

Text Summarization in the Shona Language using Natural Language Processing

Anita Sithabisiwe Manokore
Harare Institute of Technology
HIT
Harare, Zimbabwe

Monica Gondo
Harare Institute of Technology
HIT
Harare, Zimbabwe

Abstract:- The rise of digital information in many languages, including Shona, highlights the significance of developing effective text summarizing techniques to promote information accessibility and usability. This work fills a major gap in the natural language processing (NLP) for the Shona language, which is widely spoken in Zimbabwe and its surrounding areas but has received little attention. The lack of pre-trained language models specifically designed for Shona, the intricacy of Shona's morphology, and the scarcity of annotated datasets provide the main obstacles to Shona text summarization.[1] The goal of this research is to create and modify contemporary machine learning methods for efficient Shona text summarizing in order to address these issues. By gathering and analyzing texts from a variety of sources, such as news stories, scholarly papers, and social media, we produced large annotated corpora. These datasets were utilized to fine-tune existing NLP models, such as Transformer-based architectures, ensuring they account for Shona's specific language traits. To address the morphological and syntactic complexities of Shona, our solution combines statistical and rule-based techniques. When compared to baseline methods, the results show a significant improvement in the relevancy and accuracy of Shona text summaries. In addition to serving as a starting point for further NLP research in underrepresented languages, the generated models help Shona-speaking people in the areas of business, education, and media. By encouraging inclusivity and linguistic variety, showcasing the possibility for cross-lingual breakthroughs, and emphasizing the ethical development of technology, this research adds to the larger area of NLP.

Keywords:- Shona Text Summerization, Extractive Models, Abstractive Models, Hybrid Model, Natural Langauage Proc Essing.

I. INTRODUCTION

The swift progress of digital technology has significantly altered the dissemination and consumption of information, highlighting the necessity for efficient tools to handle the massive volumes of data generated on a daily basis. By reducing lengthy texts to succinct summaries, text summarization, an important aspect of natural language processing (NLP), plays an important role in improving information accessibility. While widely spoken languages

like Shona have seen tremendous advancements in text summarization approaches, low-resource languages like Shona have received less attention, leaving a large gap in NLP capabilities for these languages.

Spoken by over 10 million people, mainly in Zimbabwe, Shona is a Bantu language that offers distinct opportunities and obstacles for text summary. Its complex morphological structure, which is typified by a wide range of grammatical rules and a heavy reliance on prefixes and suffixes, renders traditional summary techniques less useful. Moreover, the creation of reliable NLP applications is made more difficult by Shona's lack of annotated datasets and pre-established models. By increasing the linguistic diversity and inclusion of computational tools, addressing these issues would not only help Shona speakers access information more easily but will also advance the area of natural language processing (NLP).

In order to close the gap in Shona text summarization, this study develops and modifies state-of-the-art machine learning approaches specifically suited to Shona's linguistic peculiarities. Our methodology entails building extensive annotated corpora from a range of text sources, such as scholarly journals, news stories, and social media posts. We make use of sophisticated NLP models, like Transformer-based architectures, which have been adjusted to take into account the unique morphological and syntactic characteristics of Shona. Through the combination of statistical and rule-based techniques, we improve the models' capacity to produce pertinent and precise summaries.

The first large-scale annotated corpora for Shona text summarization, the adaptation and optimization of cutting-edge NLP models for Shona, and the demonstration of enhanced summarization performance over baseline techniques are the three main contributions of this research. This work has practical implications for media, education, and business within Shona-speaking communities, in addition to serving as a foundation for future NLP research in underrepresented languages. In the end, our study shows how cross-lingual innovations and morally sound NLP technologies can be developed to represent the wide range of languages spoken around the world.

II. LITERATURE REVIEW

In recent years, there has been a noticeable increase in interest in text summarization using Natural Language Processing (NLP) techniques. In order to extract relevant information from patent filings, [2] used machine learning, deep learning, and artificial intelligence algorithms. This shows how useful natural language processing (NLP) is when assessing technological advantages. [3] Created a problem-based learning course that uses big data text summary to teach natural language processing (NLP). The course demonstrates how NLP can be used practically to summarize different kinds of collections. In recent years, there has been a noticeable increase in interest in text summarization using Natural Language Processing (NLP) techniques. In order to extract relevant information from patent filings, [2] used machine learning, deep learning, and artificial intelligence algorithms. This shows how useful natural language processing (NLP) is when assessing technological advantages. [3] Created a problem-based learning course that uses big data text summary to teach natural language processing (NLP). The course demonstrates how NLP can be used practically to summarize different kinds of collections. [5] Also emphasized the significance of language models (LMs) in tasks like summarizing and answering questions, highlighting the necessity for cognitive agents to successfully extract knowledge from LMs. Moreover, Wazery et al.'s study from 2022, which used deep learning models to focus on abstractive Arabic text summarization, demonstrated the sequence-to-sequence models' capacity for excellent performance. [7] Investigated how NLMs may be used to create phishing emails, highlighting the flexibility of NLP approaches in many contexts. In the context of business process management, [8] examined the efficacy of language models in natural language processing activities, such as text summarization. Furthermore, [9] evaluation of ChatGPT in automatic code summarization emphasizes how crucial it is to gauge how well language models function in particular contexts. The importance of Recursive Neural Networks (RNNs) in maximizing NLP performance was highlighted by [10], especially in tasks like sentiment analysis and language modeling. Last but not least, [11] presented Ascle, a Python NLP toolbox for creating medical text, demonstrating the adaptability of language models in a range of generating tasks, such as text summarization. Overall, the evaluation of the literature shows how widely NLP approaches are used in a variety of fields, including education, business process management, medical text production, and patent analysis. Text summarizing is one area in which these techniques are particularly useful. The findings highlight how crucial it is to use cutting-edge NLP models and algorithms to glean insightful information and improve performance across a range of jobs.

III. METHODOLOGY

By utilizing cutting-edge natural language processing (NLP) techniques, this work seeks to create an efficient text summary system specifically designed for the Shona language. The construction of the model, evaluation, and data gathering and preparation are the main components of the methodology. Finally, complete content and organizational editing before formatting. Please take note of the following items when proofreading spelling and grammar:

A. Corpus Compilation

In response to the dearth of annotated datasets specifically for Shona, we assembled a broad corpus of Shona writings from multiple sources, such as:

- News Articles: Collected from Shona news websites online.
- Academic Publications: Selected from conference papers and publications published in Shona.
- Social media: postings with substantial Shona content were taken from Facebook and Twitter.

B. Annotation

Summaries were appended to the texts that were gathered. This comprised:

- The process of manually creating reference summaries for a portion of the corpus involves native Shona speakers.
- Automated Techniques: Creating summaries by initial extraction techniques, which were subsequently checked and adjusted by specialists.

C. Text Processing

Among the preprocessing actions were:

- Tokenization: Using a Shona tokenizer, divide text into tokens (words and punctuation).
- Normalization is the process of addressing case variations, eliminating punctuation, and normalizing white spaces in order to transform text into a consistent format.
- Eliminating typical Shona stopwords that don't considerably add to the text's meaning is known as stopword removal.
- Reducing words to their base or root forms in order to manage morphological differences is known as stemming or lemmatization.

D. Model Development

➤ Baseline Models

To create a performance standard, we put baseline extractive summarization models into practice:

- Sentences with the highest term frequency across the document are chosen for frequency-based summarization.
- Position-Based Summarization: This method extracts sentences in the text according to how relevant their positions are.

➤ *Transformer-Based Models*

Transformer-based architectures customized for Shona were used to accomplish advanced summarization:

- **Fine-Tuning Pre-Trained Models:** The Shona corpus was used to refine pre-trained models such as BERT and T5. This included:
- **Transfer Learning:** Models are adjusted by honing them on Shona texts after they were first trained on languages with abundant resources.
- **Language-Specific Adjustments:** To properly represent the linguistic features of Shona, the tokenization and vocabulary should be changed.
- **Instructional Specifics:** Using the annotated corpus and supervised learning, the model parameters were optimized during the fine-tuning phase to increase summarization accuracy.

➤ *Hybrid Approach*

We integrated statistical and rule-based techniques to improve performance:

- **Rule-Based Extraction:** To help with sentence selection, rules unique to Shona language and syntax were implemented.
- **Statistical Weighting:** To rank sentences with more information, statistical techniques such as TF-IDF are applied.

E. Evaluation

➤ *Evaluation Metrics*

We used common measures to assess the summary performance: Measures the overlap in terms of n-grams, longest common subsequences, and word sequences between the generated summaries and reference summaries using ROUGE (Recall-Oriented Understudy for Gisting Evaluation).

- **Human Evaluation:** The summaries are evaluated for relevance, coherence, and grammatical accuracy by native Shona speakers.

➤ *Comparative Analysis*

The created models were contrasted with baseline techniques and earlier Shona text summarizing strategies. Important comparative elements were as follows:

- **Summary Quality:** Assessed using both human and ROUGE scores.
- **Linguistic Adequacy:** Evaluates how well Shona-specific linguistic elements, like syntax and morphological variances, are handled.

IV. RESULTS

The results of this study are summarized in Table 1, which shows the performance metrics of various models tested, including ROUGE-1, ROUGE-2, and ROUGE-L scores.

TABLE I. PERFORMANCE METRICS OF SUMMARIZATION MODELS

Model	ROUGE-1	ROUGE-2	ROUGE-L
Extractive Model	78.5%	63.2%	70.1%
Abstractive Model	82.3%	67.8%	74.5%
Hybrid Model	85.6%	71.4%	78.2%

The extractive model demonstrated substantial improvements in ROUGE-1 scores compared to traditional methods. The abstractive model showed better performance, particularly in ROUGE-2 and ROUGE-L metrics, indicating its ability to generate more coherent and contextually relevant summaries. The hybrid model, which combines the strengths of both extractive and abstractive approaches, achieved the highest scores across all ROUGE metrics.

The results of this study are summarized in Table 2. The values indicate the performance metrics of various models tested, including accuracy, precision, recall, and F1-score.

TABLE II. PERFORMANCE METRICS OF MACHINE LEARNING MODELS

Model	Accuracy	Precision	Recall	F1-Score
Convolutional NN	85.2%	84.6%	83.8%	84.2%
Recurrent NN	88.1%	87.4%	86.7%	87.0%
BERT	91.5%	91.0%	90.2%	90.6%
Hybrid Model	93.2%	92.8%	92.1%	92.4%

V. CONCLUSION

The work shows how sophisticated NLP techniques may be used to create an effective text summarizing system specifically designed for the Shona language. Summarization model performance significantly improved with the development of a large corpus of Shona literature and the optimization of pre-trained models. The hybrid model attained the highest performance indicators by fusing the extractive and abstractive techniques. This study tackles the difficulties associated with Shona text summarization, including the scarcity of annotated datasets and the intricacy of Shona morphology. Establishing a basis for forthcoming NLP investigations in marginalized languages, it fosters inclusiveness and linguistic multiplicity.

REFERENCES

- [1]. Vienna Li, Srinitha Sridharan, Sandeep Sethuraman, Georgios Avdis. "Predicting Recidivism With Machine Learning An Analysis of Risk Factors and Proposal of Preventions", Journal of Student Research, 2023
- [2]. Amy J. C. Trappey; Charles V. Trappey; Jheng-Long Wu; W. C. Wang; "Intelligent Compilation of Patent Summaries Using Machine Learning and Natural Language Processing Techniques", ADV. ENG. INFORMATICS, 2020.
- [3]. Liuqing Li; Jack Geissinger; William A. Ingram; Edward A. Fox; "Teaching Natural Language Processing Through Big Data Text Summarization with Problem-Based Learning", DATA AND INFORMATION MANAGEMENT, 2020.
- [4]. Ovishake Sen; Mohtasim Fuad; Md. Nazrul Islam; Jakaria Rabbi; Mehedi Masud; Md. Kamrul Hasan; Md. Abdul Awal; Awal Ahmed Fime; Md. Tahmid Hasan Fuad; Delowar Sikder; Md. Akil Raihan Iftee; "Bangla Natural Language Processing: A Comprehensive Analysis of Classical, Machine Learning, and Deep Learning-Based Methods", IEEE ACCESS, 2021. (IF: 3)
- [5]. III Robert E. Wray; James R. Kirk; John E. Laird; "Language Models As A Knowledge Source for Cognitive Agents", ARXIV-CS.AI, 2021.
- [6]. Y M Wazery; Marwa E Saleh; Abdullah Alharbi; Abdelmgeid A Ali; "Abstractive Arabic Text Summarization Based on Deep Learning", COMPUTATIONAL INTELLIGENCE AND NEUROSCIENCE, 2022. (IF: 3)
- [7]. Rabimba Karanjai; "Targeted Phishing Campaigns Using Large Scale Language Models", ARXIV, 2022. (IF: 3)
- [8]. Kiran Busch; Alexander Rochlitzer; Diana Sola; Henrik Leopold; "Just Tell Me: Prompt Engineering in Business Process Management", ARXIV-CS.AI, 2023.
- [9]. Weisong Sun; Chunrong Fang; Yudu You; Yun Miao; Yi Liu; Yuekang Li; Gelei Deng; Shenghan Huang; Yuchen Chen; Qunjun Zhang; Hanwei Qian; Yang Liu; Zhenyu Chen; "Automatic Code Summarization Via ChatGPT: How Far Are We?", ARXIV-CS.SE, 2023. (IF: 3)
- [10]. R. Sangeetha; J. Logeshwaran; Durgesh Srivastava; Pramod Vishwakarma; Satvik Vats; "Smart Performance Optimization of Natural Language Processing with Recursive Neural Networks", 2023 INTERNATIONAL CONFERENCE ON RESEARCH METHODOLOGIES IN ..., 2023.
- [11]. Rui Yang; Qingcheng Zeng; Keen You; Yujie Qiao; Lucas Huang; Chia-Chun Hsieh; Benjamin Rosand; Jeremy Goldwasser; Amisha D Dave; Tiarnan D. L. Keenan; Emily Y Chew; Dragomir Radev; Zhiyong Lu; Hua Xu; Qingyu Chen; Irene Li; "Ascle: A Python Natural Language Processing Toolkit for Medical Text Generation", ARXIV-CS.CL, 2023.