

1
2
3
4
5
6
7
8
9
10

A Novel Approach To Text Summarization Using Machine Learning

ABSTRACT

Text summarization is a key strategy in the domains of information retrieval and natural language processing (NLP). Its objective is to reduce a lengthy written document into a clearer, more succinct summary of the information it contains. When a text document is too lengthy or intricate to analyse completely, as in news stories, academic papers, or legal documents, this approach is extremely helpful. The major challenge of text summarising is to take the most important and relevant information from the original text and convey it in an understandable and concise way. In this study, extractive and abstractive summarising techniques are the two primary categories of text summary methods. The paper also presents several algorithms that have been proposed for text summarization, including TextRank, Seq2Seq, and BART. TextRank is a simple and fast algorithm that works well for short documents, Seq2Seq is a deep learning-based approach that generates high-quality summaries, and BART is a transformer-based algorithm that provides the best results on benchmark datasets. The obtained ROUGE Score after passing TextRank, BART, and Seq2Seq algorithm significant also..

11
12
13
14
15
16
17
18
19
20

Keywords: Text Summarization, Machine Learning, BART, TextRank, Text analysis, NLP

1. INTRODUCTION

Natural language processing (NLP), a discipline of artificial intelligence and computer science, is the study of how computers and people communicate using natural language. The objective of NLP is to create models and algorithms that can process and comprehend human language, making it possible for computers to carry out activities like sentiment analysis, text categorization, machine translation, and text summarization, among others.

21
22
23
24
25

NLP is a complex and challenging field, as human language is rich, ambiguous, and constantly evolving. To address these challenges, NLP relies on a variety of techniques, including statistical models, machine learning, deep learning, and rule-based systems. These techniques are used to process and analyze text data, extract meaningful information, and generate natural language outputs.

26
27
28
29
30
31
32
33

You can reduce a lengthy written document into a shorter, more manageable representation of its content by using text summarization. There are many different approaches to text summarising, such as extractive and abstractive techniques [24] [25]. The summary is created using extractive summarization, which pulls out the key phrases or clauses from the original text. Contrarily, abstractive summarization starts from scratch with a new summary depending on the content of the original text. The best strategy to use relies on the particular needs of the summarising work because both techniques have advantages and disadvantages [26].

34 Text summarization has numerous applications, including information retrieval, content-
35 based recommendation systems, and knowledge management. In the field of information
36 retrieval, text summarization can be used to provide a quick overview of search results,
37 allowing users to quickly identify relevant documents. In content-based recommendation
38 systems, text summarization can be used to provide users with a brief description of the
39 content of recommended documents. Finally, in knowledge management, text
40 summarization can be used to extract key information from large and complex documents,
41 making it easier to store, manage, and retrieve information [27].

42 Text summarising is a difficult topic since it calls for the capacity to extract the most crucial
43 information from the source text and to display it in a clear and simple manner.[16] A good
44 summarization should accurately capture the main ideas and information contained in the
45 original text, while excluding irrelevant or redundant information. To do this, text
46 summarising algorithms must have a thorough grasp of the architecture and content of the
47 original text and be able to recognise the most crucial clauses or phrases based on their
48 significance, applicability, and coherence [28].

49 For text summarization, a number of methods have been put forth, including TextRank,
50 Seq2Seq, and BART. The text document's structure is used by the graph-based algorithm
51 TextRank to obtain a summary. It models the text as a graph of words and sentences, and
52 calculates the importance of each node based on the number and weight of its edges.
53 Seq2Seq is a deep learning-based algorithm that uses the encoder-decoder architecture to
54 generate a summary. The transformer-based algorithm BART (Bidirectional Encoder
55 Representations from Transformers) processes the input text in a bidirectional manner.
56

57 **2. LITERATURE REVIEW**

58 This section cites earlier works that make use of the various summarizing methods. Instead
59 of sentence production for text summary, the majority of researches focus on sentence
60 extraction. The most popular approach of summarization creates extractive summaries
61 based on statistical aspects of the sentence.

62 According to Luhn[4], the words that are used the most often in a text correspond to its most
63 crucial ideas. He wanted to assess each phrase based on the frequency of each word before
64 selecting the best result. Methods based on location, title, and cue words were suggested by
65 Edmunson[16]. He argued that the summary should include the topic information, which is
66 usually found in the opening few words or paragraphs of a text. One flaw in the statistical
67 technique is that it ignores the semantic relationships between sentences. In order to give a
68 summary, Goldstein [2] developed a query-based summarizing technique that would extract
69 important lines from a text according to the query fired. There is a suggested query for the
70 extraction criterion. The more words combined in the question and a sentence, the more
71 likely it is to be included in a summary. Goldstein[2][1] used statistical and linguistic
72 characteristics to analyze the summaries of news stories in order to assess the phrases in
73 the document. One approach to summarizing is sentence extraction and grouping. In order
74 to determine how similar sentences' cumulative phrases are, sentences should first be
75 clustered depending on how far apart they are from one another semantically, according to
76 ZHANG Pei-ying and LI Cun[5]. Finally, the sentences should be selected using extraction
77 procedures. K-means algorithm is used to group the sentences together[5]. Morris and
78 Hirst[9][7] were the authors who initially developed the idea of lexical chains. Lexical chains
79 [7] take advantage of the relationships between any number of related words. By assembling
80 groups of semantically similar words, we can form lexical chains. Barzilay and Elhadad[8][6]
81 built a lexical chain by utilising WordNet to determine the semantic distance between terms.
82 The phrases associated with the chosen strong lexical chains are picked as a summary.

83 Using a liner time method, H. Gregory Silber and McCoy [10] created a method for creating
 84 lexical chains. By creating an intermediate representation, the author follows Barzilay and
 85 Elhadad's [6] approach of using lexical chains to extract crucial concepts from the original
 86 text. The method for using lexical chaining to construct an array of Meta-Chains whose size
 87 is equal to the number of noun senses in the Word Net and the document is detailed in the
 88 article [10]. Proper nouns and anaphora resolution issues with the algorithm needs to be
 89 fixed. A alternative approach to summarization is found in graph theory [11]. To create a
 90 semantic network of the original text, the author suggested a technique based on subject-
 91 object-predicate (SOP) triples from individual phrases. Every word has essential, significant
 92 concepts strewn throughout it. According to the author [11], by identifying and using the links
 93 between them, it may be feasible to recover crucial information. Pushpak Bhattacharyya [12]
 94 of the IIT Bombay, one of the researchers, proposed a Word Net-based approach for
 95 summarizing. Word-net is used to summarise the document, creating a sub-graph. The
 96 Word Net is used to assign weights to the sub-graph's nodes in respect to the synsnet. The
 97 most popular methods for text summarization incorporate one or both of the linguistic and
 98 statistical approaches.

99
 100 **3. ABOUT DATASET**

101 The CNN/Daily Mail dataset is a sizable corpus of news stories and summaries gathered for
 102 summarizing purposes. It has been frequently used to train and test summarizing models
 103 and has over 300,000 article-summary pairs. This dataset includes articles and summaries
 104 on a variety of subjects, such as politics, entertainment, sports, and more. The articles and
 105 summaries are taken from the CNN and Daily Mail news websites and other some other
 106 sources [17]. The summaries in this dataset are written by professional journalists and are
 107 typically shorter than the corresponding articles, making them ideal for training
 108 summarization models. The CNN/Daily Mail dataset has been utilized in a wide range of
 109 academic projects and has significantly advanced the field of text summarization.

110 **4. DATA PREPROCESSING**

111 Preprocessing is an important step in working with the CNN/Daily Mail dataset. The main
 112 objective of preprocessing is to clean and transform the raw data into a suitable presentation
 113 for further analysis or modeling. Here are some common preprocessing steps that are
 114 typically applied to the CNN/Daily Mail dataset:

115 **4.1 Data Cleaning-** This step involves removing any irrelevant or redundant information and
 116 handling any missing or incomplete data. This can include removing stop words, stemming,
 117 and lemmatizing the text, as well as removing any irrelevant characters or symbols. Figure 1
 118 highlights elements of the Dataset.

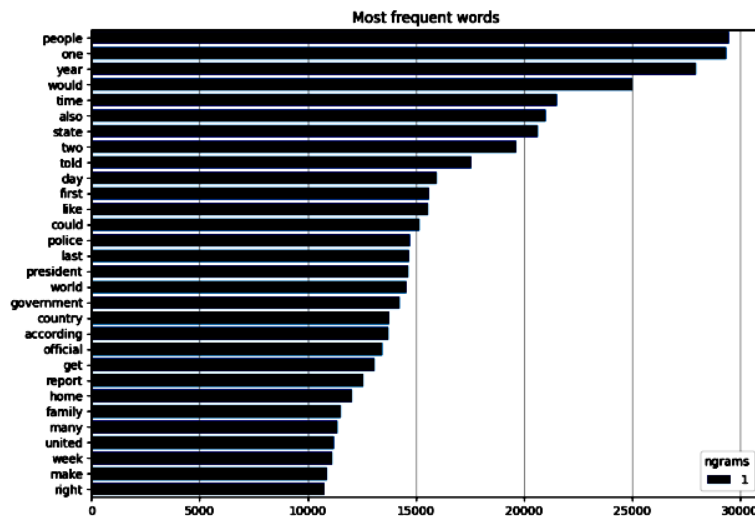
| | text | y |
|---|---------------------------------------------------|---------------------------------------------------|
| 0 | LONDON, England (Reuters) – Harry Potter star... | Harry Potter star Daniel Radcliffe gets £20M f... |
| 1 | Editor's note: In our Behind the Scenes series... | Mentally ill inmates in Miami are housed on th... |
| 2 | MINNEAPOLIS, Minnesota (CNN) – Drivers who we... | NEW: "I thought I was going to die," driver sa... |
| 3 | WASHINGTON (CNN) – Doctors removed five small... | Five small polyps found during procedure; "non... |
| 4 | (CNN) – The National Football League has ind... | NEW: NFL chief, Atlanta Falcons owner critical... |

119
 120 **Fig. 1. Elements of the Dataset**

121
 122 **4.2 Tokenization-** Tokenization comprises breaking down the text into reduced parts such
 123 as words or sentences. This is typically done using a tokenizer that can handle different
 124 types of text, such as punctuation, numbers, and special characters.

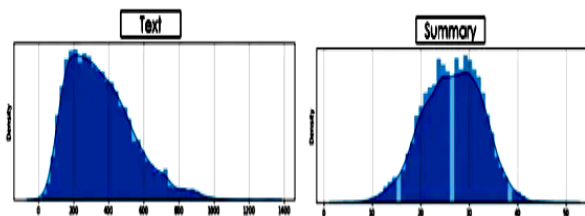
125 **4.3 Text Normalization-** Text normalization encompasses transforming the text into a
 126 standardized format [18]. This can comprise altering text to lowercase, removing diacritics,

127 and converting contractions to their full forms. Figure 2 highlights the Word Frequency of the
128 Dataset.



129
130 **Fig. 2. Word frequency**

131
132 **4.4 POS Tagging-** Part-of-speech (POS) tagging entails assigning the appropriate part of
133 speech to each word in the text. This can aid in determining the text's syntactic structure and
134 be helpful for later tasks like sentiment analysis and named entity recognition. Figure 3
135 highlights the Length Analysis of the Dataset.



136
137 **Fig. 3. Length analysis.**

138 139 **5. TYPES OF SUMMARIZATION TECHNIQUES**

140
141 **5.1 Extraction-based Summarization-** The process of extraction-based summarization is
142 locating and extracting the key expressions or sentences from the input text in order to
143 produce a summary. The ability to retain the text's original words and meaning is one of the
144 key benefits of extraction-based summarization, which might be crucial in certain situations
145 [19]. Extraction-based summarization can be useful for summarizing text that is mostly
146 factual in nature, like news items, and is also rather simple to put into practice.

147 Implementing extraction-based summarization can be done in several ways, including
148 frequency-based methods and machine learning-based methods. The most significant
149 sentences are determined using statistical criteria such as word frequency or sentence
150 length using frequency-based approaches. Machine learning-based approaches entail
151 building a model that can recognise the most significant sentences based on a variation of
152 characteristics, as well as sentence length, placement in the text, and the presence of

153 essential words or phrases. Because the extracted sentences could not flow naturally
154 together, extraction-based summarising has the potential to yield summaries that lack
155 coherence and organisation. Moreover, more complicated texts that demand in-depth
156 comprehension and interpretation may be difficult for extraction-based summarization to
157 handle.

158 **5.2 Abstraction-based Summarization-** A technique called abstraction-based
159 summarization includes creating a summary that does not have to match the text's exact
160 language but rather captures the important ideas and concepts in a broader sense. To do
161 this, natural language generation techniques are used to generate new sentences that more
162 clearly and concisely express the major concepts of the given text. Technical or scientific
163 texts are highly suited for abstraction-based summarization since it may capture more
164 complicated ideas and links between concepts [20]. Abstraction-based summarization has
165 the drawback of possibly requiring more training data and computer resources than
166 extraction-based summarization. Also, the effectiveness of the natural language generation
167 techniques used, which can be difficult to optimize, may have a significant impression on the
168 quality of the summary.

169 **6. ALGORITHM SELECTION**

170

171 The algorithm selection section is a crucial part of the text summarization project as it
172 determines the performance of the model. The primary objective of this section is to gauge
173 different machine learning algorithms and select the best performing one for the given task.
174 In this section, we discuss the various models we considered and the criteria we used to
175 select the best one.

176 •TextRank: Preprocessing the input text by removing stop words, stemming, and lower-
177 casing the text. Imagine the text as a network, with the nodes standing in for sentences and
178 the edges signifying how similar they are. using the graph's PageRank algorithm to isolate
179 the summary's most crucial phrases.

180 •Seq2Seq: Preprocessing the input and target text by tokenizing and converting them into a
181 numerical representation. Creating a sequence-to-sequence model with an encoder-decoder
182 architecture and attention mechanism. Training the model on a dataset of input and target
183 summaries. Evaluating the model performance by calculating metrics such as ROUGE and
184 BLEU.

185 •BART: Preprocessing the input and target text by tokenizing and converting them into a
186 numerical representation. Creating a BART model with an encoder-decoder architecture and
187 fine-tuning it on a dataset of input and target summaries. Evaluating the model performance
188 by calculating metrics such as ROUGE and BLEU.

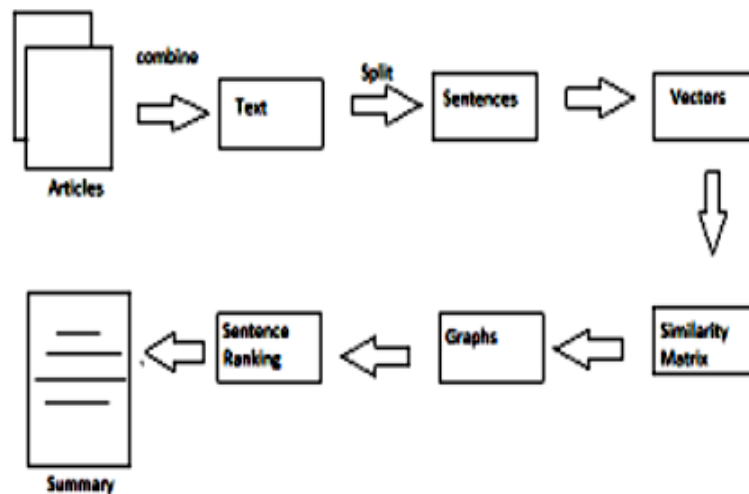
189 Overall, the methodology for developing text summarization using TextRank, Seq2Seq, and
190 BART involves a combination of data processing, model selection, training, evaluation, and
191 optimization to achieve the desired level of accuracy and performance.

192 **7. RESEARCH BACKGROUND**

193 TextRank - The fundamental concept underlying TextRank is to visualize the text as a graph,
194 where each node resembles to a phrase or a word, and the connections amongst nodes are
195 represented by the edges. The PageRank algorithm determines the ranking of the nodes,
196 and edges are weighted according to how similar the nodes they connect are. The input text
197 is first pre-processed to eliminate stop words, punctuation, and other noise before being
198 used to produce the graph. The edges connecting nodes are then determined based on how

199 semantically similar the sentences or words are to one another [21]. Cosine similarity,
200 Jaccard similarity, or other metrics can be used to determine how similar two nodes are to
201 one another.

202 The PageRank algorithm is used to assign a ranking to each node after the graph has been
203 created. The PageRank algorithm is altered for TextRank to consider how similar nodes and
204 their neighbors are. In more detail, a node's ranking is established using the average of the
205 rankings of its neighbors, weighted by how similar the nodes are to one another. The graph's
206 top-ranked nodes are selected to create the summary for text summarization. This can be
207 accomplished by choosing the top-ranked expressions or words to create a summary that
208 encapsulates the text's core concepts. Figure 4 highlights the Flow Diagram of the
209 TextRank.



210
211 **Fig.4 Flow Diagram of the TextRank**

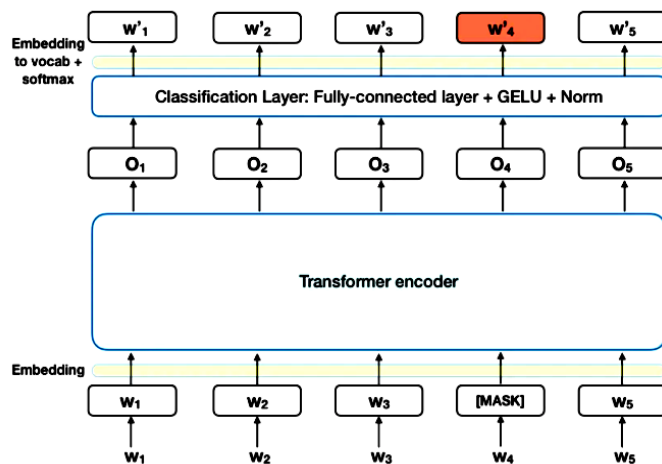
212
213 Seq2Seq - Seq2Seq models' central tenet is to discover a mapping between sequences of
214 input and output data, such as a source language's word order and a target language's word
215 order. Encoders and decoders are the two primary parts of Seq2Seq models. The input
216 sequence must be processed and encoded into a fixed-length vector representation by the
217 encoder. The decoder then receives this vector and uses the encoded input and previous
218 outputs to construct the output sequence, one token at a time.

219 A recurrent neural network (RNN), such as an LSTM network or a gated recurrent unit
220 (GRU) network, serves as the encoder in most cases. At each time step, the RNN updates
221 its hidden state as it goes over the input sequence, single token at a time. The input
222 sequence is represented in encoded form by the RNN's final hidden state. Although it often
223 has a different design than the encoder, the decoder is likewise an RNN [22]. It generates
224 the subsequent token in the output sequence using the encoded representation of the input
225 sequence as well as the previous token that was generated as input. Up until it encounters
226 an end-of-sequence token or a predetermined maximum length, the decoder keeps
227 producing output tokens.

228 BART - On the transformer architecture, BART is based. Transformers are a type of neural
229 network that models the connections between various elements of a sequence, such as the
230 words in a sentence, by using attention mechanisms. An auto-regressive decoder and a
231 bidirectional encoder are added as part of BART's extension of the transformer design. The

232 Seq2Seq model's encoder and the bidirectional encoder are comparable. It converts a string
 233 of tokens, like the words in a phrase, into a fixed-length vector representation as input. The
 234 BART encoder, in contrast to a typical Seq2Seq encoder, is bidirectional, which means that
 235 it processes the input sequence both forward and backward [23]. It has been demonstrated
 236 that doing so enhances performance on NLP tasks by enabling the encoder to record more
 237 intricate relationships between the tokens.

238 A Seq2Seq model's decoder and BART's auto-regressive decoder are comparable. It
 239 produces the subsequent token in the output sequence with the encoded representation of
 240 the input sequence as well as the previous tokens that were generated. The BART decoder,
 241 in contrast to a typical Seq2Seq decoder, is auto-regressive, which means that it generates
 242 the output tokens one at a time dependent on the tokens that were generated earlier. As a
 243 result, the decoder can recognise dependencies between output tokens and produce text
 244 that is more fluid and cohesive. A denoising autoencoder aim is used to pre-train BART
 245 using a sizable corpus of text data. By randomly masking words or rearranging expressions
 246 in the input text, the pre-training process tampers with the text before training the model to
 247 restore the original content. The model gains a broad understanding of natural language
 248 from this pre-training target, which it can then hone for particular NLP tasks. Figure 5
 249 highlights the Flow Diagram of the BART.



250
 251 **Fig.5 Flow Diagram of the BART**

252
 253 **8. RESULTS**

254 The results has been summarized below using TextRank, Bart, and Seq2Sq algorithms.

After almost dissipating on august 10, irene peaked as a category 2 hurricane on august 16. Irene persisted for 14 days as a tropical system, the longest duration of any storm of the 2005 season. Irene lasted for 14 days as a tropical system, the longest duration of any storm of the 2005 season. Hurricane irene began as a cape verde storm. Some of the models predicted that irene would make landfall in north carolina, while others continued to anticipate that irene would dissipate. hurricane irene was a long-lived cape verde hurricane during the 2005 atlantic hurricane season. A vigorous tropical wave moved off the west coast of africa on august 1, initially weakening due to cooler sea surface temperatures. However, the hurricane generated strong waves and increased the risk of rip currents along the east coast of the United States.

255
 256 **Fig 6. Predicted Summary comparison using TextRank**

257

258 Figure 6 highlights the comparison between Predicted Summary and Real Summary Using
259 TextRank.

rouge1: 0.21 | rouge2: 0.07 | rougeL: 0.07 --> avg rouge: 0.15

260
261 **Fig 7. ROUGE Score of TextRank**

262
263 Figure 7 highlights the Rouge Score for the TextRank Algorithm. ROUGE basically measure
264 the similarity between machine-generated summaries and human-written reference
265 summaries.

After almost dissipating on august 10, irene peaked as a category 2 hurricane on
august 16. irene persisted for 14 days as a tropical system, the longest duration of any
storm of the 2005 season. irene lasted for 14 days as a tropical system, the longest
duration of any storm of the 2005 season. The storm formed near cape verde on
august 4 and crossed the atlantic, turning northward around bermuda before being
absorbed by an extratropical cyclone while situated southeast of newfoundland. irene
entered a region of increased wind shear and began to weaken, and as a result it was
downgraded to a tropical storm early on august 18, when it was 520 miles (830 km)

266
267 **Fig 8. Predicted Summary comparison using Seq2Seq**

268
269 Figure 8 highlights the comparison between Predicted Summary and Real Summary Using
270 Seq2seq Algorithm

rouge1: 0.08 | rouge2: 0.01 | rougeL: 0.01 --> avg rouge: 0.05

271
272 **Fig 9. ROUGE Score of Seq2Seq**

273
274 Figure 9 highlights the Rouge Score for the Seq2Seq Algorithm.

After almost dissipating on august 10, irene peaked as a category 2 hurricane on
august 16. irene persisted for 14 days as a tropical system, the longest duration of any
storm of the 2005 season. irene lasted for 14 days as a tropical system, the longest
duration of any storm of the 2005 season. The storm formed near cape verde on
august 4 and crossed the atlantic, turning northward around bermuda before being
absorbed by an extratropical cyclone while situated southeast of newfoundland. irene
entered a region of increased wind shear and began to weaken, and as a result it was
downgraded to a tropical storm early on august 18, when it was 520 miles (830 km)

275
276 **Fig 10. Passage with Predicted Summary comparison using TextRank**

277
278 Figure 10 highlights the comparison between Predicted Summary and passage from dataset
279 Using BART algorithm.

rouge1: 0.3 | rouge2: 0.11 | rougeL: 0.11 --> avg rouge: 0.21

280
281 **Fig 11. ROUGE Score of BART**

282
283 Figure 11 highlights the Rouge Score for the BART Algorithm.

284 **9. CONCLUSION AND FUTURE SCOPE**

285 In conclusion, text summarizing is a crucial task in natural language processing that seeks to
286 produce a short and cohesive summary of a lengthy text. Text summarizing can be done in a
287 variety of ways, including with abstraction- and extraction-based methods. But recent

288 developments in deep learning models, including TextRank, Seq2Seq, and BART, have
289 produced encouraging outcomes for text summarization. A graph-based system called
290 TextRank employs the PageRank algorithm to determine the key sentences in a file and
291 produce a summary. To create abstractive summaries, Seq2Seq models—neural network
292 models—can be trained on pairs of documents and summaries. Contrarily, BART is a pre-
293 trained transformer model that has attained cutting-edge performance on a variation of NLP
294 applications, including text analysis. These methods each have advantages and
295 disadvantages of their own. TextRank is an easy-to-use technique that works well without
296 training data, although it might have trouble producing summaries that are fluent and
297 cohesive. Seq2Seq models have the ability to produce abstractive summaries that are fluent
298 and seem more natural than extractive approaches, but they need a lot of training data and
299 can occasionally produce generic summaries. Although BART is a strong and flexible model
300 that can produce abstract as well as cohesive summaries, it is computationally expensive
301 and necessitates pre-training on a substantial amount of data.

302 Overall, the method for text summarization that is chosen depends on the particulars of the
303 work at hand, including the volume and complexity of the input text, the length and format of
304 the summary that is sought, and the accessibility of training data. But these recent
305 developments in deep learning models have created new chances for more successful and
306 effective text summarization, with the potential to be useful for a variety of applications in
307 business and academics.

308 There is a large portion of scope for future development and use in the arena of text
309 summarization using deep learning models like TextRank, Seq2Seq, and BART. Multi-
310 document summarization, which involves summarizing several texts on the same subject, is
311 one promising area of future research. For specialised topics, creating domain-specific
312 summarization models is another possibility. Users' feedback on the content and
313 presentation of summaries might be collected through interactive summarizing, resulting in
314 more accurate and individualized summaries. Other significant research areas include the
315 creation of efficient evaluation measures and multilingual summarization algorithms. We can
316 anticipate more complex and successful summarization models that serve a variety of
317 sectors and applications as long as innovation and research in this area continue.

318 **COMPETING INTERESTS**

319
320 Authors have declared that no competing interests exist.

321 **REFERENCES**

- 322
- 323 1. S. Gholamrezazadeh, M. A. Salehi, and B. Gholamzadeh, "A Comprehensive Survey on
324 Text Summarization Systems." 2009 2nd International Conference on Computer Science
325 and its Applications, 2009, doi: 10.1109/csa.2009.5404226.
 - 326 2. J. Goldstein, M. Kantrowitz, V. Mittal, and J. Carbonell, "Summarizing text documents."
327 Proceedings of the 22nd annual international ACM SIGIR conference on Research and
328 development in information retrieval, 1999, doi: 10.1145/312624.312665.
 - 329 3. H. Saggion and G. Lapalme, "Generating Indicative-Informative Summaries with
330 SumUM." Computational Linguistics, vol. 28, no. 4, pp. 497-526, 2002, doi:
331 10.1162/089120102762671963.
 - 332 4. H. P. Luhn, "The Automatic Creation of Literature Abstracts." IBM Journal of Research
333 and Development, vol. 2, no. 2, pp. 159-165, 1958, doi: 10.1147/rd.22.0159.

- 334 5. P.- ying Zhang and C.- he Li, "Automatic text summarization based on sentences
335 clustering and extraction." 2009 2nd IEEE International Conference on Computer
336 Science and Information Technology, 2009, doi: 10.1109/iccsit.2009.5234971.
- 337 6. Barzilay, R., Elhadad, M, "Using Lexical Chains for Text Summarization." In Proc.
338 ACL/EACL'97 Workshop on Intelligent Scalable Text summarization, Madrid,
339 Spain,1997, pp. 10–17.
- 340 7. Y. Ko and J. Seo, "An effective sentence-extraction technique using contextual
341 information and statistical approaches for text summarization." Pattern Recognition
342 Letters, vol. 29, no. 9, pp. 1366-1371, 2008, doi: 10.1016/j.patrec.2008.02.008.
- 343 8. E. Hovy and C.-Y. Lin, "Automated text summarization and the SUMMARIST system."
344 Proceedings of a workshop on held at Baltimore, Maryland October 13-15, 1998 -, 1996,
345 doi: 10.3115/1119089.1119121.
- 346 9. J. Morris and G. Hirst, "The Subjectivity of Lexical Cohesion in Text." The Information
347 Retrieval Series, pp. 41-47, doi: 10.1007/1-4020-4102-0_5.
- 348 10. H. G. Silber and K. F. McCoy, "Efficiently Computed Lexical Chains as an Intermediate
349 Representation for Automatic Text Summarization." Computational Linguistics, vol. 28,
350 no. 4, pp. 487-496, 2002, doi: 10.1162/089120102762671954.
- 351 11. Kedar Bellare, Anish Das Sharma, Atish Das Sharma, Navneet Loiwal and Pushpak
352 Bhattacharyya, "Generic Text Summarization Using Word net. Language Resources
353 Engineering Conference."
- 354 12. J. Leskovec, M. Grobelnik, N. Milic-Frayling, "Extracting Summary Sentences Based on
355 the Document Semantic Graph." Microsoft Research, 2005.
- 356 13. Karel Jezek and Josef Steinberger, "Automatic Text Summarization (The state of the art
357 2007 and new challenges)," Znalosti, pp. 1-12, 2008.
- 358 14. M. Halliday and R. Hasan, "Cohesion in English." 2014, doi: 10.4324/9781315836010.
- 359 15. Ruqaiya Hasan, Coherence and Cohesive Harmony, "In: Flood James (Ed.),
360 Understanding Reading Comprehension: Cognition, Language and the Structure of
361 Prose." Newark, Delaware: International Reading Association, pp. 181-219, 1984.
- 362 16. W. C. Mann and S. A. Thompson, "Relational propositions in discourse." Discourse
363 Processes, vol. 9, no. 1, pp. 57-90, 1986, doi: 10.1080/01638538609544632.
- 364 17. W. C. MANN and S. A. THOMPSON, "Rhetorical Structure Theory: Toward a functional
365 theory of text organization." Text - Interdisciplinary Journal for the Study of Discourse,
366 vol. 8, no. 3, 1988, doi: 10.1515/text.1.1988.8.3.243.
- 367 18. J. Morris and G. Hirst, "The Subjectivity of Lexical Cohesion in Text." The Information
368 Retrieval Series, pp. 41-47, doi: 10.1007/1-4020-4102-0_5.
- 369 19. R. Barzilay, K. R. McKeown, and M. Elhadad, "Information fusion in the context of multi-
370 document summarization." Proceedings of the 37th annual meeting of the Association
371 for Computational Linguistics on Computational Linguistics -, 1999, doi:
372 10.3115/1034678.1034760.
- 373 20. Branimir Boguraev and Christopher Kennedy, "Salience-based Content Characterization
374 of Text Documents," In Proceedings of the ACL'97/EACL'97 Workshop on Intelligent
375 Scalable Text Summarization, 1997.
- 376 21. Li Chengcheng, "Automatic Text Summarization Based On Rhetorical Structure Theory,"
377 International Conference on Computer Application and System Modeling (ICCASM), vol.
378 13, pp. 595-598, October 2010.
- 379 22. Hongyan Jing, "Sentence Reduction for Automatic Text Summarization," In Proceedings
380 of the 6th Applied Natural Language Processing Conference, Seattle, USA, pp. 310-315,
381 2000.

- 382 23. Kevin Knight and Daniel Marcu, "Statistics-Based Summarization Step One: Sentence
383 Compression," In Proceeding of the 17th National Conference of the American
384 Association for Artificial Intelligence, pp. 703-710, 2000.
- 385 24. Hongyan Jing and Kathleen R. McKeown, "Cut and Paste Based Text Summarization,"
386 In Proceedings of the 1st Meeting of the North American Chapter of the Association for
387 Computational Linguistics, Seattle, USA, pp. 178-185, 2000.
- 388 25. K. D. Garg, V. Khullar and A. K. Agarwal, "Unsupervised Machine Learning Approach for
389 Extractive Punjabi Text Summarization," 2021 8th International Conference on Signal
390 Processing and Integrated Networks (SPIN), Noida, India, 2021, pp. 750-754, doi:
391 10.1109/SPIN52536.2021.9566038.
- 392 26. K. Lanyo and A. Wausi, "A Comparative Study of Supervised and Unsupervised
393 Classifiers Utilizing Extractive Text Summarization Techniques to Support Automated
394 Customer Query Question-Answering," 2018 5th International Conference on Soft
395 Computing & Machine Intelligence (ISCMI), Nairobi, Kenya, 2018, pp. 88-92, doi:
396 10.1109/ISCMI.2018.8703237.
- 397 27. P. Sharma and M. Chen, "Text Summarization and Keyword Extraction," 2023 14th IIAI
398 International Congress on Advanced Applied Informatics (IIAI-AAI), Koriyama, Japan,
399 2023, pp. 369-372, doi: 10.1109/IIAI-AAI59060.2023.00078.
- 400 28. S. R. R, M. R. L, D. S, A. D. P K and A. M. S, "Detection and Summarization of Honest
401 Reviews Using Text Mining," 2022 8th International Conference on Smart Structures and
402 Systems (ICSSS), Chennai, India, 2022, pp. 01-05, doi:
403 10.1109/ICSSS54381.2022.9782167.