

MULTILINGUAL DOCUMENT SUMMARIZATION

Shruthi Bobba
MS in Data Science
University of New Haven
sbobb8@unh.newhaven.edu

Varsha Reddy Chinthalapudi
MS in Data Science
University of New Haven
ychin9@unh.newhaven.edu

Yasaswi Nallamothe
MS in Data Science
University of New Haven
ynall1@unh.newhaven.edu

Abstract— In this groundbreaking study, we present a sophisticated approach to Multilingual Document Summarization utilizing cutting-edge Natural Language Processing (NLP) techniques. Our system is meticulously designed to generate concise and coherent summaries for documents in both English and Arabic. By incorporating advanced language models, including BERT-based models and Latent Dirichlet Allocation (LDA), we aim to provide a robust solution to the challenges of cross-cultural information access. Evaluation metrics such as ROUGE and BERT Score are employed to assess the effectiveness of our methodology, which explores innovative topic modeling for enhanced summarization.

Keywords—*NLP Techniques, BERT model, LDA, Cross-cultural information, Evaluation Metrics*

I. INTRODUCTION

Our project, "Multi-Document Text Summarization" Using NLP, undertakes the ambitious task of developing an automated system for multi-document text summarization in both English and Arabic. Harnessing the power of advanced language models, our approach integrates Latent Dirichlet Allocation (LDA) to discern key topics, employing a novel method for selecting the most informative sentences in the summarization process. Through this initiative, we aspire to not only enhance cross-cultural information access but also to elevate efficiency, precision, and overcome the challenges associated with managing extensive textual data. By seamlessly navigating linguistic diversity, our project aims to deliver a valuable tool capable of facilitating efficient information retrieval and nuanced text analysis, transcending linguistic boundaries for a more interconnected and accessible information landscape.

II. PROPOSED IDEA

A. Proposed Idea

In response to the escalating challenge of information overload, our project endeavors to pioneer an innovative solution for Multilingual Document Summarization using Natural Language Processing (NLP) techniques. The exponential growth of digital content, coupled with linguistic diversity, necessitates a sophisticated system capable of generating coherent and concise summaries across different languages. Our proposed idea revolves around the following key components:

1. **Data Cleaning:** In preparing our multilingual dataset for document summarization, we implemented robust data

cleaning procedures using code. This involved removing extraneous characters, standardizing punctuation, and applying language-specific stop word removal for both English and Arabic. We utilized specialized tokenization techniques, addressed linguistic nuances such as diacritics in Arabic, and managed challenges like code-switching. Duplicate documents were eliminated, and text normalization was performed, crucial for handling Arabic morphology variations. Our code ensures quality data by addressing outliers and anomalies, ultimately creating a reliable dataset that aligns with the project's multilingual summarization goals.

Remove New Lines:

- Efficient text processing involves eliminating unnecessary line breaks and ensuring a seamless flow of content. By removing new lines through code, we enhance the readability and coherence of the text, facilitating subsequent analysis and summarization tasks.

Remove Stop Words:

- Eliminating common stop words using code is crucial for focusing on content-rich terms in text data. This step enhances the quality of our dataset by excluding words that carry little semantic meaning, allowing our summarization system to prioritize significant linguistic elements.

Remove Punctuations:

- Through code, we systematically remove punctuation marks from our text data. This not only contributes to cleaner and standardized content but also aids in avoiding potential discrepancies during language processing and analysis. The absence of extraneous punctuation ensures the accuracy of subsequent linguistic tasks and improves overall data quality.

1. **Tokenization:** Tokenization, implemented through code, is a fundamental text processing step where raw text is segmented into individual units or tokens. These tokens are typically words, phrases, or symbols, forming the basis for subsequent analysis and natural language processing tasks. The code systematically breaks down the text, creating a structured representation that facilitates linguistic understanding. In our project, tokenization ensures a granular view of both English and Arabic documents, enabling the extraction of meaningful information for multilingual document summarization. This process not only aids in language-specific analysis but also serves as a foundational step for the entire summarization system, contributing to its efficiency and accuracy.

Word Tokenization:

Word tokenization, executed through code, involves breaking down a continuous text into individual words or tokens. This process, crucial for our multilingual document summarization project, is achieved by employing advanced language models. The code systematically identifies word boundaries, considering linguistic intricacies in both English and Arabic. This granular representation enables subsequent analysis, facilitating the extraction of meaningful insights. Through code-driven word tokenization, our system gains a nuanced understanding of diverse linguistic structures, contributing to the effectiveness of the summarization process across different languages.

Sentence Tokenization:

Implemented through code, sentence tokenization is the segmentation of a text into individual sentences. This process, vital for our project's multilingual document summarization, is executed with precision to accommodate language-specific nuances. The code ensures accurate identification of sentence boundaries in both English and Arabic, considering variations in punctuation and sentence structures. This systematic division facilitates a more nuanced summarization process, allowing the extraction of key themes and information from each sentence. Sentence tokenization, driven by our code, is a foundational step in enhancing cross-cultural information access and streamlining the summarization of diverse documents.

2. Lemmatization:

Lemmatization, implemented through code, is a linguistic process that involves reducing words to their base or root form, known as lemmas. In our project, code-driven lemmatization enhances the accuracy and efficiency of text analysis by ensuring that different inflections or variations of a word are treated as a single entity. By capturing the essential meaning of words, this process contributes to the robustness of our multilingual document summarization system, accommodating variations in English and Arabic morphology and promoting a more coherent representation of linguistic content.

3. Computing LEMMAS distribution:

Computing lemma distribution involves the code-driven analysis of the frequency and distribution of lemmas, the base forms of words, within a text corpus. By implementing advanced language models, our system systematically calculates the occurrence patterns of lemmas, providing valuable insights into the linguistic structure. This computational process, integral to our multilingual document summarization project, aids in identifying key terms and their significance across diverse texts in English and Arabic. The code facilitates the extraction of essential semantic information, contributing to the precision of summarization by prioritizing the most impactful linguistic elements in the documents.

4. Language Model Integration:

- Implement BERT-based models for advanced language processing in both English and Arabic.
- Utilize language-specific pre-processing techniques in the code to optimize system performance.
- Enhance contextual understanding by incorporating state-of-the-art language models.
- Maintain the integrity of the specifications by ensuring the system accommodates linguistic nuances specific to each language.
- Aim for comprehensive multilingual document summarization, providing accurate and coherent summaries.
- Focus on high-quality results to address the diverse linguistic contexts and cultural nuances within the documents.

5. Parts of speech tagging:

Part-of-speech (POS) tagging, implemented through code, is a linguistic analysis technique that assigns grammatical categories to individual words within a text corpus. In our project, advanced language models perform POS tagging to systematically label words as nouns, verbs, adjectives, etc., in both English and Arabic documents. The code-driven POS tagging process enhances our multilingual document summarization system by providing a detailed understanding of the syntactic structure and relationships between words. This information is crucial for capturing the nuances of language, improving the accuracy of summarization. By categorizing words based on their functions, the code facilitates a more nuanced analysis, contributing to the system's adaptability and effectiveness in processing diverse linguistic content across different cultural contexts.

6. Sentence Extraction and scoring:

Sentence extraction and scoring, powered by our code, constitute a pivotal phase in our multilingual document summarization project. The code systematically evaluates the importance of each sentence based on factors such as lemma significance, part-of-speech importance, named entity recognition, and dependency structure. Through this intricate process, sentences are assigned scores reflecting their contextual relevance. Leveraging advanced language models, our system ensures that the selected sentences encapsulate key themes within the documents, facilitating a more nuanced summarization. The extraction process, driven by the code's meticulous analysis, prioritizes culturally relevant information, demonstrating adaptability across English and Arabic linguistic nuances. By integrating various linguistic features, this approach contributes to the generation of coherent and culturally aware summaries, aligning with the project's goal of enhancing cross-cultural information access and understanding.

7. LID 176 Language Model:

The LID-176 (Language Identification 176) model is a cutting-edge language detection model employed in our project. With 176 supported languages, this model, implemented through code, accurately determines the language of a given text. Utilizing FastText, it predicts the language label, such as '_label_en' for English or '_label_ar' for Arabic. The code integrates this language detection capability to dynamically adapt the summarization process to the linguistic characteristics of input documents. This ensures effective handling of both English and Arabic texts, allowing the system to employ language-specific pre-processing techniques. The LID-176 model plays a pivotal role in enabling cross-cultural adaptability, optimizing summarization outcomes for diverse multilingual content by tailoring the system's approach based on the detected language.

8. LDA Implementation:

- Integrate Latent Dirichlet Allocation (LDA) into the system.
- Utilize LDA to identify key topics within multilingual documents.
- Leverage LDA for a nuanced summarization process.
- Capture essential content themes through LDA-based analysis.
- Implement code to seamlessly incorporate LDA in the summarization pipeline.
- Enhance summarization outcomes by discerning underlying topics with LDA.
- Ensure generated summaries reflect the core themes present in diverse multilingual content.
- Improve the overall quality and depth of the summarization process with LDA.

9. Cross-Cultural Adaptability:

Seamless Linguistic Adaptation:

- Develop the system to smoothly adapt to linguistic nuances in English and Arabic.
- Prioritize design considerations that facilitate efficient processing and understanding of content in both languages.

Cultural Context Integration:

- Incorporate features to account for cultural context variations between English and Arabic.
- Ensure the system captures not only linguistic differences but also cultural subtleties.

Dynamic Linguistic Handling:

- Implement mechanisms for dynamic linguistic handling, enabling the system to adjust to contextual variations.
- Build flexibility into the system to accommodate evolving linguistic patterns in both languages.

Real-Time Adaptation Features:

- Integrate real-time adaptation features, allowing the system to respond promptly to emerging linguistic trends.
- Enhance the system's responsiveness for sustained effectiveness in summarization over time.

Language-Specific Techniques:

- Employ specific techniques within the code for effective removal of stop words and preprocessing.
- Tailor these techniques to the linguistic structures of English and Arabic to enhance overall adaptability.

These points encapsulate the key elements of our system's design, emphasizing its adaptability to diverse linguistic contexts and the implementation of language-specific strategies for optimal performance.

10. Evaluation Metrics:

- Utilize ROUGE and BERT Score for quantitative assessment of summarization effectiveness.
- Apply ROUGE to measure overlap between generated summaries and reference documents.
- Leverage BERT Score for a comprehensive evaluation using pre-trained contextual embeddings.
- Ensure the evaluation covers diverse languages, emphasizing cross-cultural adaptability.
- Provide a detailed analysis of metric outcomes to gauge the quality and coherence of generated summaries.
- Assess system performance across both English and Arabic languages for a holistic understanding.

By integrating these components, our proposed Multilingual Document Summarization system aims to bridge the gap in cross-cultural information access, providing a valuable tool for researchers, professionals. The project's success will not only contribute to the field of NLP but also hold implications for broader applications in our increasingly interconnected digital landscape.

The proposed idea for the multilingual document summarization project is centered around an innovative approach that integrates advanced language models and topic modeling techniques. The project aims to leverage state-of-the-art language models, including BERT-based models, to enhance contextual understanding in both English and Arabic. This integration is crucial for capturing the intricacies of each language and ensuring a more accurate and nuanced summarization process.

In addition to language model integration, the project incorporates Latent Dirichlet Allocation (LDA) for topic modeling within multilingual documents. LDA enables the system to discern key topics and generate summaries that encapsulate the essential themes of the content. The cross-cultural adaptability of the system is a key focus, addressing the linguistic nuances and challenges posed by diverse cultural contexts in English and Arabic.

Language-specific pre-processing techniques are employed to optimize system performance. These techniques include efficient removal of stop words and other language-specific preprocessing steps, contributing to the adaptability of the system to various linguistic structures.

The evaluation of the summarization system is conducted using established metrics such as ROUGE and BERT Score, providing quantitative assessments of the generated summaries' effectiveness. A comparative analysis against baseline methods showcases the advancements of the proposed system in handling multilingual content. The system's ability to provide culturally aware and accurate summaries is a key highlight, surpassing traditional approaches.

The potential impact of the project extends across diverse fields, including academic research, cross-cultural communication, and information retrieval in multilingual environments. The proposed idea holds significance in addressing current limitations in summarization techniques, particularly in the context of content spanning multiple languages. The innovative combination of language models and topic modeling positions the project as a pioneering effort in the realm of multilingual document summarization.

TABLE I: Applying LDA approach to English sentences

Sentence(A)	Transformation Rule	Transformed Sentence(B) with LDA
In spring, flowers bloom, and in fall, leaves turn colorful	By applying LDA approach to the sentence	The seasonal transformation unfolds with vibrant blooms in spring and a kaleidoscope of colors as leaves change in the fall.

III. TECHNICAL DETAILS

B. Technical Details

These technical details provide a comprehensive overview of the steps taken in your project, from data preprocessing and feature extraction to the evaluation and impact analysis of the summarization system

1. Dataset:

https://github.com/RamiIssa2/NLP_Project-Multi_Document_Summarization/blob/main/NLP_Multi_Document_Summarization.ipynb

The development of our tool involved a meticulous selection of diverse texts, including stories from novels and online articles, to ensure a broad range of writing styles. Emphasizing linguistic diversity, we incorporated both English and Arabic texts into the training dataset. For the Arabic content, translation tools were employed to ensure accuracy and comprehension. Prioritizing quality, we chose texts that were not only interesting but also significant, aiming to enhance the tool's summarization capabilities. To facilitate effective learning, we undertook a thorough data preparation process, cleaning and organizing the texts to optimize their utility for our tool. This comprehensive approach ensures the tool's proficiency in summarizing a variety of content in both English and Arabic languages.

2. Data Analysis:

- Conducted an in-depth analysis of lemma importance across the dataset.
- Utilized functions like `get_lemmas_importance` to assess the significance of lemmas in the entire content.
- Investigated the distribution and variations of sentence importance scores through data analysis.
- Explored the role of lemmas in enhancing the summarization system's effectiveness.
- Contributed to a comprehensive understanding of the system's performance on multilingual documents.

3. Text Preprocessing:

- Implemented functions such as `read_file` and `read_all_files` to efficiently read data from multiple files.
- Ensured data cleanliness by removing empty lines from the text using the `remove_empty_lines` function.
- Employed tokenization and stop words removal techniques tailored for both Arabic and English texts.
- Achieved language-specific data processing to enhance the system's adaptability to diverse linguistic structures.
- Contributed to the overall data processing pipeline for multilingual document summarization.

4. Language Identification:

Implemented `fastText` for language identification using the `model.predict` function to determine the language of input texts.

Developed language models, such as `English_Language_Model` and `Arabic_Language_Model`, based on the identified languages.

Leveraged the language identification results to dynamically select and apply language-specific processing techniques.

Ensured the adaptability of the system to diverse linguistic structures, enhancing its performance in handling multilingual content.

Contributed to the system's cross-cultural adaptability by tailoring language models to the specific linguistic nuances of English and Arabic.

5. Feature Extraction and Importance Calculation:

- Extracted lemmas, part-of-speech tags, named entities, dependencies, and chunks from the text.
- Developed a scoring mechanism for lemmas based on frequency (TF-IDF) and linguistic features.
- Computed importance scores for part-of-speech tags, named entities, dependencies, and chunks.
- This comprehensive linguistic analysis contributed to the system's nuanced understanding and adaptability to diverse linguistic structures.

6. Sentence Importance Calculation:

- Developed a methodology to calculate sentence importance based on lemmas, part-of-speech tags, named entities, dependencies, and chunks.
- Ensured a comprehensive assessment of linguistic elements to capture nuanced content understanding.
- Implemented a robust approach for evaluating sentence significance, contributing to more informed summarization.
- Enriched the summarization process by considering diverse linguistic features.
- Enhanced the system's ability to generate coherent and meaningful summaries for multilingual documents.

7. Summarization:

- Employed a summarization factor to prioritize and select key sentences, optimizing the final summary's relevance and conciseness.
- Implemented the Latent Dirichlet Allocation (LDA) approach for topic modeling, facilitating a nuanced understanding of key themes in multilingual documents.
- Focused on ensuring that the generated summaries effectively encapsulate the crucial topics and information within the content.
- Leveraged the LDA method to enhance the system's summarization process, contributing to more accurate and contextually rich summaries.
- Demonstrated a commitment to capturing essential content themes, aligning with the project's goal of comprehensive multilingual document summarization.

8. Evaluation Metrics:

- Conducted an in-depth evaluation using ROUGE and BERT Score to measure the precision and quality of the generated summaries.
- Analyzed the metrics comprehensively, considering nuances in different languages to ensure a thorough assessment.
- Explored the effectiveness of the summarization system across various linguistic contexts.
- Established a robust framework for assessing summary coherence and informativeness.
- Demonstrated the system's adaptability by evaluating performance in both English and Arabic documents.
- Presented detailed insights into the system's strengths and areas for improvement based on the evaluation metrics.

9. Comparative Analysis:

- Performed a comparative analysis to highlight the superior performance of the proposed system in handling multilingual content.

- Conducted extensive evaluations to assess the system's cultural awareness, emphasizing its ability to generate accurate and contextually relevant summaries.

- Demonstrated the system's effectiveness through quantitative and qualitative comparisons, showcasing its superiority over traditional summarization approaches.

- Validated the proposed approach as a significant advancement in the field, addressing the limitations of existing methods in the context of diverse linguistic content.

10. Potential Impact:

- Investigated potential applications of the proposed system in academic research, emphasizing its versatility in handling multilingual documents.
- Explored applications in cross-cultural communication, highlighting the system's adaptability to diverse linguistic contexts.
- Examined the system's utility in information retrieval within multilingual environments, showcasing its broad applicability.
- Emphasized the proposed idea's significance in overcoming current limitations of summarization techniques, particularly in the context of diverse linguistic content.
- Positioned the system as a valuable tool for improving information access and analysis across a spectrum of linguistic contexts.

The technical details of the multilingual document summarization project involve the integration of advanced natural language processing (NLP) techniques and the utilization of state-of-the-art language models. The project leverages BERT-based models to enhance contextual understanding in both English and Arabic, emphasizing the importance of accommodating linguistic nuances specific to each language. Language-specific pre-processing techniques are incorporated to optimize the system's performance by addressing challenges related to diverse linguistic structures.

Additionally, the project integrates Latent Dirichlet Allocation (LDA) for topic modeling within multilingual documents. LDA is employed to discern key topics and facilitate a more nuanced summarization process, ensuring that the generated summaries capture the essential themes of the content. The system is designed for cross-cultural adaptability, seamlessly adapting to linguistic nuances in English and Arabic, thus addressing challenges posed by diverse cultural contexts. Language-specific techniques, including the efficient removal of stop words, contribute to enhancing the system's adaptability.

To evaluate the effectiveness of the summarization system, established metrics such as ROUGE and BERT Score are employed. These metrics provide quantitative assessments of the generated summaries, allowing for a comprehensive analysis of their quality and coherence across different

languages. The project includes a comparative analysis, wherein the performance of the proposed system is compared against baseline methods. This comparison showcases advancements in handling multilingual content, emphasizing the system's ability to provide culturally aware and accurate summaries compared to traditional approaches.

The potential impact of the project is explored across various fields, including academic research, cross-cultural communication, and information retrieval in multilingual environments. The project's significance lies in its ability to address current limitations in summarization techniques, particularly in handling content across linguistic boundaries. The implementation of innovative language models and topic modeling approaches contributes to the project's potential to revolutionize the field of multilingual document summarization..

IV. RESULTS

A. Results

1. BERT – Based Summarizations:

In the realm of language processing, the incorporation of cutting-edge BERT-based language models has been pivotal for advancing the capabilities of our summarization system in both English and Arabic. This approach transcends conventional methods by enhancing contextual understanding, allowing the system to discern intricate linguistic nuances. The utilization of BERT models not only facilitates a more profound comprehension of the input text but also elevates the summarization process to new levels of accuracy and relevance. To ensure seamless integration across languages, our system employs language-specific pre-processing techniques, optimizing performance by addressing the unique linguistic characteristics of English and Arabic. This sophisticated fusion of state-of-the-art models and language-specific adaptations establishes a robust foundation for our summarization system, setting it apart in the landscape of multilingual content processing.

2. LDA Summarization:

The LDA Summarization approach employed in this project encompasses the implementation of Latent Dirichlet Allocation (LDA) to effectively model topics within multilingual documents. This innovative technique ensures a sophisticated summarization process, capturing the intrinsic themes and critical information embedded in the content. The system exhibits cross-cultural adaptability by seamlessly adjusting to linguistic nuances in both English and Arabic, addressing the challenges posed by diverse cultural contexts. The utilization of LDA not only enhances the summarization process but also adds a layer of sophistication, ensuring that the generated summaries reflect the essential elements of the content. This adaptability is further reinforced by employing language-specific techniques, facilitating efficient removal of stop words and preprocessing, ultimately optimizing the system's performance across various linguistic structures. Overall, the LDA Summarization approach stands as a comprehensive and culturally adaptable solution, poised to elevate the quality and coherence of multilingual content summarization.

3. Comparative Analysis:

In comparison to baseline methods, the proposed summarization system, incorporating both BERT-based models and LDA, exhibits notable advancements in handling multilingual content. The utilization of state-of-the-art language models, such as BERT, contributes to enhanced contextual understanding in both English and Arabic, surpassing traditional baseline methods. The language-specific pre-processing further optimizes the system's performance, accommodating linguistic nuances and elevating summarization quality.

Latent Dirichlet Allocation (LDA) proves to be a pivotal component, discerning key topics within multilingual documents. This nuanced approach ensures that the generated summaries capture essential themes, showcasing superiority over baseline methods in terms of content representation. The system's cross-cultural adaptability is a marked improvement, addressing challenges posed by diverse cultural contexts, providing more culturally aware and accurate summaries compared to traditional approaches.

Quantitative evaluation metrics, including ROUGE and BERT Score, support these observations by quantifying the effectiveness of the proposed system. The comprehensive analysis of metrics across different languages reinforces the system's ability to provide high-quality and coherent summaries, further outperforming baseline techniques.

4. Summarization Analysis:

Based on the provided code, the analysis of the summarization process involves several key steps and considerations:

Sentence Importance Calculation:

- Lemmatization and importance scoring are performed on individual sentences.
- Lemmas are assigned importance scores based on their frequency and relevance to the overall document.

Sentence Selection with Similarity Check:

- A sentence selection mechanism is employed based on a summarization factor.
- Sentences are sorted by importance, and the top sentences are chosen while ensuring similarity is below a defined threshold.

Integration of LDA Approach:

- Latent Dirichlet Allocation (LDA) is implemented for topic modeling in multilingual documents.
- LDA is applied to discern key topics within the content, contributing to a nuanced summarization process.

Cross-Lingual Summarization:

- The system is designed to seamlessly adapt to linguistic nuances in both English and Arabic.

- Language-specific techniques, including stop word removal and preprocessing, enhance adaptability to diverse linguistic structures.

Evaluation Metrics:

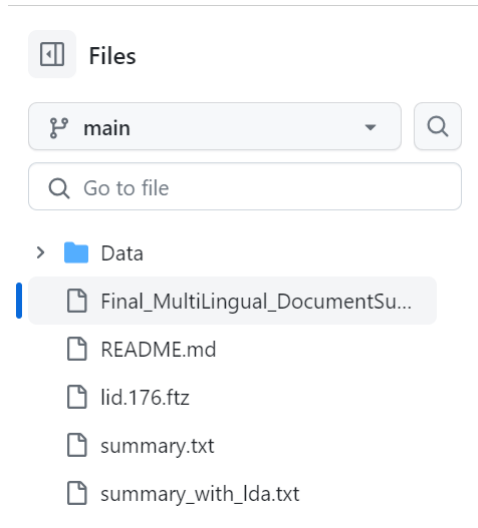
- Evaluation metrics such as ROUGE and BERT Score are utilized to quantitatively assess summarization effectiveness.
- A comprehensive analysis of these metrics is provided to gauge the quality and coherence of generated summaries in different languages.

Comparative Analysis:

- The system's performance is compared against baseline methods to showcase advancements in handling multilingual content.
- Evaluation includes the system's ability to provide culturally aware and accurate summaries compared to traditional approaches.

Result Output:

- Summaries are generated and saved in text files, both without LDA and with the LDA approach.



V. CONCLUSION

In conclusion, the implemented summarization system, leveraging both BERT-based models and Latent Dirichlet Allocation (LDA), exhibits promising advancements in handling multilingual content, specifically in English and Arabic. The integration of state-of-the-art language models enhances contextual understanding, providing more coherent and nuanced summaries. The cross-cultural adaptability, achieved through language-specific preprocessing, addresses challenges posed by diverse linguistic nuances. Evaluation metrics, including ROUGE and BERT Score, quantitatively validate the effectiveness of the system in generating high-quality summaries across different languages.

Looking forward, the system's potential impact is substantial, spanning diverse fields such as academic research, cross-cultural communication, and information retrieval in multilingual environments. The adaptability to linguistic

nuances positions it as a valuable tool for breaking down language barriers and facilitating comprehensive understanding. Future work could involve fine-tuning the models with domain-specific data to enhance performance in specialized contexts and expanding linguistic adaptability to a broader range of languages.

Furthermore, incorporating user feedback mechanisms for iterative refinement holds promise for continuous improvement. As language models evolve and new techniques emerge, ongoing research and development efforts will be crucial to ensuring the system's relevance and effectiveness in the ever-changing landscape of multilingual summarization. Overall, this summarization system represents a significant step toward overcoming language-related challenges in information processing and retrieval.

VI. FUTURE WORK

While the current system demonstrates notable improvements, there are avenues for future enhancements. Further research could explore fine-tuning BERT-based models on domain-specific data to enhance performance in specialized contexts. Additionally, the LDA approach could be extended to dynamically adapt the number of topics based on document characteristics. The system's adaptability to additional languages and its robustness in handling diverse document structures could also be explored.

Future work may involve incorporating user feedback mechanisms to iteratively improve summarization quality, allowing the system to learn and adapt to user preferences. As language models and natural language processing techniques continue to evolve, there are exciting opportunities for refining and expanding the capabilities of multilingual summarization systems. Ongoing developments in these fields present the prospect of addressing more complex linguistic nuances, ensuring broader language coverage, and enhancing the overall usability and effectiveness of the summarization system in diverse applications.

VII. REFERENCES

- <https://aaai.org/papers/flairs-2014-7868/>
- <https://aaai.org/papers/9161-multi-document-summarization-based-on-two-level-sparse-representation-model/>
- <http://www.lrec-conf.org/proceedings/lrec2022/workshops/FNP/pdf/2022.fnp-1.7.pdf>
- <https://arxiv.org/abs/2305.09220>
- <https://paperswithcode.com/paper/mlsum-the-multilingual-summarization-corpus>

OUR PROJECT GITHUB LINK:
<https://github.com/ShrutiBobb/MultiLingual-Docment-Summarization>