

ABSTRACT

Rafii Muhammad Fadhil, 51421205

PERFORMANCE COMPARISON OF THE BERTSUMEXT MODEL AND THE LEAD-3 BASELINE FOR AUTOMATIC TEXT SUMMARIZATION IN ENGLISH AND INDONESIAN USING PYTHON

Thesis, Department of Informatics, Faculty of Industrial Technology, Gunadarma University, 2025.

Keywords: Automatic Text Summarization, BERTSumExt, Lead-3, IndoSum, CNN/DailyMail, ROUGE. Python

.(xii + 73 + appendix)

The rapid growth of digital information has driven the demand for Automatic Text Summarization (ATS) systems capable of producing concise and informative summaries. This study aims to compare the performance of two extractive approaches, namely the simple baseline Lead-3 and the modern Transformer-based model BERTSumExt, implemented using Python, in summarizing English and Indonesian texts. The datasets used are CNN/DailyMail as a representative of semi-abstractive English texts and IndoSum as a standard extractive dataset in Indonesian. The research stages include text preprocessing, labeling strategies (fuzzy matching for CNN/DailyMail and exact matching for IndoSum), method implementation, limited training, and evaluation using ROUGE metrics. Experimental results show that on CNN/DailyMail, Lead-3 achieved ROUGE-1 of 0.2921, ROUGE-2 of 0.1117, and ROUGE-L of 0.1915, outperforming BERTSumExt, which only reached ROUGE-1 of 0.2388, ROUGE-2 of 0.0569, and ROUGE-L of 0.1427. On IndoSum, the difference is even more pronounced, with Lead-3 reaching ROUGE-1 of 0.6687, ROUGE-2 of 0.6046, and ROUGE-L of 0.6515, while BERTSumExt lagged behind with ROUGE-1 of 0.3327, ROUGE-2 of 0.1822, and ROUGE-L of 0.2442. These findings highlight that model effectiveness is strongly influenced by dataset characteristics; simple methods such as Lead-3 are more effective on extractive datasets, while BERTSumExt still offers semantic flexibility on semi-abstractive datasets, although it has not yet surpassed the baseline. This study contributes by presenting a cross-lingual comparative analysis of the effectiveness of simple and modern extractive methods in automatic text summarization.

References (2017-2025)