

# Extractive Text Summarization Using Word Similarity-based Spectral Clustering

FAHIM MORSHED, Institute of Information Technology, University of Dhaka, Bangladesh

MD. ABDUR RAHMAN, Centre for Advanced Research in Science, University of Dhaka, Bangladesh

SUMON AHMED, Institute of Information Technology, University of Dhaka, Bangladesh

Extractive Text Summarization is the process of picking the best parts of a larger text without losing any key information. This is really necessary in this day and age to get concise information faster due to digital information overflow. Previous attempts at extractive text summarization, specially in Bengali, either relied on TF-IDF or used naive similarity measures both of these suffers expressing semantic relationship correctly. The objective of this paper is to develop an extractive text summarization method for Bengali language, that uses the latest NLP techniques and extended to other low resource languages. We developed a word Similarity-based Spectral Clustering (WSbSC) method for Bengali extractive text summarization. It extracts key sentences by grouping semantically similar sentences into clusters with a novel sentence similarity calculating algorithm. We took the geometric means of individual Gaussian similarity values using word embedding vectors to get the similarity between two sentences. Then, used TF-IDF ranking to pick the best sentence from each cluster. This method is tested on four datasets, and it outperformed other recent models by 43.2% on average ROUGE scores (ranging from 2.5% to 95.4%). The method is also experimented on Turkish, Marathi and Hindi language and found that the performance on those languages often exceeded the performance of Bengali. In addition, a new high quality dataset is provided for text summarization evaluation. We, believe this research is a crucial addition to Bengali Natural Language Processing, that can easily be extended into other languages.

CCS Concepts: • **Do Not Use This Code → Generate the Correct Terms for Your Paper**; *Generate the Correct Terms for Your Paper*; Generate the Correct Terms for Your Paper; Generate the Correct Terms for Your Paper.

Additional Key Words and Phrases: Do, Not, Us, This, Code, Put, the, Correct, Terms, for, Your, Paper

## ACM Reference Format:

Fahim Morshed, Md. Abdur Rahman, and Sumon Ahmed. 2018. Extractive Text Summarization Using Word Similarity-based Spectral Clustering. *Proc. ACM Meas. Anal. Comput. Syst.* 37, 4, Article 111 (August 2018), 12 pages. <https://doi.org/XXXXXXX.XXXXXXX>

## 1 Introduction

Text Summarization is the process of shortening a larger text without losing any key information to increase the readability and save time for the reader. But manually summarizing very large texts is a counter-productive task due to it being more time consuming and tedious. So, developing an Automatic Text Summarization (ATS) method that can summarize larger texts reliably is really necessary to alleviate this manual labour [20]. Using ATS to summarize textual data is thus becoming very important in various fields such as news articles, legal documents, health reports, research papers, social media contents etc. ATS helps the reader to quickly and

---

Authors' Contact Information: Fahim Morshed, f.morshed.opee@gmail.com, Institute of Information Technology, University of Dhaka, Dhaka, Bangladesh; Md. Abdur Rahman, Centre for Advanced Research in Science, University of Dhaka, Dhaka, Bangladesh, mukul.arahman@gmail.com; Sumon Ahmed, sumon@du.ac.bd, Institute of Information Technology, University of Dhaka, Dhaka, Bangladesh.

---

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, or post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM 2476-1249/2018/8-ART111

<https://doi.org/XXXXXXX.XXXXXXX>

efficiently get the essential information without needing to read through large amounts of texts [6]. So, ATS is being utilized in various fields, from automatic news summarization, content filtering, and recommendation systems to assisting legal professionals in going through long documents. And researchers in reviewing academic papers by condensing vast amount of informations. It can also play a critical role in personal assistants and chatbots, providing condensed information to users quickly and efficiently [19].

There are two main types of ATS: extractive and abstractive [19]. Extractive summarization, which is the focus of this paper, works by selecting a subset from the source document, maintaining the original wording and sentence structure [14]. In contrast, abstractive summarization involves synthesising new text that reflects information from the input document but does not copy from it, similar to how a human summarizes a text [13]. Both of the method has their own advantage. The abstractive summarization can simulate the human language pattern very well thus increasing the natural flow and readability of the summary. But the extractive method requires much less computation than the abstractive method while also containing more key informations from the input [9].

The key approach to extractive summarization is implementing a sentence selection method to classify which sentences will belong in the summary. For this purpose, previously various simplistic ranking based methods were used to rank the sentences and identify the best sentences as the summary. These ranking methods used indexing [2], statistical [5] or Term Frequency-Inverse Document Frequency (TF-IDF) [4, 17, 18] based techniques to score the sentences and select the best scoring ones. But these methods fail to capture the semantic relationships between sentences of the input due to being simplistic in nature. To capture the semantic relationships between sentences, graph based extractive methods are effective due to the using of sentence similarity graph in their workflow [6]. Graph based methods represent the sentences as nodes of a graph, and the semantic similarity between two sentences as the edge between the nodes [14]. Popular graph based algorithms like LexRank [7] and TextRank [11] build graphs based on cosine similarity of the bag-of-word vectors. LexRank uses PageRank [15] method to score the sentences from the graph while TextRank uses random walk to determine which sentences are the most important to be in the summary. Graph-based methods like TextRank and LexRank offer a robust way to capture sentence importance and relationship, ensuring that the extracted summary covers the key information while minimizing redundancy [6].

Clustering-based approaches are a subset of graph-based approach to extractive text summarization. Here, sentences are grouped into clusters based on their semantic similarity to divide the document into topics, and one representative sentence from each cluster is chosen to form the summary [12]. Clustering reduces redundancy by ensuring that similar sentences are grouped together and only the most representative sentence is selected. This method is effective in summarization of documents with multiple topics or subtopics by picking sentences from each topic. An example of this method can be seen with COSUM [?] where the summarization is achieved using k-means clustering on the sentences and picking the most salient sentence from each cluster to compile in the final summary.

Despite the advancements of ATS in other languages, it remains an under-researched topic for Bengali due to Bengali being a low-resource language. Early attempts at Bengali text summarization relied on traditional methods like TF-IDF scoring to select the best scoring sentences to form the summary [1, 4, 17, 18]. These TF-IDF based approaches, while simple, faced challenges in capturing the true meaning of sentences. This is because TF-IDF based methods treats words as isolated terms resulting in synonyms of words being regarded as different terms [19]. To solve this problem, graph-based methods were introduced in Bengali to improve summarization quality by incorporating sentence similarity but they were still limited by the quality of word embeddings used

for the Bengali language. With the advent of word embedding models like FastText [8], it became possible to represent words in a vector space model, thus enabling more accurate sentence similarity calculations. However, existing models that use word embeddings, such as Sentence Average Similarity-based Spectral Clustering (SASbSC) method [16], encountered issues with sentence-similarity calculation when averaging word vectors to represent the meaning of a sentence with a vector. This method failed in most similarity calculation cases because words in a sentence are complementary to each other rather than being similar, leading to inaccurate sentence representations after averaging these different word vectors. As a result, word-to-word relationships between sentences get lost, reducing the effectiveness of the method.

In this paper, we propose a new clustering-based text summarization approach to address the challenge of calculating sentence similarity accurately. Our method improves upon previous attempts at graph-based summarization methods [3, 16] by focusing on the individual similarity between word pairs in sentences rather than averaging word vectors. We showed that the use of this novel approach greatly improved the accuracy, coverage and reliability of the output summaries due to having a deeper understanding of the semantic similarity between sentences. To calculate sentence similarity, we used the geometric mean of individual word similarities. The individual word similarities were achieved using Gaussian kernel function on a pair of corresponding word vector from each sentence. The word pairs are selected by finding the word vector with the smallest Euclidean distance from the target sentence. Thus, we get the semantic similarity between two sentences which can be used to build an affinity matrix to graphically represent the relationship between the sentences. This graph is clustered into groups to divide the document into distinct topics. One sentence from every cluster is selected to reduce redundancy and increase topic coverage. This method consistently outperforms other graph based text summarization methods such as BenSumm [3], LexRank [7], SASbSC [16] using four datasets on ROUGE metrics [10] as shown in Figure ?? and Table ?. This method performs well in other low resource languages also such as Hindi, Marathi and Turkish due to the language independent nature, as shown in the Table ?.

The main contributions of this paper are: (I) Proposed a new way to calculate similarity between two sentences. (II) Contributes a novel methodology for extractive text summarization for the Bengali language; by improving sentence similarity calculations and enhancing clustering techniques. (III) It offers a generalizable solution for creating less redundant and information rich summaries across languages. (IV) It provides a publicly available high quality dataset of 500 human generated summaries.

The rest of the paper is organized as follows: The Related works and Methodology are described in section ?? and ?? respectively. Section ?? illustrates the result of the performance evaluation for this work. Section ?? discusses the findings of the paper in more depth, and section ?? concludes the paper.

## 2 Template Overview

As noted in the introduction, the “acmart” document class can be used to prepare many different kinds of documentation — a double-anonymous initial submission of a full-length technical paper, a two-page SIGGRAPH Emerging Technologies abstract, a “camera-ready” journal article, a SIGCHI Extended Abstract, and more — all by selecting the appropriate *template style* and *template parameters*.

This document will explain the major features of the document class. For further information, the *L<sup>A</sup>T<sub>E</sub>X User’s Guide* is available from <https://www.acm.org/publications/proceedings-template>.

## 2.1 Template Styles

The primary parameter given to the “acmart” document class is the *template style* which corresponds to the kind of publication or SIG publishing the work. This parameter is enclosed in square brackets and is a part of the `\documentclass` command:

```
\documentclass[STYLE]{acmart}
```

Journals use one of three template styles. All but three ACM journals use the `acmsmall` template style:

- `acmsmall`: The default journal template style.
- `acmlarge`: Used by JOCCH and TAP.
- `acmtog`: Used by TOG.

The majority of conference proceedings documentation will use the `acmconf` template style.

- `sigconf`: The default proceedings template style.
- `sigchi`: Used for SIGCHI conference articles.
- `sigplan`: Used for SIGPLAN conference articles.

## 2.2 Template Parameters

In addition to specifying the *template style* to be used in formatting your work, there are a number of *template parameters* which modify some part of the applied template style. A complete list of these parameters can be found in the *L<sup>A</sup>T<sub>E</sub>X User’s Guide*.

Frequently-used parameters, or combinations of parameters, include:

- `anonymous, review`: Suitable for a “double-anonymous” conference submission. Anonymizes the work and includes line numbers. Use with the `\printID` command to print the submission’s unique ID on each page of the work.
- `authorversion`: Produces a version of the work suitable for posting by the author.
- `screen`: Produces colored hyperlinks.

This document uses the following string as the first command in the source file:

```
\documentclass[acmlarge]{acmart}
```

## 3 Modifications

Modifying the template — including but not limited to: adjusting margins, typeface sizes, line spacing, paragraph and list definitions, and the use of the `\vspace` command to manually adjust the vertical spacing between elements of your work — is not allowed.

**Your document will be returned to you for revision if modifications are discovered.**

## 4 Typefaces

The “acmart” document class requires the use of the “Libertine” typeface family. Your T<sub>E</sub>X installation should include this set of packages. Please do not substitute other typefaces. The “lmodern” and “ltimes” packages should not be used, as they will override the built-in typeface families.

## 5 Title Information

The title of your work should use capital letters appropriately - <https://capitalizemytitle.com/> has useful rules for capitalization. Use the `\title` command to define the title of your work. If your work has a subtitle, define it with the `\subtitle` command. Do not insert line breaks in your title.

If your title is lengthy, you must define a short version to be used in the page headers, to prevent overlapping text. The `\title` command has a “short title” parameter:

```
\title[short title]{full title}
```

## 6 Authors and Affiliations

Each author must be defined separately for accurate metadata identification. As an exception, multiple authors may share one affiliation. Authors' names should not be abbreviated; use full first names wherever possible. Include authors' e-mail addresses whenever possible.

Grouping authors' names or e-mail addresses, or providing an "e-mail alias," as shown below, is not acceptable:

```
\author{Brooke Aster, David Mehldau}
\email{dave,judy,steve@university.edu}
\email{firstname.lastname@phillips.org}
```

The `authornote` and `authornotemark` commands allow a note to apply to multiple authors — for example, if the first two authors of an article contributed equally to the work.

If your author list is lengthy, you must define a shortened version of the list of authors to be used in the page headers, to prevent overlapping text. The following command should be placed just after the last `\author{}` definition:

```
\renewcommand{\shortauthors}{McCartney, et al.}
```

Omitting this command will force the use of a concatenated list of all of the authors' names, which may result in overlapping text in the page headers.

The article template's documentation, available at <https://www.acm.org/publications/proceedings-template>, has a complete explanation of these commands and tips for their effective use.

Note that authors' addresses are mandatory for journal articles.

## 7 Rights Information

Authors of any work published by ACM will need to complete a rights form. Depending on the kind of work, and the rights management choice made by the author, this may be copyright transfer, permission, license, or an OA (open access) agreement.

Regardless of the rights management choice, the author will receive a copy of the completed rights form once it has been submitted. This form contains  $\LaTeX$  commands that must be copied into the source document. When the document source is compiled, these commands and their parameters add formatted text to several areas of the final document:

- the "ACM Reference Format" text on the first page.
- the "rights management" text on the first page.
- the conference information in the page header(s).

Rights information is unique to the work; if you are preparing several works for an event, make sure to use the correct set of commands with each of the works.

The ACM Reference Format text is required for all articles over one page in length, and is optional for one-page articles (abstracts).

## 8 CCS Concepts and User-Defined Keywords

Two elements of the "acmart" document class provide powerful taxonomic tools for you to help readers find your work in an online search.

The ACM Computing Classification System — <https://www.acm.org/publications/class-2012> — is a set of classifiers and concepts that describe the computing discipline. Authors can select entries from this classification system, via <https://dl.acm.org/ccs/ccs.cfm>, and generate the commands to be included in the  $\LaTeX$  source.

Table 1. Frequency of Special Characters

Non-English or Math	Frequency	Comments
Ø	1 in 1,000	For Swedish names
$\pi$	1 in 5	Common in math
\$	4 in 5	Used in business
$\Psi_1^2$	1 in 40,000	Unexplained usage

Table 2. Some Typical Commands

Command	A Number	Comments
\author	100	Author
\table	300	For tables
\table*	400	For wider tables

User-defined keywords are a comma-separated list of words and phrases of the authors' choosing, providing a more flexible way of describing the research being presented.

CCS concepts and user-defined keywords are required for all articles over two pages in length, and are optional for one- and two-page articles (or abstracts).

## 9 Sectioning Commands

Your work should use standard  $\LaTeX$  sectioning commands: section, subsection, subsubsection, and paragraph. They should be numbered; do not remove the numbering from the commands.

Simulating a sectioning command by setting the first word or words of a paragraph in boldface or italicized text is **not allowed**.

## 10 Tables

The “acmart” document class includes the “booktabs” package — <https://ctan.org/pkg/booktabs> — for preparing high-quality tables.

Table captions are placed *above* the table.

Because tables cannot be split across pages, the best placement for them is typically the top of the page nearest their initial cite. To ensure this proper “floating” placement of tables, use the environment **table** to enclose the table's contents and the table caption. The contents of the table itself must go in the **tabular** environment, to be aligned properly in rows and columns, with the desired horizontal and vertical rules. Again, detailed instructions on **tabular** material are found in the  *$\LaTeX$  User's Guide*.

Immediately following this sentence is the point at which Table 1 is included in the input file; compare the placement of the table here with the table in the printed output of this document.

To set a wider table, which takes up the whole width of the page's live area, use the environment **table\*** to enclose the table's contents and the table caption. As with a single-column table, this wide table will “float” to a location deemed more desirable. Immediately following this sentence is the point at which Table 2 is included in the input file; again, it is instructive to compare the placement of the table here with the table in the printed output of this document.

Always use `midrule` to separate table header rows from data rows, and use it only for this purpose. This enables assistive technologies to recognise table headers and support their users in navigating tables more easily.

## 11 Math Equations

You may want to display math equations in three distinct styles: inline, numbered or non-numbered display. Each of the three are discussed in the next sections.

### 11.1 Inline (In-text) Equations

A formula that appears in the running text is called an inline or in-text formula. It is produced by the **math** environment, which can be invoked with the usual `\begin . . . \end` construction or with the short form `$ . . . $`. You can use any of the symbols and structures, from  $\alpha$  to  $\omega$ , available in  $\text{\LaTeX}$  [? ]; this section will simply show a few examples of in-text equations in context. Notice how this equation:  $\lim_{n \rightarrow \infty} x = 0$ , set here in in-line math style, looks slightly different when set in display style. (See next section).

### 11.2 Display Equations

A numbered display equation—one set off by vertical space from the text and centered horizontally—is produced by the **equation** environment. An unnumbered display equation is produced by the **displaymath** environment.

Again, in either environment, you can use any of the symbols and structures available in  $\text{\LaTeX}$ ; this section will just give a couple of examples of display equations in context. First, consider the equation, shown as an inline equation above:

$$\lim_{n \rightarrow \infty} x = 0 \tag{1}$$

Notice how it is formatted somewhat differently in the **displaymath** environment. Now, we'll enter an unnumbered equation:

$$\sum_{i=0}^{\infty} x + 1$$

and follow it with another numbered equation:

$$\sum_{i=0}^{\infty} x_i = \int_0^{\pi+2} f \tag{2}$$

just to demonstrate  $\text{\LaTeX}$ 's able handling of numbering.

## 12 Figures

The “figure” environment should be used for figures. One or more images can be placed within a figure. If your figure contains third-party material, you must clearly identify it as such, as shown in the example below.

Your figures should contain a caption which describes the figure to the reader.

Figure captions are placed *below* the figure.

Every figure should also have a figure description unless it is purely decorative. These descriptions convey what's in the image to someone who cannot see it. They are also used by search engine crawlers for indexing images, and when images cannot be loaded.

A figure description must be unformatted plain text less than 2000 characters long (including spaces). **Figure descriptions should not repeat the figure caption – their purpose is to capture important information that is not already provided in the caption or the main text of the paper.** For figures that convey important and complex new information, a short text description may not be adequate. More complex alternative descriptions can be placed in an appendix and referenced in a short figure description. For example, provide a data table





Fig. 1. 1907 Franklin Model D roadster. Photograph by Harris & Ewing, Inc. [Public domain], via Wikimedia Commons. (<https://goo.gl/VLCRBB>).

capturing the information in a bar chart, or a structured list representing a graph. For additional information regarding how best to write figure descriptions and why doing this is so important, please see <https://www.acm.org/publications/taps/describing-figures/>.

## 12.1 The “Teaser Figure”

A “teaser figure” is an image, or set of images in one figure, that are placed after all author and affiliation information, and before the body of the article, spanning the page. If you wish to have such a figure in your article, place the command immediately before the `\maketitle` command:

```
\begin{teaserfigure}
  \includegraphics[width=\textwidth]{sampleteaser}
  \caption{figure caption}
  \Description{figure description}
```



```
\end{teaserfigure}
```

### 13 Citations and Bibliographies

The use of Bib<sub>T</sub><sub>E</sub>X for the preparation and formatting of one's references is strongly recommended. Authors' names should be complete — use full first names (“Donald E. Knuth”) not initials (“D. E. Knuth”) — and the salient identifying features of a reference should be included: title, year, volume, number, pages, article DOI, etc.

The bibliography is included in your source document with these two commands, placed just before the `\end{document}` command:

```
\bibliographystyle{ACM-Reference-Format}
\bibliography{bibfile}
```

where “bibfile” is the name, without the “.bib” suffix, of the Bib<sub>T</sub><sub>E</sub>X file.

Citations and references are numbered by default. A small number of ACM publications have citations and references formatted in the “author year” style; for these exceptions, please include this command in the **preamble** (before the command “`\begin{document}`”) of your  $\text{\LaTeX}$  source:

```
\citestyle{acmauthoryear}
```

Some examples. A paginated journal article [? ], an enumerated journal article [? ], a reference to an entire issue [? ], a monograph (whole book) [? ], a monograph/whole book in a series (see 2a in spec. document) [? ], a divisible-book such as an anthology or compilation [? ] followed by the same example, however we only output the series if the volume number is given [? ] (so Editor00a's series should NOT be present since it has no vol. no.), a chapter in a divisible book [? ], a chapter in a divisible book in a series [? ], a multi-volume work as book [? ], a couple of articles in a proceedings (of a conference, symposium, workshop for example) (paginated proceedings article) [? ? ], a proceedings article with all possible elements [? ], an example of an enumerated proceedings article [? ], an informally published work [? ], a couple of preprints [? ? ], a doctoral dissertation [? ], a master's thesis: [? ], an online document / world wide web resource [? ? ? ], a video game (Case 1) [? ] and (Case 2) [? ] and [? ] and (Case 3) a patent [? ], work accepted for publication [? ], 'YYYYb'-test for prolific author [? ] and [? ]. Other cites might contain 'duplicate' DOI and URLs (some SIAM articles) [? ]. Boris / Barbara Beeton: multi-volume works as books [? ] and [? ]. A couple of citations with DOIs: [? ? ]. Online citations: [? ? ? ]. Artifacts: [? ] and [? ].

### 14 Acknowledgments

Identification of funding sources and other support, and thanks to individuals and groups that assisted in the research and the preparation of the work should be included in an acknowledgment section, which is placed just before the reference section in your document.

This section has a special environment:

```
\begin{acks}
...
\end{acks}
```

so that the information contained therein can be more easily collected during the article metadata extraction phase, and to ensure consistency in the spelling of the section heading.

Authors should not prepare this section as a numbered or unnumbered `\section`; please use the “acks” environment.

## 15 Appendices

If your work needs an appendix, add it before the “`\end{document}`” command at the conclusion of your source document.

Start the appendix with the “`appendix`” command:

```
\appendix
```

and note that in the appendix, sections are lettered, not numbered. This document has two appendices, demonstrating the section and subsection identification method.

## 16 Multi-language papers

Papers may be written in languages other than English or include titles, subtitles, keywords and abstracts in different languages (as a rule, a paper in a language other than English should include an English title and an English abstract). Use `language=...` for every language used in the paper. The last language indicated is the main language of the paper. For example, a French paper with additional titles and abstracts in English and German may start with the following command

```
\documentclass[sigconf, language=english, language=german,
               language=french]{acmart}
```

The title, subtitle, keywords and abstract will be typeset in the main language of the paper. The commands `\translatedXXX`, `XXX` begin title, subtitle and keywords, can be used to set these elements in the other languages. The environment `translatedabstract` is used to set the translation of the abstract. These commands and environment have a mandatory first argument: the language of the second argument. See `sample-sigconf-i13n.tex` file for examples of their usage.

## 17 SIGCHI Extended Abstracts

The “`sigchi-a`” template style (available only in  $\LaTeX$  and not in Word) produces a landscape-orientation formatted article, with a wide left margin. Three environments are available for use with the “`sigchi-a`” template style, and produce formatted output in the margin:

**sidebar:** Place formatted text in the margin.

**marginfigure:** Place a figure in the margin.

**marginfigure:** Place a table in the margin.

## Acknowledgments

To Robert, for the bagels and explaining CMYK and color spaces.

## References

- [1] Sumya Akter, Aysa Siddika Asa, Md. Palash Uddin, Md. Delowar Hossain, Shikhor Kumer Roy, and Masud Ibn Afjal. 2017. 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)*. <https://doi.org/10.1109/ICIVPR.2017.7890883>
- [2] P. B. Baxendale. 1958. Machine-Made Index for Technical Literature—An Experiment. *IBM Journal of Research and Development* 2, 4 (1958), 354–361. <https://doi.org/10.1147/rd.24.0354>
- [3] Radia Rayan Chowdhury, Mir Tafseer Nayeem, Tahsin Tasnim Mim, Md. Saifur Rahman Chowdhury, and Taufiqul Jannat. 2021. Unsupervised Abstractive Summarization of Bengali Text Documents. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*. Association for Computational Linguistics, 2612–2619. <https://doi.org/10.18653/v1/2021.eacl-main.224>
- [4] Satya Ranjan Dash, Pubali Guha, Debasish Kumar Mallick, and Shantipriya Parida. 2022. Summarizing Bengali Text: An Extractive Approach. In *Intelligent Data Engineering and Analytics*. Springer Nature Singapore, 133–140.
- [5] H. P. Edmundson. 1969. New Methods in Automatic Extracting. *J. ACM* 16, 2 (apr 1969), 264–285. <https://doi.org/10.1145/321510.321519>

- [6] Wafaa S. El-Kassas, Cherif R. Salama, Ahmed A. Rafea, and Hoda K. Mohamed. 2021. Automatic text summarization: A comprehensive survey. *Expert Systems with Applications* 165 (2021), 113679. <https://doi.org/10.1016/j.eswa.2020.113679>
- [7] Günes Erkan and Dragomir R. Radev. 2004. LexRank: graph-based lexical centrality as salience in text summarization. *J. Artif. Int. Res.* 22, 1 (dec 2004), 457–479.
- [8] Edouard Grave, Piotr Bojanowski, Prakhar Gupta, Armand Joulin, and Tomas Mikolov. 2018. Learning Word Vectors for 157 Languages. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*. European Language Resources Association (ELRA). <https://aclanthology.org/L18-1550>
- [9] Vishal Gupta and Gurpreet Lehal. 2010. A Survey of Text Summarization Extractive Techniques. *Journal of Emerging Technologies in Web Intelligence* 2 (08 2010). <https://doi.org/10.4304/jetwi.2.3.258-268>
- [10] Chin-Yew Lin. 2004. ROUGE: A Package for Automatic Evaluation of Summaries. In *Text Summarization Branches Out*. Association for Computational Linguistics, 74–81. <https://aclanthology.org/W04-1013>
- [11] Rada Mihalcea and Paul Tarau. 2004. TextRank: Bringing Order into Text. In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 404–411. <https://aclanthology.org/W04-3252>
- [12] G. Bharathi Mohan and R. Prasanna Kumar. 2022. A Comprehensive Survey on Topic Modeling in Text Summarization. In *Micro-Electronics and Telecommunication Engineering*. Springer Nature Singapore, 231–240.
- [13] N. Moratanch and S. Chitrakala. 2016. A survey on abstractive text summarization. In *2016 International Conference on Circuit, Power and Computing Technologies (ICCPCT)*. 1–7. <https://doi.org/10.1109/ICCPCT.2016.7530193>
- [14] N. Moratanch and S. Chitrakala. 2017. A survey on extractive text summarization. In *2017 International Conference on Computer, Communication and Signal Processing (ICCCSP)*. 1–6. <https://doi.org/10.1109/ICCCSP.2017.7944061>
- [15] Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. 1999. The PageRank Citation Ranking : Bringing Order to the Web. In *The Web Conference*. <https://api.semanticscholar.org/CorpusID:1508503>
- [16] Sohini Roychowdhury, Kamal Sarkar, and Arka Maji. 2022. Unsupervised Bengali Text Summarization Using Sentence Embedding and Spectral Clustering. In *Proceedings of the 19th International Conference on Natural Language Processing (ICON)*. Association for Computational Linguistics, 337–346. <https://aclanthology.org/2022.icon-main.40>
- [17] Kamal Sarkar. 2012. An approach to summarizing Bengali news documents. In *Proceedings of the International Conference on Advances in Computing, Communications and Informatics (ICACCI '12)*. Association for Computing Machinery, 857–862. <https://doi.org/10.1145/2345396.2345535>
- [18] Kamal Sarkar. 2012. Bengali text summarization by sentence extraction. *CoRR abs/1201.2240* (2012). <http://arxiv.org/abs/1201.2240>
- [19] Oguzhan Tas and Farzad Kiyani. 2017. A SURVEY AUTOMATIC TEXT SUMMARIZATION. *PressAcademia Procedia* 5, 1 (2017), 205–213. <https://doi.org/10.17261/Pressacademia.2017.591>
- [20] Adhika Pramita Widyassari, Supriadi Rustad, Guruh Fajar Shidik, Edi Noersasongko, Abdul Syukur, Affandy Affandy, and De Rosal Ignatius Moses Setiadi. 2022. Review of automatic text summarization techniques & methods. *Journal of King Saud University - Computer and Information Sciences* 34, 4 (2022), 1029–1046. <https://doi.org/10.1016/j.jksuci.2020.05.006>

## A Research Methods

### A.1 Part One

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Morbi malesuada, quam in pulvinar varius, metus nunc fermentum urna, id sollicitudin purus odio sit amet enim. Aliquam ullamcorper eu ipsum vel mollis. Curabitur quis dictum nisl. Phasellus vel semper risus, et lacinia dolor. Integer ultricies commodo sem nec semper.

### A.2 Part Two

Etiam commodo feugiat nisl pulvinar pellentesque. Etiam auctor sodales ligula, non varius nibh pulvinar semper. Suspendisse nec lectus non ipsum convallis congue hendrerit vitae sapien. Donec at laoreet eros. Vivamus non purus placerat, scelerisque diam eu, cursus ante. Etiam aliquam tortor auctor efficitur mattis.

## B Online Resources

Nam id fermentum dui. Suspendisse sagittis tortor a nulla mollis, in pulvinar ex pretium. Sed interdum orci quis metus euismod, et sagittis enim maximus. Vestibulum gravida massa ut felis suscipit congue. Quisque mattis elit a risus ultrices commodo venenatis eget dui. Etiam sagittis eleifend elementum.

Nam interdum magna at lectus dignissim, ac dignissim lorem rhoncus. Maecenas eu arcu ac neque placerat aliquam. Nunc pulvinar massa et mattis lacinia.

Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009