

# Comparative Analysis of Different Text Summarization Techniques Using Enhanced Tokenization

Tanzirul Islam\*, Mofazzal Hossain<sup>†</sup>, MD. Fahim Arefin<sup>‡</sup>

Dept of Computer Science and Engineering, Green University of Bangladesh, Bangladesh\*<sup>†‡</sup>  
Email: tanziruljoy\*@gmail.com, mdmofazzalhossainhossain01<sup>†</sup>@gmail.com, fahim<sup>‡</sup>@cse.green.edu.bd

**Abstract**—As a huge amount of data is being generated everyday, text summarization is a must-have technique to obtain the required information concisely. Summaries reduce reading time. When it comes to researching documents, summaries make the job easier. The challenge of creating a short and fluent summary while retaining important information content and overall meaning is known as automatic text summarization. As a huge amount of data is being generated everyday, text summarization is a must-have technique to obtain the required information concisely. It is simple to deal with summarization in the other languages like English, Turkish, Arabic. But due to the diverse and complex nature of the Bangla language, not much has been done on the technique of summarizing the Bangla text. Given the importance of text summarization, this paper focused on the creation of an extraction-based summary approach that works on Bangla text documents. Here we apply different kind of model for generating a summary for a single bangla text document. As compared to other outcomes, our experimental results are outstanding and people who read the summary evaluated them. Further development of these methods will undoubtedly deliver fascinating results. This can also contribute significantly in the effort to build smart machines, which form the basis of industry 4.0.

**Index Terms**—Natural Language Processing, Text Pre-processing, Cosine Similarity based summarization, Extractive summarization, Text Rank based summarization, Word count and heapQ based summarization

## I. INTRODUCTION

Natural Language Processing (NLP) is a type of artificial intelligence that enables computers to read, comprehend, and interpret human language. It assists computers in measuring sentiment and determining which portions of human language are significant. However, because of the vast amount of unstructured data, the lack of formal rules, and the lack of real-world context or intent, human language is exceedingly difficult for computers to interpret. Automatic text summarization is the task of producing a concise and fluent summary while preserving the key information content and overall meaning. It is challenging for machines because humans can readily choose the keywords from a given text and summarize them, but computers cannot. As summarization requires linguistic comprehension, reasoning, and the application of common

sense information in the same way that humans do, text summarization is quite complex. Summaries facilitate the selection process since they minimize reading time while investigating documents.[1] As a result, there is a strong need to simplify much of this text material into shorter, focused summaries that capture the essential information, both so we can traverse it more efficiently and to determine whether the bigger papers include the information we want. Text can be summarized in two ways: extractive and abstractive. [2]

In extraction-based summarizing, a subset of words representing the most significant points from a piece of text is extracted and merged to form a summary. Consider it similar to a highlighter, which chooses the most important information from a source text. Advanced deep learning algorithms are used in abstraction-based summarization to paraphrase and reduce the original material in the same way that people do.[3] Consider it a pen that generates fresh sentences that may or may not be part of the underlying material. Bangla is one of the world's most widely spoken languages. However, text summarization in Bangla has been limited by various of issues, including the lack of digital standard text databases and inflectional morphology in Bangla. As a consequence, we performed an analysis of several text summarization approaches utilizing enhanced tokenization and obtained some interesting findings. To obtain the desired outcome, we must first go through certain procedures like tokenization, stop words removal, POS tagging, and model construction.

Where The initial step in any NLP pipeline is to tokenize the data. [4] It has a significant impact on the entire pipeline. A tokenizer is a program that converts unstructured data and natural language text into discrete information pieces. Stop word removal is one of the most often utilized preprocessing stages in various NLP applications. The goal is to simply remove words that appear often across all of the documents in the corpus and POS tagging is the process of transforming a sentence into a form – a list of words, a list of tuples (each tuple having a form (word, tag)). The tag in the case of is a part-of-speech tag that indicates whether the word is a noun, adjective, verb, or other.

We deal with all of these phases in our technique. We

implemented most of the steps of text summarization in order to increase efficiency and produce a more fluent and effective summary.

## II. LITERATURE REVIEW

The last several years a large amount of work have been done with extractive method and many papers have been published on it. And Bengali being the 7<sup>th</sup> most spoken language all over the world.

Shofi Ullah, Sagar Hossain, K. M. Azharul Hasan [5] worked on cosine similarity based graph ranking and Relevance based scoring and ranking approach for the summarization of bangla text. And the Method consist of: POS tagging, redundancy removal algorithm & cosine similarity Based graph. Sarkar [6] proposed a summarization technique implementing term frequency and semantic sentence similarity based summarizing approach to summarize a single Bangla text. This method consisted of four major steps: preprocessing, stemming, sentence ranking and summary generation. Sentences are ranked based on two important features: thematic term and position. In paper [7], the reasearchers take a document as input, find the word frequency, determine the sentence similarity, weight the sentences, sort the sentences according to their rank and finally select the sentences with the higher rank for generating the summary.

Talukder et. al. [8] worked with bi-directional RNNs with LSTM in encoding layer and attention model at decoding layer. In this model, they implement an abstractive summarizer and reduce the given text and run time of the model. And the model works as sequence to sequence model to generate summary. Mahsa Afsharizadeh, Hossein Ebrahimpour-Komleh, Ayoub Bagheri [9] developed a query-oriented text summarization technique which extracted the most informative sentences. Each sentence had 11 of the greatest qualities taken from it. They used the ROUGE criterion to illustrate that using more appropriate features results in better summaries. They use certain appropriate features, with the first set of features identifying informative sentences and the second set of features assisting in the extraction of query relevant sentences. To improve text summarizing, their system reliably picks the most informative sentences and then provides the summary.

Narendra Andhale, L.A. Bewoor's work [4] includes a clear examination of both the extractive and abstractive approaches to text summarization, as well as their respective limitations. The examined literature points to a difficult area where these methodologies could be used to provide interesting, well-compressed, and legible summaries.

Md Mahbubur Rahman, Md. Aktaruzzaman Pramanik, Rifat Sadik, Monikrishna Roy, Partha Chakraborty [2] used a transformer-based bangla document classification system. They used the most recent transformer or attention mechanism-based models to classify Bangla text con-

tent. For Bengali text classification, they used the BERT (Bidirectional Encoder Representations from Transformers) and ELECTRA (Efficiently Learning an Encoder that Classifies Token Replacements Accurately) models. Sumya Akter, Aysa Siddika Asa, Md. Palash Uddin, Md. Delowar Hossain, Shikhor Kumer Roy, Masud Ibn Afjal,[10] described a text summarizing algorithm that extracts key lines from a single or many Bengali papers. They used the K-means clustering approach to provide the final summary for single or many articles. They compared the results to those of other extractive summarization techniques and measured the run-time complexity, demonstrating that the new technique's performance has improved. Sheikh Abujar, Mahmudul Hasan, M.S.I Shahin Syed Akhter Hossain [3] focused with the extraction method for summarizing Bangla text. Basic extractive summarization was used for the new proposed model in their proposed amendment. A set of rules for analyzing Bangla text was also developed from the heuristics. The material is carefully examined using the Bangla sentence clustering method. Their methodology groups all of the data together and presents the information in a clear and concise manner. Prachi Shah and Nikita P. Desai[1] examined text summarizing techniques for a variety of Indian and foreign languages, including English, European, and others. They also offered a method for summarizing Hindi text using machine learning techniques. This research focused on surveying and analyzing the performance of automated text summarizers for a variety of Indian and international languages. The top ranked sentences are considered in the summary, and their approach assigns ranks based on sentence score.

In this section, we summarized the state of the art works in the field of text summarization in Bangla. We discuss our proposed methodology in the next section.

## III. PROPOSED METHODOLOGY

Before summarizing a text, we need to complete a few pre-processing steps, we highlight those in Fig. 1.

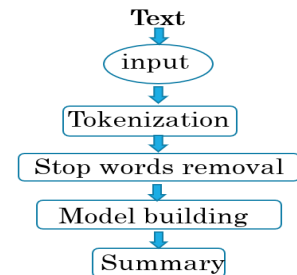


Fig. 1: Work Flow of summary generation

For every approach, we have to complete some pre-processing. Those are tokenization, removing stopwords, and lemmatization.

(a) Tokenization: Tokenization is one of the common task when it comes to working with text data. Tokenization is essentially splitting a phrase, sentence, or an entire text document into smaller units such as individual words or sentence. units are called tokens.

(b) Removing Stopwords: Stop words are the words which have less significance in texts and have very high or very low occurrence in the document and they are removed. More than 350 stop words available for bengali texts such as: অতএব, আগামী, অবধি, আর, আরও, ইত্যাদি, এত, এতটাই, অথচ, অথবা

(c) Lemmatization: Lemmatization is a method for combining various inflected forms of words into a single root form with the same meaning. It's similar to stemming in that it results in a stripped-down word with dictionary meaning. There are several words in different formations like বাংলাদেশী, বাংলাদেশে, বাংলাদেশীয়, বাংলাদেশকে etc. They basically denote the same meaning বাংলাদেশে but with the inclusion of terms like 's', 'es' it shows the different form. To remove redundancy, the conversion of several formations into an actual word is also needed.

In this paper, we used four different approaches, for generating a summary for a single text document which are as follows:

- 1) Cosine Similarity based summarization
- 2) Extractive summarization
- 3) Text Rank based summarization
- 4) Word count and heapQ based summarization

#### A. Cosine similarity-based summarizing

It is mostly used for hierarchical clustering and multidimensional scaling, after which we evaluate sentence similarity and merge the ranking sentences for summary those that are approximately one or one. To begin cleaning, we remove the white space, bangla digits, upper and lower case English letters, and punctuation-marks. After cleaning the text, we split it into sentences and then ranked each phrase on a scale of (0-1) it's would be like

বিশ্বকাপে ৮ ম্যাচে ৩৪ ২৮ গড়ে ২৪০ রান করছিলেন  
0.08186748647115961  
২টি সিক্বেচর ও ৯টি ফিফটি ছিল তাঁর 0.07580509024530348  
৪২ ম্যাচে ২৮ ৪৮ গড়ে ৮৮৩ রান করছিলেন 0.0744032242252412  
এ ছাড়া অস্ট্রেলিয়ার বপিক্ষে ৪০ রান করছিলেন যশপাল  
0.07402828995191987  
তাত্তে ছিল চারটি ফিফটি 0.0665731254875899  
তাত্তে ৩৩ ৪৫ গড়ে ১৬০৬ রান করছিলেন 0.0651854654535155  
ছিল দুটি ফিফটি 0.062214838422544744

Here we can see that every sentence are ranked to (0-1) and since we are using the positive quadrant, which means no sentence may have a negative value, before scoring all sentences from the input text. Finally, we compute the summary frequency (how many sentences are produced) from the rated sentences.

#### B. Extractive Summarization:

There are two methods for extractive summarization: First of all The extractive method will take the same words, phrases, and sentences from the original text. Extractive methods can be considered as important sentence selection in the given text - one is percent summary and the other one is word count summary.

Percent summary: For percent summary we use ratio = 0.5 which means that the ratio between two sentences. \* the similarity measure of each sentence with ratio 0.5(50% after that ranked those sentences for summary

Word count summary: for Word count summary we use word-count = 100 which means that the summary will be within 100 words.

- Here we just count those word are repeated many time in a sentence.
- Then select those sentences for summary within 100 words

#### C. Textrank based summarization

It focuses on the ranking of text sentences and is computed recursively based on information accessible throughout the text. TextRank is a graph-based algorithm, easy to understand and implement. It does the necessary preprocessing before calculating the similarity between texts. Then it ranks the sentences based on their importance before providing the summary. TextRank works in the following steps:

- 1) Tokenize documents into sentences.
- 2) Preprocess each sentence in the document.
- 3) Calculate the Similarity between sentences.
- 4) Rank the sentences with higher significance.

[মোল্লা মোহাম্মদ হাসান আখুন্দকে প্রধানমন্ত্রী করে ৭ সেপ্টেম্বর আফগানিস্তানরে অন্তর্বর্তীকালীন সরকাররে ঘোষণা দিয়েছে তালবান।', 'মন্ত্রিসভার সবাই পুরুষ এবং আগ্নেয়াস্ত্রসহ কন্দ্ৰীয় ব্যাংকরে গভর্নররে ছবি আমরা সামাজিক মাধ্যমে দেখেছি।', 'তারও আগে ১৫ আগস্ট রাজধানী কাবুলসহ দেশটির নয়ন্তরণ নিয়ে তারা।'] This is the first step to tokenize documents into sentences and preprocess each sentence in the text document. The output summary will consist of the most representative sentences and will also be returned as a string, divided by newlines. If the split parameter is set to True, a list of sentences will be returned. The length of the output can be specified using the ratio and word-count parameters: ratio should be a number between 0 and 1 that determines the percentage of the number of sentences of the original text to be chosen for the summary (defaults to 0.2). word-count determines how many words will the output contain. If both parameters are provided, the ratio will be ignored, and finally given the summary for the given text those sentence are higher significance, and this method is more popular nowadays for both Abstractive and Extractive summarization.

iv. The input text is cleaned first by word count and heapQ-based summarization. Following cleaning, After cleaning the text like (white space, bangla-digits Like(1,2,3,4), english letter (Aa-zZ) both upper case and lower case, punctuation-mark. We remove the stop word then we divide it to sentences, Then scoring all sentence using word count (it means that it count all word of a sentence) it eliminates stopwords, tokenizes the text, and scores all sentences based on word count. Then scoring all sentence we use heapQ for our summary that take the largest value only to make the summary. For example-

মরীচিকা চরকরি আট পর্বের অরজিনাল ওয়েবের সরিজে 2.75  
ওয়েবের সরিজিটির পরচালক শহিব শাহীন জানালেন দর্শককে এই সাড়া  
তাকে আপলুত করে 4.5  
তিনি বলেন টরলোরটি সরিজিরে একটি চুম্বক অংশ 1.75  
আমার বিশ্বাস পুরো সরিজি দর্শককে আরও ভালো লাগবে 3.75  
বাংলাদেশের কনটেন্ট হিসেবে এই সরিজি দর্শককে আলাদাভাবে ভাবাবে  
3.25  
তবে পুরো কাজ তুলে আনাটা তাঁর জন্য চ্যালেঞ্জিং ছিল বলে জানালেন  
শহিব শাহীন 4.5

.....  
Here we can see that every sentence are tokens with their individual significance number and among them which sentence is provide highest number that sentence will be the part of the summary.

#### IV. EXPERIMENTAL RESULT

In this section we discussed all methods input and output as sample text document each and every method have an own text document and we represented all of methods. First, we collected some articles from Prothom Alo newspaper. Then we summarize it, and then we compare each and every method's output to each other, according to what we consider to be each method's summary. We give these outputs to some of our friends and senior brothers to know which output is nearer to a human summary or more accurate as a summary.

In TABLE I, we present a sample result for Cosine Similarity based summarization. At first we take a text document as input and after processing, we get our summary for the given text document.

For Extractive summarization we represented percent summary and word count summary in TABLE II. In percent summary, we use ratio = 0.5 which refers to the ratio between main text and summary length. We can change the ratio as for our given text here we use 0.5 because we use here mid-sized text document. If it became large we reduce the value of it as for our user requirement. It will find the ratio between two sentence from the given text and after that rank those sentences to generate the percent summary. In extractive method, it takes the same words, phrases, and sentences from the original text. For Word count summary, we use word-count = 100 which means that the summary will be within 100 words for the given text document. We can change the value of it as for our user requirement. Extractive summarization is easy to implement

TABLE I: sample result for Cosine Similarity based summarization

Input	সাবকে ভারতীয় ক্রিকেটের ও ১৯৮৩ বর্ষিকাপজয়ী যশপাল শর্মা আজ মঙ্গলবার মৃত্যুবরণ করছেন। ৬৬ বছর বয়সী এই বিশ্বকাপজয়ী ক্রিকেটের মৃত্যুর কারণ হার্ট অ্যাটাক বলে জানিয়েছে ভারতীয় সংবাদমাধ্যম। যশপালসুত্রী ও তিনি সন্তান রখে গেছেন। সত্তরতর দশককে শেষ থেকে আশরি দশককে মধ্যভাগ পর্যন্ত ভারতীয় দলের অংশ ছিল যশপাল। পাঞ্জাবের এই মডিল অন্ডাররে অভ্যিকে হয়ছিল পাকিস্তানে বপিক্ষে ১৯৭৮ সালে শিয়ালকোটের সো ওয়ানডরে পর আরও ৪১টি ওয়ানডে খেলেন যশপাল। ৪২ ম্যাচে ২৮.৪৮ গড়ে ৮৮৩ রান করছিলেন। তাত ছিল চারটি ফিফটি। বো যশপাল তাঁর ক্যারিয়ারের সেরা মুহূর্ত কাটিয়েছেন ১৯৮৩ বর্ষিকাপে। বর্ষিকাপে ৮ ম্যাচে ৩৪.২৮ গড়ে ২৪০ রান করছিলেন। ছিল দুটি ফিফটি। প্রথম ম্যাচে যশপালের ক্যারিয়ার-সেরা ৮৯ রানের ইনহিসে ওয়েস্ট ইন্ডিজকে চমকে দিয়েছিল ভারত। সমেফাইনালে ইংল্যান্ডের বপিক্ষে ৬১ রানের ম্যাচ জেতানো ইনহিসে ছিল তাঁর। এ ছাড়া অস্ট্রেলিয়ার বপিক্ষেও ৪০ রান করছিলেন যশপাল। ভারতের হয়ে ৩৭টি টেস্টেও খেলেন এই ব্যাটসম্যান। তাত ৩৩.৪৫ গড়ে ১৬০৬ রান করছিলেন। ২টি সেঞ্চারি ও ৯টি ফিফটি ছিল তাঁর। প্রথম শ্রণের ক্যারিয়ারে ১৬০ ম্যাচে ২১ সেঞ্চারিতে ৮ হাজার ৯৩৩ রান যশপালের।
Output	বর্ষিকাপে ৮ ম্যাচে ৩৪.২৮ গড়ে ২৪০ রান করছিলেন। ২টি সেঞ্চারি ও ৯টি ফিফটি ছিল তাঁর। ৪২ ম্যাচে ৮.৪৮ গড়ে ৮৮৩ রান করছিলেন। এ ছাড়া অস্ট্রেলিয়ার ২বপিক্ষেও ৪০ রান করছিলেন যশপাল। তাত ছিল চারটি ফিফটি

and understand and it is more familiar for python language. It is more interesting because we get two summary for one method at a time for a single text document. So here is an example of this method and it works on any kind of text document.

In TABLE III and TABLE IV, we try to get a summary with different kind of approaches which is different from other approaches. The output summary will consist of the most representative sentences of the original text for the summary (defaults at 0.2). Here, word-count determines how many words will the output contain to get the summary. Our first step is to tokenize all sentence from the given text and every sentences are ranked with their individual significance number and among them which sentence is provide highest number that sentence will be the part of the summary. And it work for a single text document to generate this summary.

We highlighted our results in Fig. 2. Cosine Similarity based summarization was liked by 55 percent of them, Extractive summarization-70 percent, Text Rank based

TABLE II: sample result for Extractive summarization

Input	<p>মনরে সুখে মায়ের হাতের হালুয়া খয়ে রীতমিত্তে ওজন বাড়িয়ে ফেলেছেন বলডিডরে সুপারফটি সুপারস্টার অক্ষয় কুমার। মাত্র কদিনেই তার পাঁচ কজি ওজন বড়ে গেছে। তবে ওজন বাড়ায় মোটেও বচিলতি নন অক্ষয়। বরং তিনি বিজয়ে খুশি। কারণ, পরের ছবির প্রয়োজনে এমনতিই তাঁর ওজন বাড়ানোটা জরুরি ছিল। এই মুহূর্তে আগামী ছবি 'রক্ষাবন্ধন'-এর শূটিংয়ে ব্যস্ত অক্ষয়। আনন্দ এল রাই পরিচালিত এই ছবটির শূটিং চলছে মুম্বাইতে। সনিমোয় তিনি পাঁচ বোনরে একমাত্র ভাই।</p> <p>ভাই-বোনরে সম্পর্করে এক সুন্দর রসায়ন নিয়ে এক পারিবারিক ছবি নির্মাণ করতে চলেছেন আনন্দ। প্রদায় অক্ষয়রে বোনরে চরিত্রে পাঁচ নবাগতা অভিনেত্রীকে দেখা যাবে। জানা গেছে, এই ছবির চরিত্ররে প্রয়োজন অক্ষয়রে ওজন বাড়ানোর প্রয়োজন ছিল। আর তিনি স্বাভাবিক পদ্ধতিতে নিজের ওজন বাড়িয়েছেন। এ জন্য এই বলডিড সুপারস্টার প্রাণভরে খয়েছেন মায়ের হাতের হালুয়াসহ আরও নানান বাহারি খাবার। এত দিন যা তিনি মুখে তুলতে পারেননি। দীর্ঘদিন পর হালুয়া খেতে পরে দারুণ খুশি অক্ষয়। আর খুব কম সময়ের মধ্যে পাঁচ কজি ওজন বাড়িয়ে ফেলেছেন।</p>
Output	<p>Percent summary: মনরে সুখে মায়ের হাতের হালুয়া খয়ে রীতমিত্তে ওজন বাড়িয়ে ফেলেছেন বলডিডরে সুপারফটি সুপারস্টার অক্ষয় কুমার। মাত্র কদিনেই তার পাঁচ কজি ওজন বড়ে গেছে। কারণ, পরের ছবির প্রয়োজনে এমনতিই তাঁর ওজন বাড়ানোটা জরুরি ছিল। সনিমোয় তিনি পাঁচ বোনরে একমাত্র ভাই। প্রদায় অক্ষয়রে বোনরে চরিত্রে পাঁচ নবাগতা অভিনেত্রীকে দেখা যাবে। জানা গেছে, এই ছবির চরিত্ররে প্রয়োজন অক্ষয়রে ওজন বাড়ানোর প্রয়োজন ছিল। আর তিনি স্বাভাবিক পদ্ধতিতে নিজের ওজন বাড়িয়েছেন। আর খুব কম সময়ের মধ্যে পাঁচ কজি ওজন বাড়িয়ে ফেলেছেন।</p> <p>Word count summary: মনরে সুখে মায়ের হাতের হালুয়া খয়ে রীতমিত্তে ওজন বাড়িয়ে ফেলেছেন বলডিডরে সুপারফটি সুপারস্টার অক্ষয় কুমার। মাত্র কদিনেই তার পাঁচ কজি ওজন বড়ে গেছে। তবে ওজন বাড়ায় মোটেও বচিলতি নন অক্ষয়। কারণ, পরের ছবির প্রয়োজনে এমনতিই তাঁর ওজন বাড়ানোটা জরুরি ছিল। এই মুহূর্তে আগামী ছবি 'রক্ষাবন্ধন'-এর শূটিংয়ে ব্যস্ত অক্ষয়। সনিমোয় তিনি পাঁচ বোনরে একমাত্র ভাই। প্রদায় অক্ষয়রে বোনরে চরিত্রে পাঁচ নবাগতা অভিনেত্রীকে দেখা যাবে। জানা গেছে, এই ছবির চরিত্ররে প্রয়োজন অক্ষয়রে ওজন বাড়ানোর প্রয়োজন ছিল। আর তিনি স্বাভাবিক পদ্ধতিতে নিজের ওজন বাড়িয়েছেন। দীর্ঘদিন পর হালুয়া খেতে পরে দারুণ খুশি অক্ষয়। আর খুব কম সময়ের মধ্যে পাঁচ কজি ওজন বাড়িয়ে ফেলেছেন।</p>

TABLE III: sample result for Text Rank based summarization

Input	<p>২৮ বছর অপেক্ষার পর পাওয়া পরম আরাধ্য এক শরীপো। আরজনেটিনায় এখন খুশি জোয়ার। আবগে বাঁধ না মানাই স্বাভাবিক। দলরে মডিফলিডার রদর্গিগোদি পলও আবগে লাগাম দতি পারনেনি। আরজনেটিনার হয়ে বাতসিত্তা-সমিওনরো সেই ১৯৯৩ সালে যখন মহাদেশেরো হয়েছিল, দি পলরে তখন জন্মও হয়নি। সমর্থক হিসেবে হোক বা খেলোয়াড় হিসেবে, দলরে সাফল্য দেখেননি কখনো। কোপা আমেরিকা জেতার পর উচ্ছ্বসিত দি পলরে হয়তো স্বাভাবিক জাগ্রন একটু হলেও লোপ পেয়েছিল। শরীপো নিয়ে উদ্যাপন করতে গিয়ে ব্রাজিলিকে টেনে আনতে চাইছিলেন। কোপার ফাইনাল দেখতে সাড়ে ছয় হাজারেও বেশি দর্শককে টুকতে দেওয়া হয়েছিল মারাকানায়, যাঁদের অধিকাংশই ছিলেন ব্রাজিলি-সমর্থক। ভাঙা হৃদয় নিয়ে দেখেছিলেন আরজনেটাইনদরে শরীপো উৎসব। এই অবস্থায় তাঁদের কাটা যায় নুনের ছটা দেওয়ার যাবতীয় বন্দোবস্ত করে ফেলেছিলেন দি পল। উদ্যাপন করতে গিয়ে ব্রাজিলিবিদ্রোহী এক গান গাইতে শুরু করেছিলেন। কিন্তু মসি সটো হতে দেবেন কনে। আঙুল উঁচিয়ে দি পলকে সতর্ক করে দিয়েছেন মসি, যেনে ব্রাজিলি দর্শকদের সামনে তাঁদের দেশকে হয়ে করে কোনো গান গাওয়া না হয়। ফুটবলরে জয়-পরাজয় মাঠেই থাকুক, স্টোর ওপর ভিত্তি করে দুই দলের মধ্যে পারস্পরিক শ্রদ্ধাবোধের জায়গাটুকু যেনে নষ্ট না হয়, সে দিকে সতর্ক ছিলেন মসি। দি পলকে সাফ জানিয়ে দিয়েছেন, 'আমার চোখের সামনে এসব গাওয়া যাবে না।' ভিডিওটা সামাজিক যোগাযোগমাধ্যমে ছড়িয়ে পড়েছে এর মধ্যেই। অর্থনায়ক হিসেবে মসির সুনাম চলছে সব জায়গায়। শুধু মসিই নন, দি পলকে সেই গান গাইতে মানা করলে আরকে সনিয়ির সতীর্থ সর্হেও আগুয়রোও। তবে মাঠে দর্শকদের সামনে না গাইলেও ড্রসেইরুমে দি পলকে আটকে রাখা যায়নি। সতীর্থদের সঙগে ওই গান গয়েছেন, তালে তালে নচেছেন। সে ভিডিওতে অবশ্য মসিকে অংশ নতি দেখা যায়নি।</p>
Output	<p>উদ্যাপন করতে গিয়ে ব্রাজিলিবিদ্রোহী এক গান গাইতে শুরু করেছিলেন। আঙুল উঁচিয়ে দি পলকে সতর্ক করে দিয়েছেন মসি, যেনে ব্রাজিলি দর্শকদের সামনে তাঁদের দেশকে হয়ে করে কোনো গান গাওয়া না হয়। শুধু মসিই নন, দি পলকে সেই গান গাইতে মানা করলে আরকে সনিয়ির সতীর্থ সর্হেও আগুয়রোও। তবে মাঠে দর্শকদের সামনে না গাইলেও ড্রসেইরুমে দি পলকে আটকে রাখা যায়নি।</p>

TABLE IV: sample result for Word count and heapQ based summarization

Input	<p>গত জুনরে শুরুতে 'চরক' ইউটিউব চ্যানেলে 'মরীচিকা'র ট্রলোর প্রকাশের পরপর ফেসবুকে ছড়িয়ে পড়ে। পুরো সরিজি দখোর জন্য 'মরীচিকা' চরকরি</p> <p>আট দর্শকরে মধ্যে বাড়তি আগ্রহ তরৈ হয়। পরবরে অরজিনাল ওয়বে সরিজি। ওয়বে সরিজিটির পরচালক শহিব শাহীন জানালনে,দর্শকরে এই সাড়া তাঁকে আপ্লুত করে। তিনি বলেন, 'ট্রলোরটি সরিজিরে একটা চুম্বক অংশ। আমার বশ্বাস পুরো সরিজি দর্শকরে আরও ভালো লাগবে।</p> <p>বাংলাদেশেরে কনটেন্ট হিসেবে এই সরিজি দর্শকরে আলাদাভাবে ভাববে।'</p> <p>তবে পুরো কাজ তুলে আনাটা তাঁর জন্য চ্যালেঞ্জিং ছিল বলে জানালনে শহিব শাহীন। শূটিংয়ে নমে নরিধারতি শডিউলরে আরও দড়ে গুণ শডিউল বাড়তি দতিে হয়েছে। সেই যুদ্ধে শলিপিরাও সহযোগিতার হাত বাড়িয়ে দিচ্ছেনে।</p> <p>মরীচিকার চার অভিনয়শলিপি মাহি, নশিা, সয়িম ও জোভান বলেন, 'শলিপিদের সহযোগতি না পলে কাজটি হয়তো শেষে করতো পারতাম না।</p>
Output	<p>পুরো সরিজি দখোর জন্য দর্শকরে মধ্যে বাড়তি আগ্রহ তরৈ হয়। ওয়বে সরিজিটির পরচালক শহিব শাহীন জানালনে দর্শকরে এই সাড়া তাঁকে আপ্লুত করে। তবে পুরো কাজ তুলে আনাটা তাঁর জন্য চ্যালেঞ্জিং ছিল বলে জানালনে শহিব শাহীন। আমার বশ্বাস পুরো সরিজি দর্শকরে আরও ভালো লাগবে</p>

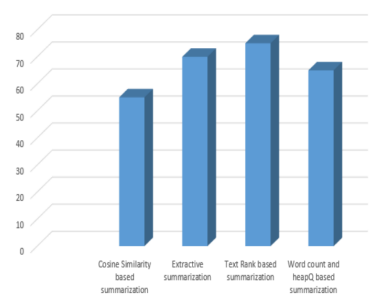


Fig. 2: Result Analysis

summarization-75 percent, and Word count and heapQ based summarization was liked by 65 percent of them.

## V. CONCLUSION

In this paper, we have implemented and compared text summarization using different techniques. To compare, We summarize a single bangla text document with all approaches, namely Cosine Similarity based Summarization, Extractive summarization, Text Rank based summarization, Word count and heapQ based summarization. The first one is based on performing hierarchical clustering and multidimensional scaling after that we check sentence similarity then join the ranked sentence for summary those

nearly one or one. In extractive summarization two methods are implemented one is percent summary and the other one word count summary. Textrank based summarization It focused on the ranking of text sentences and computed recursively based on information accessible throughout the text. The input text is cleaned first by word count and heapQ-based summarization, it eliminates stopwords, tokenizes the text, and scores all sentences based on word count. After scoring all the sentences, heapQ was used to get the final summary. In future we improve our work with different kind of approaches and develop an android app for the user.

## ACKNOWLEDGMENT

This work was supported in part by the Center for Research, Innovation, and Transformation of Green University of Bangladesh.

## REFERENCES

- [1] P. Shah and N. P. Desai, "A survey of automatic text summarization techniques for indian and foreign languages," in 2016 International Conference on Electrical, Electronics, and Optimization Techniques (ICEEOT), 2016, pp. 4598--4601.
- [2] N. Andhale and L. Bewoor, "An overview of text summarization techniques," in 2016 International Conference on Computing Communication Control and automation (ICCCBEA), 2016, pp. 1--7.
- [3] S. Abujar, M. Hasan, M. Shahin, and S. A. Hossain, "A heuristic approach of text summarization for bengali documentation," in 2017 8th International Conference on Computing, Communication and Networking Technologies (ICCCNT), 2017, pp. 1--8.
- [4] M. Afsharizadeh, H. Ebrahimpour-Komleh, and A. Bagheri, "Query-oriented text summarization using sentence extraction technique," in 2018 4th International Conference on Web Research (ICWR), 2018, pp. 128--132.
- [5] A. Sarkar and M. S. Hossen, "Automatic bangla text summarization using term frequency and semantic similarity approach," in 2018 21st International Conference of Computer and Information Technology (ICIT), 2018, pp. 1--6.
- [6] S. Ullah, S. Hossain, and K. M. Azharul Hasan, "Opinion summarization of bangla texts using cosine similarity based graph ranking and relevance based approach," in 2019 International Conference on Bangla Speech and Language Processing (ICBSLP), 2019, pp. 1--6.
- [7] S. Abujar, A. K. M. Masum, M. Mohibullah, Ohidujaman, and S. A. Hossain, "An approach for bengali text summarization using word2vector," in 2019 10th International Conference on Computing, Communication and Networking Technologies (ICCCNT), 2019, pp. 1--5.
- [8] M. A. I. Talukder, S. Abujar, A. K. M. Masum, F. Faisal, and S. A. Hossain, "Bengali abstractive text summarization using sequence to sequence rnns," in 2019 10th International Conference on Computing, Communication and Networking Technologies (ICCCNT), 2019, pp. 1--5.
- [9] M. M. Rahman, M. Aktaruzzaman Pramanik, R. Sadik, M. Roy, and P. Chakraborty, "Bangla documents classification using transformer based deep learning models," in 2020 2nd International Conference on Sustainable Technologies for Industry 4.0 (STI), 2020, pp. 1--5.
- [10] S. Akter, A. S. Asa, M. P. Uddin, M. D. Hossain, S. K. Roy, and M. I. Afjal, "An extractive text summarization technique for bengali document(s) using k-means clustering algorithm," in 2017 IEEE International Conference on Imaging, Vision Pattern Recognition (icIVPR), 2017, pp. 1--6.