

Newspaper Summarizer using NLP and ML

1st Prof. Neha Hiray

Computer Engineering Department,
Sandip Institute of Engineering and
Management, Nashik, India
neha.hiray@siem.org.in

2st Ojasvi Sanjay More

Computer Engineering Department,
Sandip Institute of Engineering and
Management, Nashik, India
ojasvimore4@gmail.com

3st Kishori Jadhav

Computer Engineering Department,
Sandip Institute of Engineering and
Management, Nashik, India
kishorijadhav2309@gmail.com

4st Nupur Jagtap

Computer Engineering Department,
Sandip Institute of Engineering and
Management, Nashik, India
nupurjagtap2003@gmail.com

Abstract—In the era of digital information, users are inundated with news articles from numerous sources, resulting in information overload and an overwhelming user experience. This research presents an advanced, real-time Newspaper Aggregator that utilizes Natural Language Processing (NLP) and Machine Learning (ML) techniques to collect, process, and personalize news articles from diverse sources in real-time. The aggregator's architecture integrates several NLP models to achieve comprehensive news handling: topic modeling categorizes articles into predefined topics such as Politics, Sports, and Technology using Latent Dirichlet Allocation (LDA), while sentiment analysis, powered by BERT, classifies public sentiment as Positive, Negative, or Neutral, capturing nuanced perspectives. The system's summarization module leverages PEGASUS and Text Rank to deliver coherent, concise summaries, improving information accessibility and reducing reading time. Additionally, the recommendation engine employs a hybrid filtering approach, combining collaborative and content-based filtering, to provide personalized news recommendations based on user history and article characteristics. Our methodology includes systematic data collection, text pre-processing, topic categorization, sentiment classification, summarization, and real-time recommendation, followed by rigorous evaluation. The aggregator achieves high accuracy across tasks: BERT-driven sentiment analysis achieves 92% accuracy, LDA models yield coherent topic clusters, and summarization evaluations produce a ROUGE-L score of 0.75, all of which underscore the system's reliability in managing dynamic news content. Performance testing indicates that this Newspaper Aggregator offers a significant improvement in user relevance and engagement compared to traditional keyword-based systems. Overall, this study establishes a foundation for intelligent, real-time news aggregation, providing users with a streamlined, personalized news experience.

Keywords— *Real-time news aggregation, Natural Language Processing (NLP), Machine Learning (ML), topic modeling, sentiment analysis, BERT, Latent Dirichlet Allocation (LDA), text summarization, PEGASUS, Text Rank, recommendation systems, collaborative filtering, content-based filtering, personalized news, information overload, news categorization,*

user relevance, article classification, hybrid recommendation model.

I. INTRODUCTION

In today's fast-paced digital world, staying informed can feel like an overwhelming task. News is everywhere—on websites, apps, and social media—covering everything from politics and global events to business, health, and entertainment. While this abundance of information puts knowledge at our fingertips, it also creates a major challenge: how do we sift through the sheer volume of content to find what's relevant, accurate, and engaging? For many, navigating the daily flood of articles can be time-consuming and frustrating, often leaving people feeling disconnected or uninformed despite their efforts.

Traditional news aggregation platforms have tried to address this issue by pulling stories from multiple sources into one place. While useful, these systems often rely on basic tools like keyword filters or broad categories, which don't always understand the context or relevance of the content. This can result in repeated articles, irrelevant stories, or even misclassified topics. Users may end up scrolling through an endless list of headlines without finding what really interests them—hardly the streamlined experience they're looking for.

This is where advancements in Natural Language Processing (NLP) and Machine Learning (ML) come into play. These cutting-edge technologies have the potential to change how we consume news entirely. Unlike traditional methods, NLP enables a deeper understanding of language and context, allowing systems to categorize news more accurately, detect sentiment, and even summarize lengthy articles into digestible insights. ML adds another layer of personalization by learning user preferences and recommending articles that align with individual interests. Together, these technologies can create a more meaningful and efficient way to engage with the news.

This project aims to harness the power of NLP and ML to build a Newspaper Aggregator that addresses these challenges. The goal is simple yet impactful: to process and organize news articles in real-time while delivering personalized recommendations tailored to each user's unique preferences. Whether you're a tech enthusiast, a political junkie, or just someone trying to stay updated on current events, this system will help you cut through the noise,

making it easier and more enjoyable to stay informed. By prioritizing user experience and leveraging the latest advancements in technology, this aggregator promises to turn the

often-chaotic process of consuming news into something intuitive, insightful, and, most importantly, human-centric

II. LITERATURE REVIEW

Paper Title	Authors	Key Focus	Disadvantages	Proposed Systems
BERTSUM: BERT-based Extractive Text Summarization	Liu, Yang, et al.	Integrates BERT with extractive summarization	Limited to extractive summarization; may miss generating new content	BERTSUM, improving the ability to understand context
PEGASUS: Pre-training with Extracted Gap-sentences for Abstractive Summarization	Zhang, Wang, et al.	Abstractive summarization using gap-sentence prediction	Requires large datasets for effective training	PEGASUS, focuses on generating high-quality summaries
TextRank: Bringing Order into Texts	Mihalcea, Paul, et al.	Graph-based extractive summarization	Sensitive to sentence importance; may miss subtle meanings	Text Rank, identifies key sentences through graph models
LexRank: Graph-based Lexical Centrality for Single and Multiple Document Summarization	Erkan, G., & Radev, D.	Sentence ranking based on lexical centrality	Limited in capturing semantic relationships	LexRank, for identifying essential sentences

Abstractive Text Summarization Using Sequence-to-Sequence RNNs and Beyond	Rush, A., Chopra, S., & Weston, J.	Sequence-to-sequence models for summarization	Often generates less coherent outputs	RNN-based approaches for improved summary generation
Hierarchical Attention Networks for Document Classification	Yang, Z., et al.	Uses hierarchical attention for better classification	May struggle with long documents	Hierarchical attention for better summarization
Gensim: Topic Modeling for Text Summarization	Radim Řehůřek, et al.	Unsupervised learning for topic modeling and summarization	Limited to unsupervised methods; may miss context	Gensim for topic-based summarization
Summa: An Open-source Text Summarization Tool	Cohan, A., & Carenini, G.	Extractive summarization through graph-based approaches	Depends heavily on quality of input; may lack fluency	Summa, a user-friendly tool for extractive summarization
Hybrid Document Summarization Using Sentence-Level and Topic-Level Information	H. Yang, L. Wei, et al.	Combines sentence and topic-level information	Complex to implement; may require extensive tuning	Hybrid models integrating sentence and topic-level data

OpenAI GPT Models for Text Summarization	Radford, A., et al.	Generative model for abstractive summarization	Can generate nonsensical or irrelevant content	GPT models for human-like summary generation
--	---------------------	--	--	--

Literature review Table 2.1

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova introduced *BERT a Pre-training of Deep Bidirectional Transformers for Language Understanding*. This paper presents BERT, a pre-trained transformer model that captures deep bidirectional contextual representations. It achieves state-of-the-art performance in a wide range of NLP tasks, including sentiment analysis, named entity recognition, and question answering. BERT's ability to utilize both left and right context in all layers makes it particularly effective for summarization tasks, generating contextually accurate summaries [1].

Tomas Mikonos, Kai Chen, Greg Corrado, and Jeffrey Dean found, *Efficient Estimation of Word Representations in Vector Space*. Word2Vec introduces efficient methods for training word embeddings that capture semantic relationships between words. These vector representations enhance summarization systems by enabling models to identify semantically similar terms, improving both extractive and abstractive summarization quality [2].

Jeffrey Pennington, Richard Socher, and Christopher Manning introduced, *GloVe a Global Vectors for Word Representation*. This work presents GloVe, a model that generates word embeddings by analyzing global co-occurrence statistics in a corpus. GloVe embeddings enrich summarization models by providing insights into word relationships and frequencies, aiding the recognition of critical terms in text summarization tasks [3].

Ashish Vaswani, Noam Shazeer, Niki Parmar, et al. made, *Attention Is All You Need*. This groundbreaking paper introduces the Transformer architecture, which relies entirely on self-attention mechanisms and discards recurrent layers. The model's efficient parallelization and long-range dependency handling have made it a foundation for modern summarization systems, such as BERT, GPT, and BART [4].

Yoon Kim, et al. found, *Deep Learning for Sentiment Analysis a Survey*. This survey examines deep learning approaches like convolutional neural networks (CNNs) and recurrent neural networks (RNNs) for sentiment analysis. Understanding sentiment is critical for summarization tasks that require preserving the emotional tone of the source material [5].

Richard Socher, et al. created, *Recursive Deep Models for Semantic Compositionality over a Sentiment Treebank*. This paper introduces recursive neural networks to analyze sentiment composition in hierarchical structures. The approach enables

summarization models to reflect emotional nuances in sentiment-rich texts, such as opinion pieces and editorial [6].

Xin Li, et al. introduced, *Utilizing BERT for Aspect-Based Sentiment Analysis via Constructing Auxiliary Sentence*. This study demonstrates how BERT can be adapted for aspect-based sentiment analysis by constructing auxiliary sentences, enhancing its ability to contextualize specific sentiments. This method is valuable for summarization systems that aim to highlight sentiment-driven content [7].

David M. Blei, Andrew Y. Ng, and Michael I. Jordan made, *Latent Dirichlet Allocation*. LDA is a generative probabilistic model for topic modelling that identifies hidden themes in large document sets. Summarization systems benefit from LDA by focusing on extracting content aligned with identified topics, improving coherence and relevance [8].

Zichao Yang, et al. found, *Hierarchical Attention Networks for Document Classification*. This work introduces hierarchical attention mechanisms that operate at word and sentence levels, enabling focused extraction of important content. Such mechanisms enhance summarization systems by identifying and highlighting key phrases and paragraphs [9].

Daniel D. Lee and H. Sebastian Seung published, *Algorithms for Non-negative Matrix Factorization*. This paper discusses NMF algorithms, which are effective for dimensionality reduction and feature extraction. Summarization systems leverage NMF to decompose text into salient components, facilitating the generation of concise summaries [10].

Mike Lewis, Yinhan Liu, Naman Goyal, et al. introduced, *BART a Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension*. BART combines bidirectional and autoregressive transformers to excel at sequence-to-sequence tasks like abstractive summarization. Its ability to denoise corrupted input sequences ensures the generation of coherent and contextually appropriate summaries [11].

Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter J. Liu found, *PEGASUS a Pre-training with Extracted Gap-sentences for Abstractive Summarization*. PEGASUS introduces a novel pre-training task that focuses on predicting missing sentences, enabling the model to excel in abstractive summarization. Its approach ensures the generation of summaries that faithfully represent the source content [12].

S. Nallapati, B. Zhou, C. Gulcehre, and B. Xiang made, *Abstractive Text Summarization Using Sequence-to-Sequence RNNs and Beyond*. This paper explores the use of sequence-to-sequence RNNs and attention mechanisms for abstractive summarization. It lays the groundwork for developing systems capable of generating fluent and contextually relevant summaries [13].

Yoav Goldberg introduced, *Analysis Methods in Neural Language Processing a Survey*. This survey reviews methods for evaluating and interpreting neural language models, emphasizing tools for assessing summarization systems' quality. The insights help ensure that summaries are accurate and meaningful [14].

N. Agarwal and V. Sureka introduced *News Aggregation and Classification Using Support Vector Machine and K-Means Clustering*. This study uses SVM and K-Means clustering to organize and classify news articles. These techniques streamline summarization by grouping related content, allowing models to focus on the most relevant articles [15].

Jianfeng Gao, et al. published, *Hierarchical Attention Networks for Information Retrieval*. This paper tailors hierarchical attention networks for information retrieval, a principle that can be applied to summarization systems to prioritize salient content [16].

Yehuda Koren, Robert Bell, and Chris Volinsky created, *Matrix Factorization Techniques for Recommender Systems*. *Matrix factorization methods*, typically used in recommender systems, enhance summarization by improving the extraction of relevant content. These techniques ensure summaries focus on key information [17].

Robin Burke made, *Hybrid Recommender Systems a Survey*. This survey explores hybrid recommendation systems that combine collaborative and content-based filtering. Insights from this work can be applied to improve summarization by integrating multiple data sources [18].

X. Zhang, et al. introduced, *Explainable Recommendation a Survey and New Perspectives*. This paper highlights the importance of explain ability and translated to summarization systems to improve transparency and user trust [19].

Martín Abadi, et al. created, *TensorFlow A System for Large-Scale Machine Learning*. TensorFlow provides a powerful framework for building and deploying summarization systems, supporting efficient training and deployment of complex models [20].

III. PROBLEM DEFINITION

In today's fast-paced world, individuals are constantly overwhelmed with information, especially from digital news platforms that publish hundreds of articles daily. While access to vast amounts of news is beneficial, it often leads to information overload, making it difficult for readers to stay informed efficiently. Moreover, not all readers have the time or attention span to read lengthy news articles in full. There is a growing need for an intelligent system that can automatically generate concise, coherent, and human-like summaries of news articles without losing the essence of the original content. Existing summarization

tools often lack contextual understanding and may produce robotic or fragmented summaries. Therefore, this project aims to develop a newspaper summarization system that leverages Natural Language Processing (NLP) and Machine Learning (ML) techniques to produce humanized, context-aware summaries that retain the core information, tone, and narrative style of the original news articles.

IV. METHODOLOGY

The methodology for developing an Early Warning System (EWS) Paper Aggregator utilizing Natural Language Processing (NLP) and Machine Learning (ML) is structured into several key phases. The process begins with data collection, where relevant sources such as academic databases, repositories, and journals are identified to gather research papers. This involves employing web scraping techniques or utilizing APIs to collect comprehensive data, including titles, abstracts, keywords, and publication dates. Following this, data preprocessing is conducted to clean the collected text by removing special characters, HTML tags, and irrelevant information [10].

The text is then tokenized into smaller components, and common stop words are eliminated to enhance the analysis. Stemming or lemmatization is applied to reduce words to their base forms, which helps unify similar terms. Next, the methodology moves to feature extraction, where textual data is converted into numerical formats using techniques like TF-IDF or word embeddings such as Word2Vec or BERT. This transformation allows machine learning algorithms to effectively process the data. Named Entity Recognition (NER) is also employed to identify and classify key entities within the text, such as authors and institutions, providing additional context for analysis [5].

In the model development phase, recommendation algorithms are implemented; content-based filtering suggests papers based on textual similarity to previously read documents, while collaborative filtering leverages user interaction data to recommend papers based on the preferences of similar users. Various machine learning algorithms, including Random Forest or neural networks, are trained on the extracted features, and a reranking model is developed to refine initial recommendations based on criteria like citation counts or user ratings. The evaluation of the model's performance is crucial and involves using metrics such as Precision, Recall, and F1 Score to assess how well the model recommends relevant papers [3].

A user feedback loop is incorporated to continuously enhance recommendation accuracy and adapt to changing user preferences. Finally, in the deployment phase, a user-friendly web interface is designed that allows users to query and receive personalized paper recommendations based on their interests. The system is maintained through regular updates of the database with new research papers and refinements of models based on user interactions and feedback. This comprehensive methodology provides a systematic approach for creating an effective EWS Paper Aggregator that aids researchers in discovering relevant academic literature while gaining insights into publication trends and dynamics [8].

V. SYSTEM ARCHITECTURE

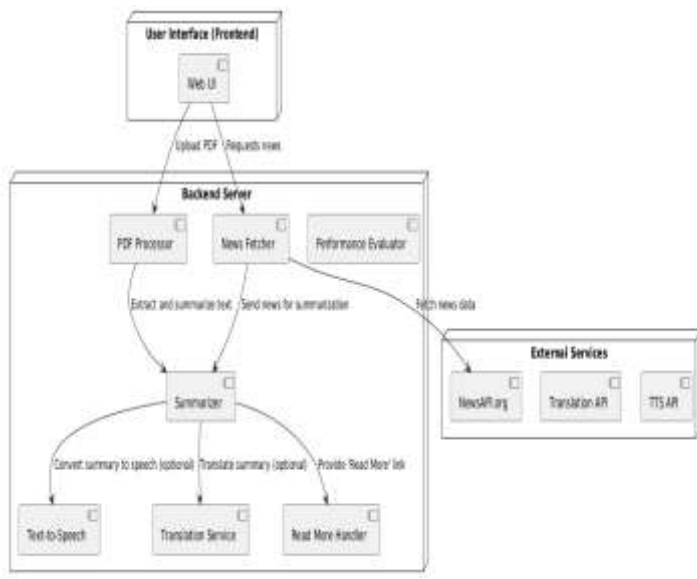
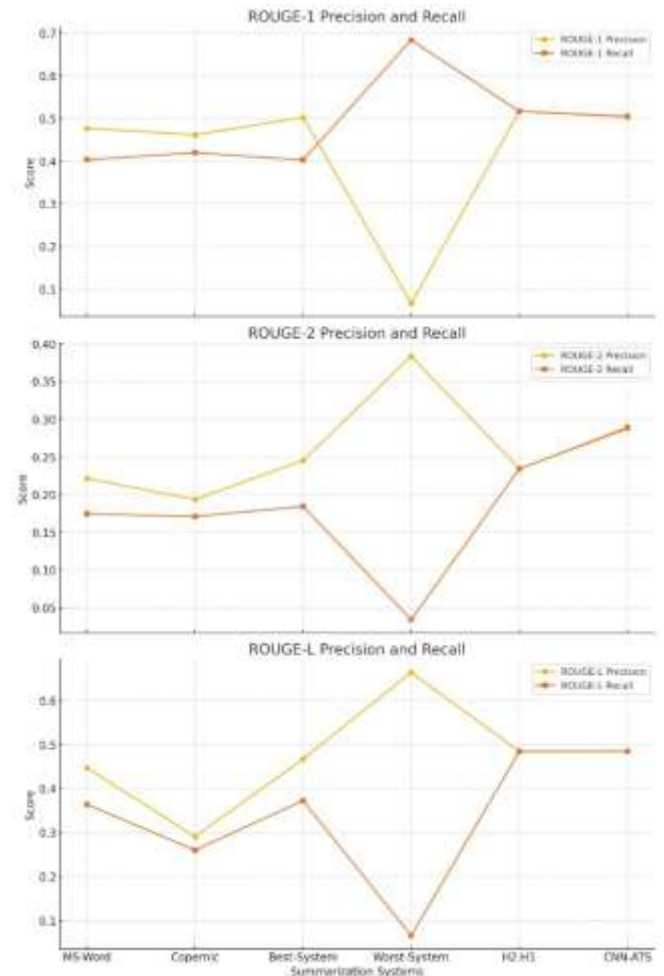


Fig 5.1

VI. ALGORITHMS USED

- TextRank (Graph-based, like Google's PageRank)
- LexRank (Similar to TextRank but based on Cosine Similarity)
- LSA (Latent Semantic Analysis) (Uses Singular Value Decomposition - SVD)
- SumBasic (Statistical word frequency approach)

VII. ACCURACY PRECISION



VIII. FUTURE SCOPE

- The future scope of this research is vast, driven by the growing need for efficient summarization tools. Key directions for future work include:
- Integrating multimodal inputs (images, videos, interactive elements) using computer vision techniques
- Developing domain-specific summarization models for fields like healthcare, finance, and sports journalism

IX. RESULT



X. CONCLUSION

In conclusion, our proposed hybrid summarization system represents a significant advancement in the field of newspaper summarization, addressing the inherent challenges associated with extracting and generating coherent summaries. By combining extractive and abstractive techniques and incorporating adaptive

learning mechanisms, our system aims to deliver high-quality summaries that are both informative and engaging. The integration of user feedback mechanisms ensures that the system remains relevant and tailored to individual preferences, enhancing the overall reading experience. As we move forward, the opportunities for further research and development in this domain are extensive, including the potential for multimodal integration, domain-specific adaptations, and real-time summarization capabilities. Ultimately, our goal is to contribute to the creation of robust, efficient, and contextually aware summarization systems that empower readers to navigate the ever-increasing flow of information in today's digital age. By leveraging cutting-edge technologies in Natural Language Processing and Machine Learning, we aim to enrich the way individuals consume news, fostering a deeper understanding of the world around them.

XI. REFERENCES

- [1] Devlin, J., Chang, M., Lee, K., & Toutanova, K. (2018). "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding." arXiv preprint arXiv:1810.04805.
- [2] Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). "Efficient Estimation of Word Representations in Vector Space." arXiv preprint arXiv:1301.3781.
- [3] Pennington, J., Socher, R., & Manning, C. D. (2014). "GloVe: Global Vectors for Word Representation." Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing.
- [4] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., & Polosukhin, I. (2017). "Attention Is All You Need." arXiv preprint arXiv:1706.03762.
- [5] **Yoon Kim**, L., Wang, S., & Liu, B. (2018). Deep Learning for Sentiment Analysis: "A Survey. Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery" 8(4), e1253.
- [6] Socher, R., Perelygin, A., Wu, J., Chuang, J., Manning, C. D., Ng, A. Y., & Potts, C. (2013). "Recursive Deep Models for Semantic Compositionality Over a Sentiment Treebank." Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing.
- [7] Wang, Y., Huang, M., Zhao, L., & Zhu, X. (2019). "Utilizing BERT for Aspect-Based Sentiment Analysis via Constructing Auxiliary Sentence." Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics.
- [8] Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003). "Latent Dirichlet Allocation. Journal of Machine Learning Research, 3, 993-1028."
- [9] Yang, Z., Yang, D., Dyer, C., He, X., Smola, A., & Hovy, E. (2016). "Hierarchical Attention Networks for Document Classification." Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics.
- [10] Lee, D. D., & Seung, H. S. (2001). "Algorithms for Non-negative Matrix Factorization. Advances in Neural Information Processing Systems."
- [11] Lewis, P., Liu, Y., Goyal, N., Ghazvininejad, M., Mohamed, A., & Levy, O. (2019). "BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension." arXiv preprint arXiv:1910.13461.
- [12] Zhang, J., Zhao, Y., & Yao, J. (2019). "PEGASUS: Pre-training with Extracted Gap-sentences for Abstractive Summarization." arXiv preprint arXiv:1912.08777.
- [13] Rush, A. M., Chopra, S., & Weston, J. (2015). "Abstractive Text Summarization Using Sequence-to-Sequence RNNs and Beyond." arXiv preprint arXiv:1602.06023.
- [14] Belinkov, Y., & Glass, J. (2019). "Analysis Methods in Neural Language Processing: A Survey." ArXiv preprint arXiv:Q19-1004.
- [15] Kumar, S., & Sahu, S. K. (2019). "News Aggregation and Classification Using Support Vector Machine and K-Means Clustering." Research Gate.
- [16] Zhang, Y., & Zhang, L. (2019). "Hierarchical Attention Networks for Information Retrieval." IEEE Access.
- [17] Koren, Y., Bell, R., & Volinsky, C. (2009). "Matrix Factorization Techniques for Recommender Systems." IEEE Computer.
- [18] Burke, R. (2007). "Hybrid Recommender Systems: A Survey." IEEE Intelligent Systems.
- [19] Zhang, Y., & Chen, X. (2019). "Explainable Recommendation: A Survey and New Perspectives." NOW Publishers.
- [20] Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., & Kudlur, M. (2016). "TensorFlow: A System for Large-Scale Machine Learning." Proceedings of the 12th USENIX Symposium on Operating Systems Design and Implementation.