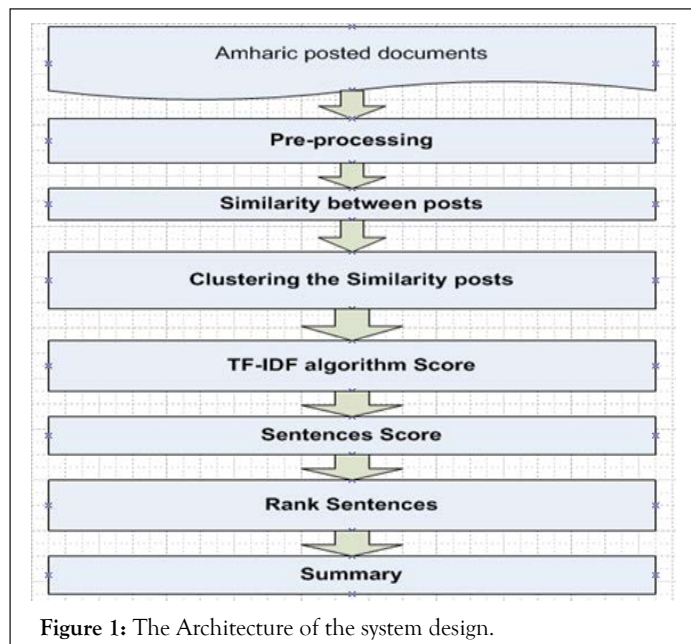


Figure 1: The Architecture of the system design.

document set for the input similarity results, posted documents in the sentences and give clusters based on similarity values of news items posted Amharic texts on social media and each set clusters the similar posted documents for the user posts on the social media [29]. When we summarize the user posts by clustered document for each post of clusters by similar distance values Amharic text documents for the social media on the Twitter and Facebook posted data. Day to day user posts to arrange with the same cluster like political #Protests with #Protests similar contents, #Droughts with #Droughts, #Floods with #floods, #sports with # sports, #health with #health the same posted documents or sentences, and others user posts to arrange or order the tweets for the algorithm and summarize for each user posts. In our research, the data onto social media could not be identified automatically protests at droughts, protests at sports posted, or other data sets. So, we could have a different corpus for data sets to train and test the performance of Amharic text documents posted on social media.

**TF-IDF algorithm score:** The TF-IDF value was composed of two primary parts. The Term Frequency (TF) and Inverse Document Frequency (IDF). TF component assigns more weight to words that occur frequently to a document. Because important words are often repeated in one document and a single document that encompasses all the documents together [4]. Therefore, TF-IDF gives the most weight towards that occurs most frequently within a small number of documents and the least weight to terms that occur infrequently within the majority of the documents (Equation 7-9).

$$Tf(t, d) = \frac{\text{Term counts for the stem}}{\text{Total terms in the count}} \quad (7)$$

$$IDF(t, Docs) = \log_{10} \left( \frac{\text{Docs}}{t \text{ appeared in all Docs}} \right) \quad (8)$$

$$Tfidf(t, d, Docs) = tf(t, d) \times idf(t, Docs) \quad (9)$$

Where t is term or word, d documents, Docs=all documents in the corpus.

**Sentence scoring:** This step after the tf-idf stem word score, each term of obtaining the score its values if this sentence determines the score of each sentence and several possibilities exist. The score can also be made to the number of sentences in which the words in the sentence appear in the document. Scoring each sentence is to rank each sentence we need to score each sentence using the tf-idf values calculated before. Rather than simply taking the sum of all the values of a given word with one sentence sequentially. These include only summing TF-IDF value score of the term where the word is a noun, verb, and others, but stop words could not be added for the sentences that could be scaled down the stop words in a sentence (Equation 10).

$$\text{Score}(q, d) = \sum_{t \in q} tf - idf_{t,d} \quad (10)$$

**Sentence ranking:** After the sentences score for the list of stem words, the sentences will be ranked according to the sentence scores values and any other measures like the position of a sentence in the document can be used to control the ranking. After applying the sentences scored we can finally sort our sentences in descending (top value to low value) orderly the sentences score values to sort. For example, even though the scores are high, we would be putting first ranks sentences and comparing to each score value of the sentence's higher score values to order the sentences Step 7.

**Summarizing extraction:** After ranking the sentences list, if the user to selects the input on the size of the summary, the sentences will be picked from the ranked lists orderly. The news items text summarization is the process of automatically creating a compressed version of a given text that provides useful information about the user in social media. The length of a summary depends on the user's needs.

## RESULTS OF EVALUATION

### Similarity measures between sentences

Sentences are made up of words, so it is reasonable to represent a sentence using the words with the sentence. The most relevant research area, to our task, is the Amharic text summarization news items posted on social media. The cosine similarity measure is used to measure similarity between a pair of sentences and also after similarity to group the similar clusters to form. Given two pairs of sentences:

$$S1=\{w11, w12 \dots w1m1\} S2=\{w21, w22 \dots w2m2\}$$

Where  $wij$  is the  $j$ th word of  $Si$  ( $i=1, 2$ ),  $mi$  is the number of words in  $Si$ . A joint-word set  $S=S1 \cup S2$  is then formed into distinct words with  $S1$  and  $S2$ .

$S=S1 \cup S2=\{w1, w2 \dots Wm\}$ , that has  $m$  distinct words. The joint word sets  $S$  contains all distinct words of  $S1$  and  $S2$ . And also we could be intersected the similarity in words with the two sentences word similarity (Equation 11).

$$S = \begin{pmatrix} S_{11} & S_{1,2} \dots & S_{1,N} \\ \vdots & \ddots & \vdots \\ S_{p,1} & S_{p,2} \dots & S_{p,N} \end{pmatrix}; S_{\mu,j} = \begin{cases} \omega_{\mu,j} & \text{When the word } j \in \bar{s}_{\mu} \\ 0 & \text{otherwise} \end{cases} \quad (11)$$

Where each  $\mu$  row contains the  $\omega_{\mu,j}$  weighting of the word  $j$  in a sentence.

Example:

$S1=9\text{-Feb-2016}$  Ethiopians protest in Washington DC #OromoProtests #Ethiopia Protests #Oromo #Ethiopia "Stop the abuse of citizens in Oromia, Gondar and Gambella!"

$S2=2\text{-Feb-2016}$  Ethiopians Demonstrate in Washington DC "Stop Violence Against Citizens in Oromia, Gondar, and Gambella" etc.

Two sentences similarity is calculated as a formula:  $S1=\text{Sentence1}$ ,  $S2=\text{Sentence2}$ ,  $\cap$ =Intersection,  $\cup$ =Union  $\text{Sim}=\text{Similarity}$  between two sentences (Equation 12).

$$\text{Sim}(S1, S2) = \frac{2 * \text{match}(s1 \cap s2)}{\text{len}(s1 \cup s2)} \quad (12)$$

The syntactic similarity with the two strings of capturing the similarity between words was concerned to work [39]. The work of the author to compute the similarity between two sentences and basically to capture the semantic similarity between the two-word senses for the strings of the path length similarity.

### Clustering-based algorithms using k-means

Clustering is to group similar posted sentences into their classes. It is a process of creating groups of similar items or posted documents. Clustering to find clusters of data posts that are similar in some nearest distance values of one another. The members of a cluster are more like each other than they are like members of other clusters. The main aims at clustering algorithms are similar to one another within the same cluster and dissimilar to the objects in other clusters. The Kmeans data sets to access and Microsoft Excel have been used for initial preprocessing and storing the data, and in particular to analyze the data sets.

#### K-means algorithm

**Input:** input the number of cluster  $k$  with a centroid (mean) Process:

**Step 1:** partition the data into  $k$ -cluster or  $k$  nonempty subsets

**Step 2:** compute the mean for each partition

**Step 3:** assign each object of the cluster of its nearest centroid (mean)

**Step 4:** Step 2 and Step 3 are continuous until no change in the mean values and also the sum of the sequence of distance is minimized within the clusters

**Output:** number of clusters of partitioned data objects (Figure 2).

**Hybrid TF-IDF algorithm calculates:** The TF-IDF is one of the techniques that are to pick the most frequently occurring terms (words with high term frequency or  $tf$ ). However, the most frequent word is less useful since some words like the stop word occur very frequently in all documents. Hence, we also want a measure of how unique a word is i.e. how infrequently the word occurs to all documents (Inverse Document Frequency or IDF). Hence, the product of TF-IDF of a word gives a product of how frequently, this word is in the document multiplied by how unique the word is written the total corpus of documents.

**Term Frequency (FT):** Which measures how frequently a term occurs in a document. Since every document is different in length, it is possible that a term would appear much more time for the longer document than shorter ones.  $TF(t, d) = (\text{Number of times term } t \text{ appears in a document}) / (\text{Total number of terms in the document})$ .

#### Algorithm for term frequency:

Read the list of term and frequently key values in HashMap

Iterate  $itr=tf$  for the Keyset to iterate for each

While ( $itr$  to hash next)

String  $term=itr$  next to string

Find  $tf \rightarrow$  to calculate the number of times term  $t$  appears in a document count over the total number of terms in the document  $tfr$  to put in term and double values

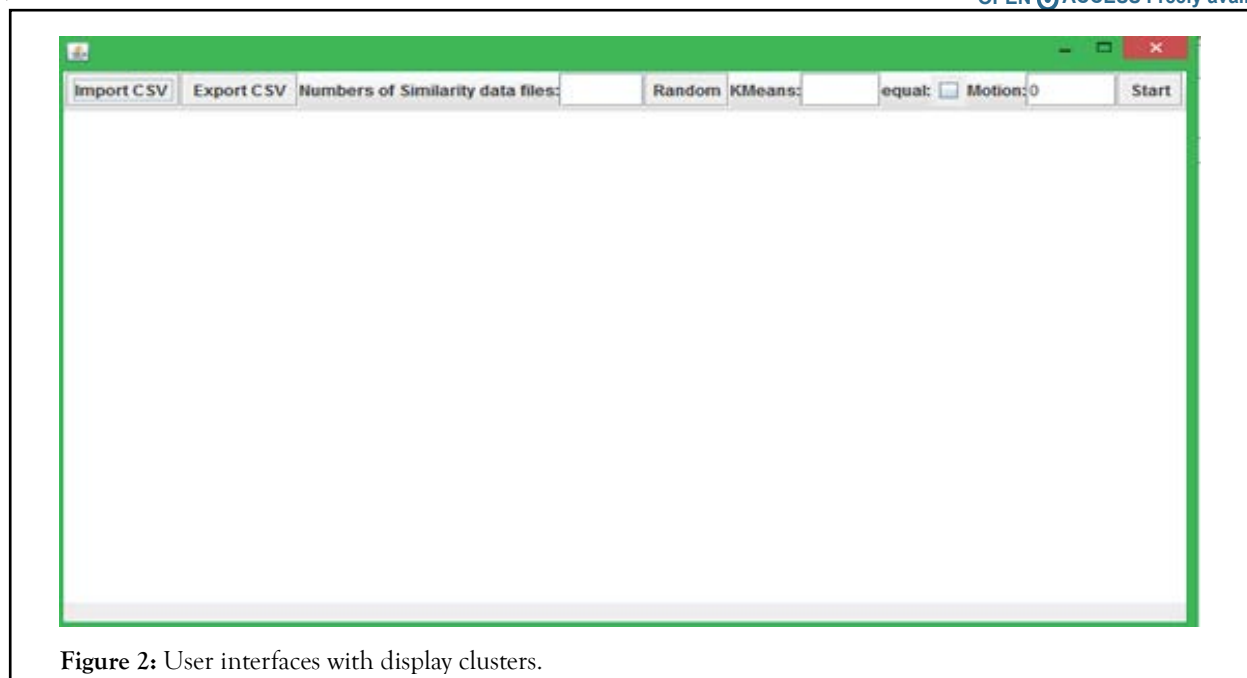


Figure 2: User interfaces with display clusters.

Display the output terms and tfr values results

End While

End

The term frequency calculates  $tf(t, d) = \text{terms count} / \text{total of the terms of document1 words}$ . Example the term ("Opposition") to the word that occurs in document one is 1060 and the total document contain words 4,802,014, to calculate tiff as follows. So,  $tf(\text{"Opposition"}, d) = 1060 / 4,802,014 = 2.2074071420866327E-4$

**Inverse Document Frequency (IDF):** The inverse document frequency is the number of times a word occurs to a corpus of documents. IDF, which measures how important a term is while computing TF, all terms are considered equally important. However, it is known that certain terms, such as "Was", "are", "and", "this", "that", may appear a lot of times but have little importance. Thus we need to weigh down the frequent terms while scaling up the rare ones, by computing the following:  $IDF(t, Docs) = \log_{10}(\text{Total number of documents of the corpus} / \text{Number of documents with term } t \text{ in it appears on all documents})$ .

#### Algorithm for inverse document frequency

Find if idf string and double values found in the hash map

Iterate  $ir=idf$  for Keyset to iterate

While ( $itr$  to hash next)

String  $term=ir$  next to string

Double  $idf=\text{math.log10}(\text{total documents in the corpus divided by the term } t \text{ appears in all documents})$

Idfr1 to put in term and idf results

Display output in the form of term and idf values

End While

End

The inverse document frequency calculated the term of the protest documents to appear,  $IDF(\text{term}, docs) = \log_{10}(\text{docs} / \text{term appears in all docs})$ . For example the term "Opposition" to calculate its

edge of all documents to appear the term within three documents to find and we had been three documents in the corpus.  $IDF(\text{"Opposition"}, 3) = 1 + \log_{10}(3/3) = 1 + 0 = 1$ . The IDF value is zero means if the term occurs to all documents in the corpus occurred or frequently appear in all documents in the corpus, so this term is the usefulness of the documents the IDF to scale down the term values. TF-IDF (Term Frequency-Inverse Document Frequency): It is a way to score the importance of words (or "terms") in a document based on how frequently they appear across multiple documents. The TF-IDF is used to weigh words, according to how important they are. Words that are used frequently in many documents will have a lower weighting while infrequent one term will have a higher weighting. Intuitively, if a word appears frequently in a document, it's important. It gives the word a high score. But if a word appears in many documents, it is not a unique identifier. It gives that the word a low score in the documents. Therefore, common words like "And", "about", "or", which appear in many documents, will be scaled down.

#### Algorithm for term frequency-inverse document frequency

Read the results of the tfr1 from the hash map the term frequency

Read the results of the idfr1 from a hash map in the inverse document frequency

Iterate  $ir=idf$  for keyset to iterator

While ( $itr$  has next)

String  $term=itrs$  next to string

Double  $idf=tfr1$  get the term

Double  $idf=idfr1$  get the term

Calculate the  $tf-idf=tf*idf$

Tfidf1 to put the term and its Tf-idf values for each

Display the output the term and tf-idf1 values

End

The tf-idf value score, example the term ("Opposition") to the word occur in the document the multiplication of tf and idf. Tf-



$IDF(term,d)=tf(term,d)*idf(term,docs)=2.2074071420866327E-4*1.0=2.2074071420866327E-4$

**News items posted text summaries**

We implement the news items posted texts in the current time to post texts in social media like protests, droughts, sports, floods, and others to post in the social media on Facebook, Twitter, and others social media, but we focused on the post texts on Facebook and Twitter post-Amharic texts to a summary by using date, month and years to summaries the poster texts to give to the readers within a short period to summaries the posted documents. News items posted to identified by the time of the intervals of the period posted texts in the social networks/media like to use the date, weeks, months, years, one year's, etc. posted texts to summarize the texts. For example, if the user selected for political protests summary to need within the monthly to the group for multi documents within a one month and to want to summarize by entering the month and year into the system to find automatically the system group posted text document based on the user asked the month to group the same month and years rank each posted sentences for a multi documents post and also give to the summarize sentences to the reader. If the compression ratio to calculate the selected posted sentence for monthly in to post text documents divided by the total monthly posted text documents the give the percentage. Example protests posted texts to display each clustered text with a console in Java programming (Net Beans).

**For clustered document 1:**

How many sentences do you want to select in doc 1 to the summarization: 2?

Enter date of summary time Month-Year: Mar-2016

Total numbers of posted text documents in doc 1:104

The compression ratio of doc 1 is: 98.07692307692308

- 8-Mar-2016 Journalists arrested in Oromia Released #Ethiopia #FreePress #OromoProtests #Journalism: 1.8259625047899546E-5
- 12-Mar-2016 with the protests in the Oromia region of EthiopiaRelated: #Diaspora #Oromo #OromoProtests in Oromia, DC and other parts of the country prime minister Hailemariam Desalegn has said that the states will take full responsibility for the public outcry.1.60538780117E-5

**For clustered document 2:**

How many sentences do you want to select in doc 2 to the text Summarization: 1?

Enter the date of summary time Month-Year: Jun-2016 Total number posted text documents in doc 2: 40

The compression ratio of doc 2 is: 97.5

- 15-Jun-2016 Bench Maji Zone Head Office Office Packed Education Parents protest: Disruption of High School Principal Office in Bench Maji Zone, Southern Region: 0.102

=====

**For clustered document 3:**

How many sentences do you want to select in doc 3 to the summarization: 3?

Enter date of summary time Month-Year: Aug-2016 Total numbers of posted text documents in doc 3: 90

The compression ratio of doc 3 is: 96.666667

- 23-Aug-2016 Opposition to the government's use of force ESAT (August 16, 2008) The Ethiopian government's crackdown on protesters in Oromia and Amhara states that the next step is to expand Ethiopia's political space: 2.014347526732104 E-4
- In this article, Tom Malinowski described the recent protests in Amhara and Oromia as a major challenge for Ethiopia and the United States.
- It is learned that the people of Bure town and neighboring villages who went out on their initiative to protest the burning of the Baden logo of the organization were called "TPLF thief": 1.598495843953218 E-4

==the three documents are finished the process==

**System summarization for posting texts for each experiment**

As to describe the four experimentations for each post documents by summarizing extracted rate under each experiment, the four experiments are referred to as E1, E2, E3, and E4. All of these experiments would be conducted for all the posted Amharic input texts for datasets with the summary extraction percentages of 10%, 20%, and 30%. Those summarized extraction percentages are selected to be observed the effect of the summarization processes on different ranges of summarization posted text documents. The total numbers of automatic summaries of each post texts in protests, droughts, sports, and floods are 120,862, 43,774, 10,299, and 1,209 respectively posted Amharic texts as training sets on social media, having 30 posts texts automatic summaries in each experiment for a test set for each clustered post text. For this is 30 posts texts are selected randomly in our training data set corpus of the clustered post documents as to a summary.

System summary is the number of sentences extracted for each post text using the selected percentage is to determine by multiplying the number of sentences in the post texts in a sentence level by the percentage and rounding it's to the nearest zero post texts in a sentence. That is if the original one clustered posts texts have 67 sentences, within a 10% summary that has 7 sentences, within 62 sentences in a 10% summary that has 6 sentences, and so on to find for each sentence in a document. It is depending on the extraction rate post texts and the number of sentences in the document, this Summarizer selects the first "n text file" sentences of the documents.

$n = extraction\ rate \times\ numer\ of\ sentences$

Those lists of the number of sentences extracted using each percentage rate for each posted text file presented for each (Table 2).

**Table 2:** The number of sentences extracted using the selected percentages for each posted texts.

File name	Post texts	No of totals sentences	No. sentences in system and Ideal summary under E1,E2,E3,and E4		
			10%	20%	30%
<b>Protests posts</b>					
Prot1.txt	30	53	5	11	16
Prot2.txt	30	62	6	12	19
Prot3.txt	30	53	5	11	16
Average score	30	56	5	11	17
<b>Droughts posts</b>					
Drot1.txt	30	67	7	13	20
Drot2.txt	30	62	6	12	19
Drot3.txt	30	59	6	12	18
Average score	30	63	6	12	
<b>Sport posts</b>					
Spot1.txt	30	42	4	8	13
Spot2.txt	30	38	4	8	11
Spot3.txt	30	34	3	7	10
Average score	30	38	4	8	11
<b>Floods posts</b>					
Flood1.txt	30	56	6	11	17
Flood2.txt	30	43	4	9	13
Flood3.txt	30	30	3	6	9
Average score	30	42	4	9	13

**Table 3:** The average performance comparison results for all experiments.

Group File	Average of F-measure in percentage (%) for all E1, E2,E3,E4																	
	At 10%			E1			E2			E3			E4			Differ-ence		
Protests(E1)	74%	73.81	-0.19	74	81	7	74	86.11	12.11									
Droughts(E2)										73.81	81	7.19	73.8	86.1	12.3			
Sports(E3)																81	86.11	5.11
Floods(E4)																	86.11	
<b>At20%</b>																		
Protests(E1)	85.4	83.54	-1.81	85.35	89.25	3.9	85.35	89.56	4.21									
Droughts(E2)										83.54	89.3	5.71	98.3	89.56	6.02			
Sports(E3)																89.25	89.56	0.31
Floods(E4)																	89.56	
<b>At30%</b>																		
Protests(E1)	87.1	84	-3.07	87.07	91.37	4.3	87.07	93.52										
Droughts(E2)										84	91.37	7.37	84	93.5	9.52			
Sports(E3)																91.37	93.52	
Floods(E4)																	93.52	

**Manual evaluation for each experiment**

From in the manual evaluations are prepared for 40 summaries texts to produce from 30 posted texts in news items posted Amharic texts on the social media in the protests, droughts, sports and floods post clusters texts at the three summary extraction percentage rate used in this the experiments. As discussed in the earlier sections, the six manual evaluation, quality measured criteria were used to grade the summaries on the scales to give from 1-5. These grades show that the level of the criteria under the respect of the summary information is achieved in and for each summary in the documents.

**Comparison of each experiment**

The four experiments are implemented in the posted texts on social media in protests post (E1), droughts post (E2), sports post (E3), and floods post (E4) for Twitter and Facebook the users posted Amharic texts. We were tested for each post to take 30 posted texts for each text file and 40 summaries to select in ranked sentences for each experiment both the automatic summary and human summary. The main difference for each experiment for the use of the porter stemmer in the different post texts to remove the affixes by using the Java tools. For this each summary to evaluated by the human/ideal summary in all experiments to generate the summaries of

the post texts in the social media. During the experiments in the four experimentations that observed the difference between their performance on the effectiveness and efficiency to compared to extract the summary and extraction percentages. For the systems evaluated in 1<sup>st</sup>, 2<sup>nd</sup>, 3<sup>rd</sup> and 4<sup>th</sup> experiments are effective for all groups of post texts on the extraction percentages and the size of the summary the average score for each experiment are greater than 70% to score the performances. In all the experiments the evaluation results could be shown the automatic summaries for using the Amharic stemmers to perform the extraction rate of the size of the posted texts with the better consistency in higher rates. We compared the first experiments to the second experiments are more condensed or compressed for the summary, but in the third and the fourth experiments are less compressing the posted documents at a 10% extraction rate. For the others, at 20% extraction rates for the size of the summary in the first Experiments E1 is better results than E2 in the group of post texts, such as 1.81%, 3.9%, 4.21% and in the droughts 5.71%, 6.02% and for in sports post groups is 0.31%. We compare the first experiments for the other experiments in the extraction rate at 20%, it is greater than E2 for compressing the summary, but it is less than others 20% extraction rate for the post texts. The size of the summary increase, the extraction rate is increasing in the summary. In all the experiments at 20% summary the highest performance to compress the summary in the E4, E3, E1, E2 and 0.31%, 5.71%, 1.81%, and 3.9% respectively. At a 30% extraction rate for the summary to compare for each experiment E1 is the best result outperformed than E2 in the groups of posts for 3.07%, 4.3%, 6.45%, and in the droughts post 7.33%, 9.52%, and in the sports 2.15%. In the experiments at 30% summary the highest performance to compress the summary in the E4, E3, E1, E2 2.15%, 7.37%, 3.07%, and 4.3% respectively. The average score for Experiment E4 is more condensed the posted documents for the summary than the other experiments, which is 93.52% at a 30% extraction rate (Table 3).

## DISCUSSION

In this paper, we reported the evaluation of the results of our research news items posted Amharic text summarization on social media for posting text on Twitter and Facebook and how to process, the steps minimize the error and to increase the products of the results to summarize Amharic text our research algorithm is hybrid tf-IDF. The evaluation of the Amharic text summarizes the Intrinsic evaluations are made by comparing for automatic summary within the human/reference summary made manually prepared the summary posted texts on social media. It is measuring the automatic system summary performance on its own posted texts. We evaluated the process of the system as an intrinsic method by using the subjective (qualitative) and objective (quantitative) for posting texts in social media to summaries by automatic and manual summary. The subjective evaluations are used to measure the Informativeness, coherences structure, and linguistic quality measuring criteria of the automatic summary generated for the posted documents. The other evaluation method is objective evaluations that are measured by the summarizer of the performance in the extraction summary and also the identification of the post documents of salient sentences with the given post texts. For this to measure the performance of the standard for precision and recall measure in the given input texts, human summary, and

Summarizer extract percentages how to close the extracts for each other for the system. Whereas, extrinsic evaluations are completed based on evaluating how the automatic summaries of the systems are good enough to be accomplished the purpose of some other specific tasks. As the evaluation of the experiments for each training data set to compress the text summarization by automatically if the user selected the size of the summary of the document and to condense the document. We tested the best compression rate for each training data for protests posts higher or better compress ratio sequentially to condense the objective summary. For the average score for the performance measuring the three extraction percentages at 10%, 20%, and 30% are presented for the protests post texts in the protest group 74%, 85.35% and also 87.07% summarize the protests post texts in social media.

For the droughts training data sets, we tested the best summarize document sequentially to condense the objective summary. For the average score for the performance measuring the three extraction percentages at 10%, 20%, and 30% are presented for the droughts post texts in the droughts group 73.81%, 83.54%, and also 84% summarize the droughts post texts in social media. For sports training data sets, we tested the summarize cluster document sequentially to condense the objective summary. For text summarization in this post texts. For the average score for the performance measuring the three extraction percentages at 10%, 20%, and 30% are presented for the sports post texts in the sports group 81%, 89.25%, and also 91.37% summarize the sports post texts in social media. In Floods training data sets, we tested the summarize document sequentially to condense the objective summary. For the average score for the performance measuring the three extraction percentages at 10%, 20%, and 30% are presented for the floods post texts in the floods group 86.11%, 89.56%, and also 93.52% summarize the floods post texts in social media.

## Future work

- Apply proposed algorithm to summary the identifying the news item posted texts automatically grouped on social media post news rather than to prepare manual the corpus separately posted news items on summarization for others local language.
- Applying for the tasks of multi-documents posted Amharic texts and sequential update summarization tasks are one possible future work to research social media posts.
- Apply more Amharic lexicon rules and the dictionary file to use a dictionary to control over and under stemming the Amharic words.

## CONCLUSION

In this paper, we introduced Amharic text summarization for news items posted on social media. In the information overload on social media are different news items posted in Amharic texts. For this is a single or multi-document to need summarization that contains in the frequent terms or repeated posting documents to summarize using hybrid TF-IDF algorithm is introduced. As discussed in the previous chapters, in this research, a hybrid TF-IDF algorithm was recommended for Amharic and for under-resourced languages in social media posted text documents. The TF-IDF is used to process for counting words in a document as well as throughout a corpus

of documents to the end of sorting documents in statistically relevant ways. With the growth of the web for social media a huge amount of information is posted and the posted documents coming from different sources at this time, it becomes very difficult for summarizing the documents to find the specific information they want to summarize. It provides support for those users to get that relevant information in posted Amharic texts and a summary is today mandatory to get relevant important information for the documents for social media.

## ACKNOWLEDGEMENT

I would like to thank our Advisors and families for their given comments, and also those who helped us during this research study.

## REFERENCES

- Guo Q, Diaz F, Yom-Tov E. Updating users about time critical events. 2013: 483-494.
- Agency E N. Amharic Text classification. 2000.
- Xu T, Oard DW, McNamee P. HLTCOE at TREC 2013: Temporal Summarization. TREC. 2013.
- Chakrabarti D, Punera K. Event summarization using tweets. Int AAAI Conf Web Soc M. 2011: 66-73.
- Chua FC, Chong F, Asur S. Automatic summarization of events from social media. Int AAAI Conf Web Soc M. 2013.
- Sayyadi H, Hurst M, Maykov A. Event detection and tracking in social streams. Int AAAI Conf Web Soc M. 2009: 311-314.
- Sakaki T, Okazaki M, Matsuo Y. Earthquake shakes twitter users: Real-time event detection by social sensors. Int AAAI Conf Web Soc M. 2010: 851-860.
- Eidheim HL. Temporal Summarization of Time Critical Events: A system for summarizing events over time in a continuous stream of documents (Master's thesis, NTNU). 2015.
- McBurney PW, McMillan C. Automatic source code summarization of context for java methods. IEEE Trans Softw. 2015; 42: 103-119.
- Premlatha KR, Geetha TV. Extracting temporal patterns and analyzing peak events. 2010.
- Makkonen J, Ahonen-Myka H. Utilizing temporal information in topic detection and tracking. 2003. 393-404.
- Tamiru M, Libsie M. Automatic amharic text summarization using latent semantic analysis. 2009.
- Teklewold AA. Automatic summarization for amharic text using open text summarizer. 2013.
- Yirdaw ED. Topic-based amharic text summarization. 2011.
- Dessalegn KD, Tachbelie MY. Graph-based automatic amharic text summarizer. 2017.
- Saggion H, Radev DR, Teufel S, Lam W, Strassel SM. Developing infrastructure for the evaluation of single and multi-document summarization systems in a cross-lingual environment. LREC. 2002: 747-754.
- Radev DR. Centroid-based summarization of multiple documents: Sentence extraction, utility-based evaluation, and user studies. 2000.
- Barzilay R, Elhadad N. Inferring strategies for sentence ordering in multidocument news summarization. J Artif Intell Res. 2002; 17: 35-55.
- Orasan C. Comparative evaluation of modular automatic summarisation systems using CAST. 2006.
- Marcu D. Discourse-based summarization in duc-2001. Proc Int Conf Doc Anal Recognit. 2001.
- Inouye D, Kalita JK. Comparing twitter summarization algorithms for multiple post summaries. IEEE. 2011: 298-306.
- He Z, Chen C, Bu J, Wang C, Zhang L, Cai D et al. Document summarization based on data reconstruction. AAAI. 2012.
- Paroubek P, Chaudiron S, Hirschman L. Principles of evaluation in natural language processing. 2007; 48 (1): 7-31.
- Tao K, Abel F, Gao Q, Houben GJ. Tums: Twitter-based user modeling service. ESWC. 2011: 269-283.
- Chung W, Chen H, Chaboya LG, O'Toole CD, Atabakhsh H. Evaluating event visualization: A usability study of COPLINK spatio-temporal visualizer. Int J Hum Comput. 2005; 62 (1): 127-157.
- Busch M, Gade K, Larson B, Lok P, Luckenbill S, Lin J. Earlybird: Real-time search at twitter. IEEE. 2011.
- Potthast M, Stein B, Loose F, Becker S. Information retrieval in the commentsphere. ACM Trans. Intell Syst Technol. 2012; 3: 1-21.
- Novalija I, Papler M, Mladenici D. Towards social media mining: Twitterobservatory. 2014.
- Sharifi B, Hutton MA, Kalita JK. Experiments in microblog summarization. IEEE. 2010: 49-56.
- Meladianos P, Nikolentzos G, Rousseau F, Stavarakas Y, Vazirgiannis M. Degeneracy-based real-time sub-event detection in twitter stream. 12th Int AAAI Conf Web Soc Media ICWSM. 2018.
- Ibekwe-Sanjuan F, Silvia F, Eric S, Eric C. Annotation of scientific summaries for information retrieval. 2011.
- Bontcheva K, Derczynski L, Funk A, Greenwood MA, Maynard D, Aswani N. Twitie: An open-source information extraction pipeline for microblog text. Int Conf Recent Adv Nat. Lang. 2013: 83-90.
- Truong B, Caragea C, Squicciarini A, Tapia AH. Identifying valuable information from twitter during natural disasters. ASIS&T. 2014; 51 (1): 1-4.
- Seki Y. Sentence extraction by tf/idf and position weighting from Newspaper Articles. 2003.
- Nagwani NK, Verma S. A frequent term and semantic similarity based single document text summarization algorithm. Int J Comput Appl. 2011; 17 (2): 36-40.
- Saggion H, Poibeau T. Automatic text summarization: Past, present and future. 2013.
- Neto JL, Freitas AA, Kaestner CA. Automatic text summarization using a machine learning approach. 2002.
- Eyassu S, Gambäck B. Classifying amharic news text using self-organizing maps. 2005: 71-78.
- Dao TN, Simpson T. Measuring similarity between sentences. 2005.