

Optimal Keyword Selection by Hybrid Optimization with Itemset Mining for Text Summarization in Biomedical Sector

Nishant N. Pachpor, Assistant Professor, IIMS, Chinhewad Pune, India

Salim G. Shaikh, Associate Professor, Department of Computer Engineering, Kalsekar Technical Campus, Panvel India

Ashwini Brahma, Associate Professor, IIMS, Chinhewad Pune, India.

Sachin Misal, Associate Professor, IIMS, Chinhewad Pune, India.

Manuscript Received: Oct 15, 2024; Published: Oct 28, 2024

Abstract

Biomedical literature is growing exponentially, creating challenges for researchers and clinicians to access relevant information efficiently. Automatic biomedical text summarization is a promising solution to this issue, allowing the extraction of essential information while reducing redundancy. Traditional summarization techniques often rely on shallow linguistic features or simple term frequencies, which fail to capture the complex relationships in biomedical texts. This paper explores advanced methods for biomedical text summarization, including graph-based approaches, frequent item set mining, and deep learning models such as BERT. The proposed framework introduces a hybrid optimization technique combining Colliding Bodies Optimization (CBO) and Cuckoo Search Optimization (CSO) for optimal keyword selection, alongside a Recurrent Neural Network (RNN) for sentence categorization. Extensive experimentation using UMLS-based concept extraction, keyword selection, and Apriori-based itemset mining demonstrates that the proposed method significantly outperforms existing models in both the informativeness and coherence of generated summaries. The results reveal that combining deep language models with domain-specific knowledge enhances summarization quality and can be effectively applied to diverse types of biomedical text.

Keywords: Biomedical text summarization, deep learning, hybrid optimization, graph-based summarization, recurrent neural network, Colliding Bodies Optimization (CBO), Cuckoo Search Optimization (CSO).

Introduction

Large volumes of valuable biomedical information are available in the form of scientific articles, medical records, web documents, and clinical reports [13] [14]. However, retrieving, interpreting, and integrating relevant information from the numerous and huge sources of biomedical text are demanding tasks that call for developing automatic tools [9]. Text summarization systems can ease the task of retrieving relevant information by condensing the amount of text that the user has to read [10]. A text summary is a condensed version of the original document that conveys the most relevant information [11]. More specifically, summarization systems can be of high importance in the biomedical domain where clinicians and researchers need to efficiently read and manage a large number of text documents [12]. Automatic text summarization tools can help professionals and users who seek information to effectively and efficiently identify and process useful and relevant information from the vast volume of scientific literature. This also applies to biomedical documents. Automatic biomedical text summarization is an efficient and reliable method aiming at condensing a full-text biomedical paper while preserving its most significant points. Therefore, text summarization plays a pivotal role in alleviating the problem of accessing accurate and up-to-date information relevant to biomedical researchers and physician's needs.

Initial work in text summarization relied on simple term frequency features to identify the most important content of a text document [15] [17]. Since then, many summarization methods have incorporated a wide variety of features and heuristics into the process of content selection. The most widely-used features include the position of sentences, the lengths of sentences, the presence of cue phrases, keywords extracted from the text, the presence of numerical content, the title word, the centroid-based cohesion, the co-occurrence feature and other related features. The majority of text summarization methods do not consider the characteristics of the domain or the type of documents [16]. They mostly work with units extracted directly from the document itself, such as terms, sentences or paragraphs, etc. Then they rely on data mining or information retrieval techniques to analyze effectively this data [19]. However, in the biomedical domain like any other specific domain, these techniques may not seem to be working well because the literature of this domain has its properties and they should be

considered during the summarization process [18] [20]. For this reason, researchers in this domain used domain knowledge resources like ontologies, thesaurus, and taxonomies, etc...to provide meaning to biomedical texts, and then linking information within each text to specifications contained in these resources using a markup language and return concepts that express the semantic meaning of texts.

The summarization problem has been modeled using different linguistic, probabilistic, machine learning, and graph-based approaches [21]. Previous studies indicate that the graph-based approach has great potential to be adopted for summarization of different types of general and domain-specific documents [22]. However, there are still two main challenges that need to be addressed in graph-based summarization. First, when the input text is modeled as a graph, it is crucial to effectively capture the linguistic, semantic, and contextual relationships between the sentences. It is also necessary to have an accurate quantification of the strength of the relations. Secondly, an efficient ranking strategy is required to identify the most important nodes within the graph that correspond to the most related sentences of the document. Deep neural network-based language models [23] can be utilized to address many of the challenges associated with using domain knowledge in context-aware biomedical summarization. A deep language model is pretrained on large corpora of text data and learns how to represent units of text, generally words, in a vector space. These vectorized representations of text can accurately capture a large amount of semantic and syntactic information of words [24] [25]. The pretrained embeddings can be either fine-tuned on a downstream task or used directly as numerical features.

Related Works

In 2020, Moradi et al. [1] have addressed the challenges in the context of biomedical text summarization. The efficacy of a graph-based summarizer was evaluated using different types of context-free and contextualized embeddings. The word representations were produced by pre-training neural language models on large corpora of biomedical texts. The summarizer modelled the input text as a graph in which the strength of relations between sentences was measured using the domain specific vector representations. The usefulness of different graph ranking techniques was also assessed in the sentence selection step of this summarization method. Using the common Recall-Oriented Understudy for Gisting Evaluation (ROUGE) metrics, the performance of this summarizer was evaluated against various comparison methods. The results showed that when the summarizer utilized proper combinations of context-free and contextualized embeddings, along with an effective ranking method, it could outperform the other methods. It was demonstrated that the best settings of this graph-based summarizer could efficiently improve the informative content of summaries and decreased the redundancy.

In 2017, Nasr et al. [2] have proposed a novel graph-based summarization method that took advantage of the domain-specific knowledge and a well-established data mining technique called frequent itemset mining. This summarizer exploited the Unified Medical Language System (UMLS) to construct a concept-based model of the source document and mapping the document to the concepts. Then, it discovered frequent itemsets to take the correlations among multiple concepts into account. The method used these correlations to propose a similarity function based on which a represented graph was constructed. The summarizer then employed a minimum spanning tree based clustering algorithm to discover various subthemes of the document. Eventually, it generated the final summary by selecting the most informative and relative sentences from all subthemes within the text.

In 2019, Rouane et al. [3] have proposed a novel biomedical text summarization system that combined two popular data mining techniques: clustering and frequent itemset mining. Biomedical paper was expressed as a set of biomedical concepts using the UMLS metathesaurus. The K-means algorithm was used to cluster similar sentences. Then, the Apriori algorithm was applied to discover the frequent itemsets among the clustered sentences. Finally, the salient sentences from each cluster were selected to build the summary using the discovered frequent itemsets.

In 2020, Moradi et al. [4] have proposed a novel summarization method that utilized contextualized embeddings generated by the Bidirectional Encoder Representations from Transformers (BERT) model, a deep learning model that recently demonstrated state-of-the-art results in several natural language processing tasks. The different versions of BERT was combined with a clustering method to identify the most relevant and informative sentences of input documents. The summarizer obtained state-of-the-art results and significantly improved the performance of biomedical text summarization in comparison to a set of domain-specific and domain-independent methods. The largest language model not specifically pretrained on biomedical text outperformed other models.

However, among language models of the same size, the one further pretrained on biomedical text obtained best results.

In 2018, Moradi and Ghadiri [5] have described a Bayesian summarization method for biomedical text documents. The Bayesian summarizer initially mapped the input text to the Unified Medical Language System (UMLS) concepts; then it selected the important ones to be used as classification features. Six different feature selection approaches were introduced to identify the most important concepts of the text and selected the most informative contents according to the distribution of these concepts. With the use of an appropriate feature selection approach, the Bayesian summarizer could improve the performance of biomedical summarization. Using the Recall-Oriented Understudy for Gisting Evaluation (ROUGE) toolkit, extensive evaluations were performed on a corpus of scientific papers in the biomedical domain. The results showed that when the Bayesian summarizer utilized the feature selection methods that do not use the raw frequency, it could outperform the biomedical summarizers that rely on the frequency of concepts, domain-independent and baseline methods.

In 2018, Milad Moradi [6] has proposed a novel summarization method named Clustering and Itemset mining based Biomedical Summarizer (CIBS). The summarizer extracted biomedical concepts from the input documents and employed an itemset mining algorithm to discover main topics. Then, it applied a clustering algorithm to put the sentences into clusters such that those in the same cluster shared similar topics. Selecting sentences from all the clusters, the summarizer could produce a summary that covered a wide range of topics of the input text. Using the Recall-Oriented Understudy for Gisting Evaluation (ROUGE) toolkit, the performance of the CIBS method was evaluated against four summarizers including a state-of-the-art method. The results showed that the CIBS method could improve the performance of single- and multi-document biomedical text summarization. It was shown that the topic-based sentence clustering approach could be effectively used to increase the informative content of summaries, as well as to decrease the redundant information.

In 2017, Moradi and Ghadiri [7] have proposed a summarization method that combined itemset mining and domain knowledge to construct a concept-based model and to extract the main subtopics from an input document. Our summarizer quantified the informativeness of each sentence using the support values of itemsets appearing in the sentence. To address the concept-level analysis of text, this method initially mapped the original document to biomedical concepts using the Unified Medical Language System (UMLS). Then, it discovered the essential subtopics of the text using a data mining technique, namely itemset mining, and constructed the summarization model. The employed itemset mining algorithm extracted a set of frequent itemsets containing correlated and recurrent concepts of the input document. The summarizer selected the most related and informative sentences and generated the final summary.

In 2020, Du et al. [8] have proposed a novel model called BioBERTSum to better capture token-level and sentence-level contextual representation, which used a domain-aware bidirectional language model pretrained on large-scale biomedical corpora as encoder, and further fine-tuned the language model for extractive text summarization task on single biomedical document. Especially, a sentence position embedding mechanism was adopted, which enabled the model to learn the position information of sentences and achieved the structural feature of document. To the best of the knowledge, this was the first work to use the pre-trained language model and fine-tuning strategy for extractive summarization task in the biomedical domain. Experiments on PubMed dataset showed that this proposed model outperformed the recent SOTA (state-of-the-art) model by ROUGE-1/2/L.

Problem Definition

Biomedical text summarization can be used to extract the key information from the vast biomedical documents yet it cannot make the usage of domain-aware external knowledge nor learn document-level and sentence-level features. These issues require to be handled in the future. Some of the major features and challenges are listed in Table 1. Graph ranking method [1] minimizes the redundancy and enhances the informative content of summaries. But, the usage of new deep language models like Multi-Task Deep Neural Network (MT-DNN) is not investigated and it cannot be applied to various summarization tasks like query-oriented or multi-document summarization. Frequent itemset mining [2] considers several correlations between multiple concepts and return more semantics and it also handles the intrinsic ambiguities related to the biomedical literature in a very effective manner. Still, various knowledge sources are not combined to enhance the summarizer's performance and the advantages of removal of the stop words, stemming, and typical pre processing steps are not considered. Clustering

and frequent itemset mining [3] exceeds the baselines and various summarizers and also improves the quality of the produced summaries. Yet, the duplicate information is not reduced and it does not combine the techniques of word embedding in the text representation. Bidirectional Encoder Representations from Transformers (BERT) [4] produces a novel generation of context-aware biomedical text summarizers and permits the summarizer to quantify the relatedness of sentences on the basis of various dimensions of the vector space. But, the standard corpus of documents is lacked along with their summaries and the usage of contextualized language models is not examined for various types of biomedical text summarization. Bayesian summarization method [5] chooses the most informative contents on the basis of the distribution and the significant concepts of the text are also recognized. Still, it does not consider various discriminative classifiers and it also does not focus on query-focused and multi-document summarization. Clustering and Itemset mining based Biomedical Summarizer [6] enhances the information coverage and the amount of redundant information present in the summary is also reduced. Yet, it does not recognize the necessary concepts inside a group of documents using novel measures and methods and a trade-off is not established among the redundancy and the information coverage. Itemset mining and domain knowledge [7] quantifies the informativeness of every sentence by the summarizer with the help of the support values and an accurate concept-oriented model is produced for the summarization. But, it does not extend the GraphSum to deal with concepts and it also does not extend the itemset-based summarization modeling for the query-focused biomedical text summarization. BioBERTSum [8] uses the fine-tuning strategy and pre-trained language model and it also achieves the structural feature of the document. Still, it does not monitor the generation process by more expert knowledge and the abstractive and extractive methods are not integrated. These challenges are laid as a basic motivation to introduce a novel biomedical text summarization method.

Table 1: Features and challenges of state-of-the-art biomedical text summarization methods

Author [citation]	Methodology	Features	Challenges
Moradi <i>et al.</i> [1]	Graph ranking method	<ul style="list-style-type: none"> The informative content of summaries is enhanced. It reduces the redundancy. 	<ul style="list-style-type: none"> It cannot be applied to various summarization tasks like query-oriented or multi-document summarization. It does not examine the usage of new deep language models like Multi-Task Deep Neural Network (MT-DNN)
Nasr <i>et al.</i> [2]	Frequent itemset mining	<ul style="list-style-type: none"> The intrinsic ambiguities related to the biomedical literature are handled in an effective manner. It considers various correlations between multiple concepts and return more semantics. 	<ul style="list-style-type: none"> It does not consider the advantages of removal of the stop words, stemming, and typical pre processing steps. It does not combine various knowledge sources to enhance the summarizer's performance.
Rouane <i>et al.</i> [3]	Clustering and frequent itemset mining	<ul style="list-style-type: none"> The quality of the produced summaries is improved. It exceeds the baselines and various summarizers. 	<ul style="list-style-type: none"> The techniques of word embedding are not combined in the text representation. It does not minimize the duplicate information.
Moradi <i>et al.</i> [4]	Bidirectional Encoder Representations from Transformers (BERT)	<ul style="list-style-type: none"> It permits the summarizer to quantify the relatedness of sentences on the basis of various dimensions of the vector space. It produces a novel generation of context-aware biomedical text summarizers. 	<ul style="list-style-type: none"> It does not examine the usage of contextualized language models for various types of biomedical text summarization. It lacks the standard corpus of documents along with their summaries.



Moradi and Ghadiri [5]	Bayesian summarization method	<ul style="list-style-type: none"> It recognizes the significant concepts of the text. On the basis of the distribution, the most informative contents are chosen. 	<ul style="list-style-type: none"> It does not concentrate on query-focused and multi-document summarization. Various discriminative classifiers are not considered.
Milad Moradi [6]	Clustering and Itemset mining based Biomedical Summarizer	<ul style="list-style-type: none"> It minimized the amount of redundant information present in the summary. It enhances the information coverage. 	<ul style="list-style-type: none"> It does not establish a trade-off among the redundancy and the information coverage. The necessary concepts inside a group of documents are not recognized using novel measures and methods.
Moradi and Ghadiri [7]	Itemset mining and domain knowledge	<ul style="list-style-type: none"> It produces an accurate concept-oriented model for the summarization. With the help of the support values, the informativeness of every sentence is quantified by the summarizer. 	<ul style="list-style-type: none"> The itemset-based summarization modelling was not extended for the query-focused biomedical text summarization. GraphSum is not extended to deal with concepts.
Du <i>et al.</i> [8]	BioBERTSum	<ul style="list-style-type: none"> The structural feature of the document is achieved. It uses the fine-tuning strategy and pre-trained language model. 	<ul style="list-style-type: none"> It does not integrate the abstractive and extractive methods. The generation process is not monitored by more expert knowledge.

Research Objectives

This research works covers the following objectives.

- To perform a detailed review on text summarization in biomedical and other relevant applications to tackle the challenging part of this research.
- To extract the relevant keywords from the text using the new hybrid optimization concept.
- To accomplish the sentence categorization whether it could be suitable for summarized content using an improved deep learning.
- To develop the new variant or hybrid optimization algorithm for implementing the efficient text summarization model.
- To confirm the efficiency of the proposed model over the conventional models by analysing the relevant performance measures.

Research Methodology and Proposed Model

Automatic text summarization offers an efficient solution to access the ever-growing amounts of both scientific and clinical literature in the biomedical domain by summarizing the source documents while maintaining their most informative contents. Some of biomedical text summarization systems put the basis of their sentence selection approach on the frequency of concepts extracted from the input text. However, it seems that exploring other measures rather than the raw frequency for identifying valuable contents within an input document, or considering correlations existing between concepts, may be more useful for this type of summarization. This proposal tactics to develop the text summarization model in biomedical sector. The different phases of the proposed text summarization will be (a) pre-processing, (b) mapping text to biomedical concept, (c) keywords extraction, (d) optimal keyword selection, (e) sentence or itemset categorization, (f) itemset mining, and (g) summary construction. Initially, the pre-processing of the text will be done by tokenization and stopword removal. Once the pre-processing is finished, mapping text to biomedical concept will be done UMLS Metathesaurus, which is a large, multi-lingual, and multi-purpose lexicon containing millions of biomedical and health related concepts, their relationships and their synonymous names. Further, the keywords will be extracted using the Bag

of n-grams, is a natural extension of bag of words that simply generates the sequence of n tokens or words. Based on the relevance of keywords, the optimal keyword selection will be performed based on the hybridization of two optimization algorithms like Colliding Bodies Optimization (CBO) [26] and Cuckoo Search Optimization (CSO) [27]. A deep learning model called Recurrent Neural Network (RNN) will be adopted for the sentence or itemset categorization. Once the itemsets are categorized, frequent itemset mining will be accomplished by the Apriori Algorithm. Further, the corresponding sentences will be evaluated and the summary will be constructed. The results demonstrate that this combination can successfully enhance the summarization performances, and the proposed system outperforms other tested summarizers. The proposed model is shown in Fig. 1.

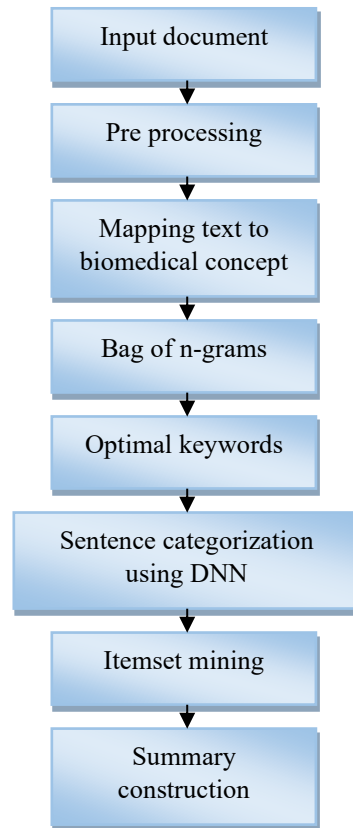


Figure 1: Proposed Text Summarization Model

Conclusion

The rapid expansion of biomedical literature necessitates efficient methods for summarizing vast amounts of data. The research presented in this paper introduces a novel biomedical text summarization framework that combines multiple state-of-the-art techniques, such as deep learning, graph-based ranking, and frequent itemset mining. The hybrid optimization approach, involving CBO and CSO algorithms, significantly improves the selection of relevant keywords, while the use of RNN for sentence categorization ensures the capture of meaningful text structures. Through rigorous evaluation using the UMLS concept extraction and Apriori itemset mining, the proposed method demonstrates superior performance over conventional approaches. Future work will focus on extending the system to handle multi-document summarization and explore the integration of abstractive methods with extractive techniques to further refine the summarization process. These advancements represent a substantial step forward in addressing the critical challenge of navigating large-scale biomedical documents.

References

- [1] Milad Moradi, Maedeh Dashti, and Matthias Samwald, "Summarization of biomedical articles using domain-specific word embeddings and graph ranking", *Journal of Biomedical Informatics*, vol. 107, 2020.
- [2] Mozghan Nasr Azadani, Nasser Ghadiri, Ensieh Davoodijam, "Graph-based biomedical text summarization: An itemset mining and sentence clustering approach", *Journal of Biomedical Informatics*, October 2017.
- [3] Oussama Rouane, Hacene Belhadef, and Mustapha Bouakkaz, "Combine clustering and frequent itemsets mining to enhance biomedical text summarization", *Expert Systems with Applications*, vol. 135, pp. 362-373, November 2019.
- [4] Milad Moradi, Georg Dorffner, and Matthias Samwald, "Deep contextualized embeddings for quantifying the informative content in biomedical text summarization", *Computer Methods and Programs in Biomedicine*, vol. 184, February 2020.
- [5] Milad Moradi, and Nasser Ghadiri, "Different approaches for identifying important concepts in probabilistic biomedical text summarization", *Artificial Intelligence in Medicine*, vol. 84, pp. 101-116, January 2018.
- [6] Milad Moradi, "CIBS: A biomedical text summarizer using topic-based sentence clustering", *Journal of Biomedical Informatics*, vol. 88, pp. 53-61, December 2018.
- [7] Milad Moradi, and Nasser Ghadiri, "Quantifying the informativeness for biomedical literature summarization: An itemset mining method", *Computer Methods and Programs in Biomedicine*, vol. 146, pp. 77-89, July 2017.
- [8] Yongping Du, Qingxiao Li, Lulin Wang, and Yanqing He, "Biomedical-domain pre-trained language model for extractive summarization", *Knowledge-Based Systems*, vol. 199, July 2020.
- [9] R. Mishra, J. Bian, M. Fiszman, C.R. Weir, S. Jonnalagadda, J. Mostafa, "Text summarization in the biomedical domain: a systematic review of recent research", *J. Biomed. Inform.*, vol. 52, pp. 457-467, 2014.
- [10] M. Gambhir, V. Gupta, "Recent automatic text summarization techniques: a survey", *Artif. Intell. Rev.*, vol. 47, pp. 1-66, 2016.
- [11] S. Afantenos, V. Karkaletsis, and P. Stamatopoulos, "Summarization from medical documents: a survey", *Artificial Intelligence in Medicine*, vol. 33, no. 2, pp. 157-177, 2005.
- [12] J.-g. Yao, X. Wan, and J. Xiao, "Recent advances in document summarization", *Knowledge and Information Systems*, vol. 53, no. 2, pp. 297-336, 2017.
- [13] Baralis, E., Cagliero, L., Mahoto, N., & Fiori, A. "GraphSum: Discovering correlations among multiple terms for graph-based summarization", *Information Sciences*, vol. 249, pp. 96-109, 2013.
- [14] Erkan, G., & Radev, D. R. "Lexrank: Graph-based lexical centrality as salience in text summarization", *Journal of Artificial Intelligence Research*, vol. 22, pp. 457-479, 2004.
- [15] A. Mahajani, V. Pandya, I. Maria, and D. Sharma, "A Comprehensive Survey on Extractive and Abstractive Techniques for Text Summarization," Singapore, pp. 339-351, 2019.
- [16] L. H. Reeve, H. Han, and A. D. Brooks, "The use of domain-specific concepts in biomedical text summarization," *Information Processing & Management*, vol. 43, pp. 1765-1776, 2007.
- [17] Reeve LH, Han H, Brooks AD. "The use of domain-specific concepts in biomedical text summarization", *Inf Process Management*, vol. 43, pp. 1765-1776, 2007.
- [18] Plaza L, Carrillo-de-Albornoz J. "Evaluating the use of different positional strategies for sentence selection in biomedical literature summarization", *BMC Bioinf*, vol. 14, no. 1, 2013.
- [19] R. Ferreira, L. de Souza Cabral, F. Freitas, R. D. Lins, G. de França Silva, S. J. Simske, "A multi-document summarization system based on statistics and linguistic treatment," *Expert Systems with Applications*, vol. 41, pp. 5780-5787, 2014.
- [20] Y. Sankarasubramaniam, K. Ramanathan, and S. Ghosh, "Text summarization using Wikipedia," *Information Processing & Management*, vol. 50, pp. 443-461, 2014.
- [21] H. Saggion, "SUMMA: A robust and adaptable summarization tool," *Traitement Automatique des Langues*, vol. 49, 2008.
- [22] M. A. Fattah and F. Ren, "GA, MR, FFNN, PNN and GMM based models for automatic text summarization," *Computer Speech & Language*, vol. 23, pp. 126-144, 2009.
- [23] K. Sarkar, M. Nasipuri, S. Ghose, "Using machine learning for medical document summarization", *Int. J. Database Theory Applications*, vol. 4, no. 1, pp. 31-48, 2011.



- [24] P. Gayathri, N. Jaisankar, "Towards an efficient approach for automatic medical document summarization", *Cybern. Inf. Technol.*, vol. 15, no. 4, pp. 78-91, 2015.
- [25] Y. Shang, Y. Li, H. Lin, Z. Yang, "Enhancing biomedical text summarization using semantic relation extraction", *PLoS One*, vol. 6, no. 8, 2011.
- [26] A. Kaveh, and V.R. Mahdavi, "Colliding bodies optimization: A novel meta-heuristic method", *Computers & Structures*, vol. 139, pp. 18-27, July 2014.
- [27] A.S. Joshi, Omkar Kulkarni, G.M. Kakandikar, and V.M. Nandedkar, "Cuckoo Search Optimization- A Review", *Materials Today: Proceedings*, vol. 4, no. 8, pp. 7262-7269, 2017.