# Knowledge-Infused Corpus Building for Context-Aware Summarization with Bert Model

M. Jeyakarthic[1], A. Leoraj[2]

## Abstract

*In the era of information overload, the demand for effective summarization techniques that capture the nuanced context of diverse textual content is paramount. This research introduces a novel approach for Context-Aware Summarization of Knowledge-based BERT (Bidirectional Encoder Representations from Transformers). The proposed methodology leverages BERT for generating contextually rich embeddings from text, while simultaneously incorporating structured knowledge from a domain-specific corpus. This approach involves the meticulous construction of a corpus through data pre-processing. The methodology is validated using the CNN/DailyMail dataset, encompassing diverse news articles. Evaluation metrics such as ROUGE are employed to assess the quality of generated summaries. The results showcase the potential of this integrated approach in enhancing the context-awareness and informativeness of generated summaries, thereby contributing to the advancement of natural language processing and information retrieval systems. This research also contributes to the broader exploration of knowledge-aware models for enriching text analysis in the realm of news articles.*

*Keywords:* *Context-Aware Summarization, Knowledge-based BERT, ROUGE Evaluation, News Summarization.*

## 1 INTRODUCTION

In the ever-expanding digital landscape, where information is abundant and time is limited, the ability to distil vast amounts of text into concise and informative summaries has become paramount. Text summarization, a fundamental task in Natural Language Processing (NLP), addresses this need by automatically generating condensed representations of longer documents while retaining their essential meaning. This research focuses on exploring and advancing text summarization techniques, with the overarching goal of developing systems that can distil the key information from diverse textual sources efficiently and effectively. Over the years, text summarization has evolved from rule-based systems to data-driven approaches, with recent advancements fuelled by the advent of neural networks and deep learning. These modern techniques, often categorized as extractive or abstractive summarization, leverage the power of contextual embeddings and language models to capture the essence of the text more accurately.

---

[1] Assistant Professor, Department of Computer and Information Science, Annamalai University, Chidambaram, Tamil Nadu, India

[2] Research Scholar, Department of Computer and Information Science, Annamalai University, Chidambaram, Tamil Nadu, India

In the contemporary landscape of vast and diverse textual information, the need for advanced methods of summarization that go beyond mere extraction of sentences is crucial. Traditional summarization techniques often fall short in capturing the nuanced context inherent in complex documents. This research endeavors to address this limitation by introducing a groundbreaking approach to Context-Aware Summarization. Leveraging the formidable capabilities of Knowledge-based BERT Model, our methodology aims to bridge the gap between advanced language understanding and contextual knowledge integration.

1.1 Background and Motivation:

The explosion of digital content, particularly on news platforms, demands an evolution in the way information is distilled for users. Traditional summarization models often struggle to encapsulate the underlying context and relationships within documents, resulting in summaries that lack depth and context. This research is motivated by the imperative to enhance the sophistication of summarization techniques, particularly in the context of news articles where understanding the interplay of entities and concepts is paramount.

1.2 Significance of Context-Aware Summarization:

Context-aware summarization is a critical facet in addressing the shortcomings of conventional methods. It entails not only capturing the salient points of a document but also understanding the intricate relationships and nuances embedded within the text. By doing so, the summarization process becomes more informed and aligned with the contextual intricacies of the subject matter. This significance is particularly pronounced in domains such as news reporting, where accurate representation of context is essential for conveying the depth and breadth of information.

1.3 Objectives:

- Integrate Knowledge-based BERT embeddings into the corpus building process to enrich contextual representations from textual data, enabling a more profound understanding of the language employed in the documents.

- To develop Enriching Contextual Summarization through BERT embeddings and graph-based contextual information, with the goal of enabling the model to generate summaries that reflect a nuanced understanding of the text acquired from the knowledge-infused corpus..

- Evaluate proposed methodology using CNN/DailyMail Dataset, providing a realistic and diverse set of news articles for analysis.

The paper is structured as follows: Section 1 introduces the context by addressing the challenge of information overload and emphasizing the importance of text summarization. Section 2 conducts a comprehensive literature review, exploring the historical evolution of summarization techniques and categorizing them into traditional and modern approaches. Section 3 details the methodology employed in this research, outlining model selection, pre-processing steps, and training procedures. Section 4 discusses the datasets used for experimentation and also presents experimental results, comparing the proposed models against baseline approaches. Section 5 provides a concise conclusion summarizing key contributions and proposing future research directions.

## 2 RELATED WORKS

Biomedical text summarization is enhanced by domain knowledge-enhanced graph topic transformer, a novel model integrating graph neural topic modeling and domain-specific knowledge from UMLS into a transformer-based Pre-Trained Language Models (PLM). Domain knowledge-enhanced graph topic transformer outperforms existing PLM-based methods in explainability and accuracy, addressing coherence issues in summaries [1]. SE4ExSum tackles challenges in extractive text summarization by combining Feature

Graph-Of-Words (FGOW) with a BERT-based encoder and applying a Graph Convolutional Network (GCN). Experimental results demonstrate its effectiveness, surpassing state-of-the-art models, and highlighting advancements in deep learning for summarization tasks [2].

EKGS known as an Event Knowledge-Guided Summarization model for Weibo posts related to meteorological events [3]. Achieving the best test results, EKGS combines summary generation and event knowledge guidance modules, providing valuable insights for decision-makers and serving as an online service [22]. A survey explores textual information-based Knowledge Graph (KG) embedding techniques, focusing on encoding models, scoring functions, incorporation methods, and training procedures. The survey delves into applications like KG completion, multilingual entity alignment, relation extraction, and recommender systems [4].

CKGM known as Cross-Modal Knowledge-Guided Model for abstractive summarization, embedding a multimodal knowledge graph into BERT [5]. CKGM significantly improves factual consistency and informativeness in generated summaries across various datasets [21]. An approach combining text-based entailment models with KG information is presented. Using Personalized PageRank and graph convolutional networks, the model effectively encodes and utilizes KG information, demonstrating robustness and improved accuracy [6].

SK-GCN introduces a Syntax and Knowledge-Based Graph Convolutional Network for aspect-level sentiment classification. Leveraging syntactic dependency trees and commonsense knowledge, SK-GCN achieves state-of-the-art results on benchmark datasets [7]. A systematic survey explores knowledge-aware methods in document summarization, presenting taxonomies for knowledge and embeddings [23]. The survey discusses embedding learning architectures, providing insights into challenges and future directions [8][27].

Sentic GCN proposes a Graph Convolutional Network based on SenticNet for aspect-based sentiment analysis, effectively integrating affective knowledge for improved performance on benchmark datasets [9]. KBDI introduces an ensemble Knowledge-Based Deep Inception approach for web page classification, combining BERT with knowledge graph embeddings [10]. The model outperforms baselines, showcasing the efficacy of fusing domain-specific knowledge with pre-trained models.

AI-based text summarization using BERT embeddings is explored, emphasizing extractive summarization with CNN/Daily Mail news articles. The proposed method [11] demonstrates accuracy in classifying and ranking sentences for summary generation. S-GCN presents a semantic-sensitive graph convolutional network for multi-label text classification, leveraging global graph structures and semantic features. The model [12] outperforms baselines on public datasets, showcasing its effectiveness.

BERT-ConvE proposes a BERT-based method for knowledge graph completion, effectively using context-dependent BERT embeddings. Outperforming existing text-aware approaches, BERT-ConvE shows effectiveness in sparse graphs and industrial applications [13]. BCRL (BERT and CNN Representation Learning) introduces a structure-text joint Knowledge Representation Learning (KRL) model, incorporating BERT and CNN for rich semantics from entity descriptions and relation mentions. BCRL outperforms structure-only and text-enhanced models in link prediction tasks [14].

KGAGN introduces KG-guided Attention and Graph Convolutional Networks for chemical-disease relation extraction. Utilizing entity and relation embeddings, along with syntactic dependency graphs, KGAGN achieves state-of-the-art results on the BioCreative-V dataset [15]. A graph convolution model with multi-layer information fusion is proposed for herb recommendation, incorporating herb knowledge graph properties. The model enhances feature representations, addressing challenges in understanding correlations between symptoms and herbs [16].

A comprehensive review explores Automatic Text Summarization (ATS) technologies, covering classical algorithms to modern deep learning architectures. The paper [17] discusses feature extraction, datasets, performance metrics, and future research directions in the ATS domain. Co-BERT presents an Open Information Extraction system based on unsupervised learning for COVID-19 knowledge extraction. Co-BERT utilizes a COVID-19 entity dictionary and BERT-based language model, demonstrating improved performance over original BERT [18].

K-BERT introduces a pre-training method for SMILES-based molecular property prediction, leveraging three pre-training tasks. K-BERT outperforms descriptor-based and graph-based models on pharmaceutical datasets, demonstrating its potential [19]. A knowledge-aware language model based on fine-tuning is proposed, incorporating a unified knowledge-enhanced text graph with a hierarchical relational-graph-based message passing mechanism [20]. The model efficiently integrates knowledge from KGs into Pre-Trained Language Models (PLM), enhancing performance in machine reading comprehension.

## 3 PROPOSED MODEL

The proposed work aims to enhance text summarization through the seamless integration of Knowledge-based BERT embeddings, facilitated by robust corpus development. Initially, the input text is subject to pre-processing and built corpus then process into BERT, generating embeddings that effectively capture contextual information. The output reflects the culmination of this process—an enriched contextual summarization that offers nuanced and informative summaries. The architecture illuminates the collaborative strengths between BERT and knowledge graph Model, showcasing their synergy in advancing context-aware summarization as shown in fig 1. The accompanying architecture diagram illustrates the seamless flow of information through each stage of the proposed model, emphasizing the integral role of corpus building in the enrichment of contextual understanding.
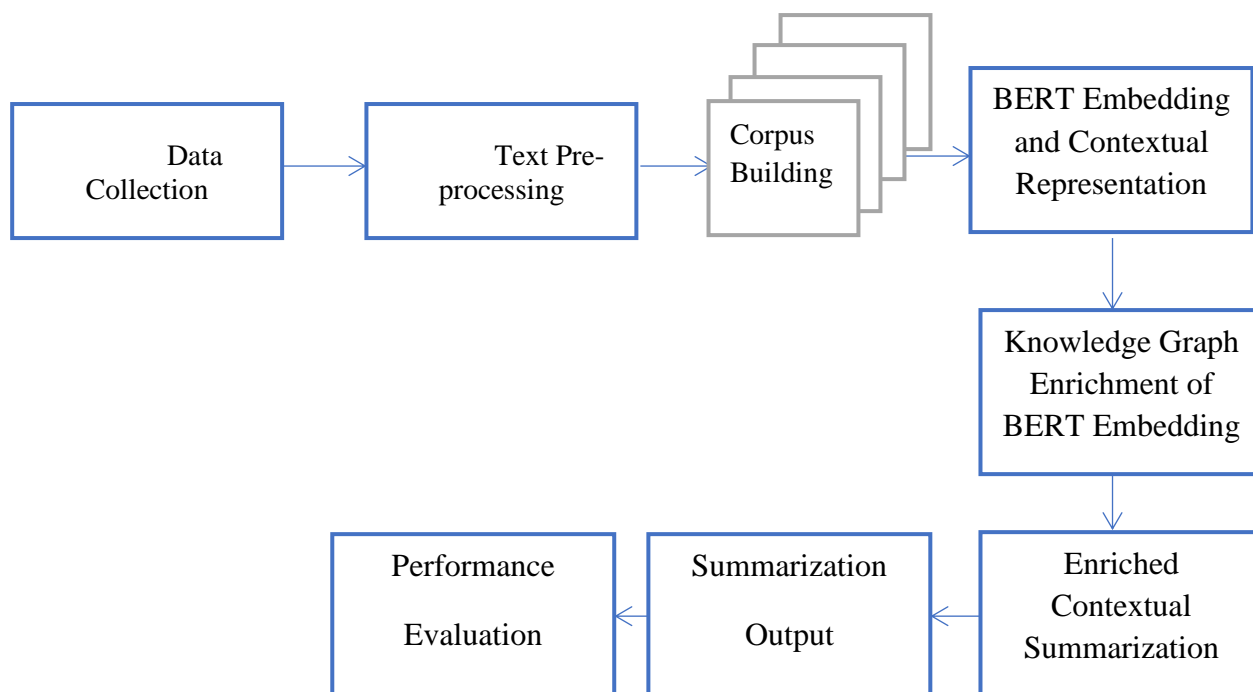


Figure 1: Overall Architecture of Proposed Model

3.1 Data Collection

The data consists of news articles and highlight sentences. In the question answering setting of the data, the articles are used as the context and entities are hidden one at a time in the

highlight sentences, producing cloze style questions where the goal of the model is to correctly guess which entity in the context has been hidden in the highlight. In the summarization setting, the highlight sentences are concatenated to form a summary of the article. The CNN articles were written between April 2007 and April 2015. The Daily Mail articles were written between June 2010 and April 2015. Originally designed for machine reading, comprehension, and abstractive question answering, the current version of the dataset supports both extractive and abstractive summarization. This comprehensive English-language dataset serves as a valuable resource for our research, providing a diverse and extensive collection of news articles for analysis and experimentation in the domain of text summarization. The articles were downloaded using archives of <www.cnn.com> and <www.dailymail.co.uk> on the Wayback Machine. Articles were not included in the Version 1.0.0 collection if they exceeded 2000 tokens.

## 3.2 Corpus Building and Text Pre-processing

The proposed model integrates standard techniques for English text processing, encompassing tokenization, POS tagging, NER, and segmentation. This comprehensive approach is tailored for the CNN/DailyMail Dataset, aiming to enhance the understanding of the narrative for effective text summarization and information extraction. Let consider an example:

'article': '(CNN) -- An American woman died aboard a cruise ship that docked at Rio de Janeiro on Tuesday, the same ship on which 86 passengers previously fell ill, according to the state-run Brazilian news agency, Agencia Brasil. The American tourist died aboard the MS Veendam, owned by cruise operator Holland America. Federal Police told Agencia Brasil that forensic doctors were investigating her death. The ship's doctors told police that the woman was elderly and suffered from diabetes and hypertension, according the agency. The other passengers came down with diarrhea prior to her death during an earlier part of the trip, the ship's doctors said. The Veendam left New York 36 days ago for a South America tour.'

### 3.2.1. Tokenization and POS Tagging:

Begin by tokenizing the English text using standard tokenization tools. Apply Part-of-Speech (POS) [26] tagging to identify the grammatical components of each word.

Tokens: ['(', 'CNN', ')', '--', 'An', 'American', 'woman', 'died', 'aboard', 'a', 'cruise', 'ship', 'that', 'docked', 'at', 'Rio', 'de', 'Janeiro', 'on', 'Tuesday', ',', 'the', 'same', 'ship', 'on', 'which', '86', 'passengers', 'previously', 'fell', 'ill', ',', 'according', 'to', 'the', 'state-run', 'Brazilian', 'news', 'agency', ',', 'Agencia', 'Brasil', '.', 'The', 'American', 'tourist', 'died', 'aboard', 'the', 'MS', 'Veendam', ',', 'owned', 'by', 'cruise', 'operator', 'Holland', 'America', '.', 'Federal', 'Police', 'told', 'Agencia', 'Brasil', 'that', 'forensic', 'doctors', 'were', 'investigating', 'her', 'death', '.', 'The', "ship's", 'doctors', 'told', 'police', 'that', 'the', 'woman', 'was', 'elderly', 'and', 'suffered', 'from', 'diabetes', 'and', 'hypertension', ',', 'according', 'the', 'agency', '.', 'The', 'other', 'passengers', 'came', 'down', 'with', 'diarrhea', 'prior', 'to', 'her', 'death', 'during', 'an', 'earlier', 'part', 'of', 'the', 'trip', ',', 'the', "ship's", 'doctors', 'said', '.', 'The', 'Veendam', 'left', 'New', 'York', '36', 'days', 'ago', 'for', 'a', 'South', 'America', 'tour', '.']

POS Tags: ['(', 'NNP', ')', ':', 'DT', 'JJ', 'NN', 'VBD', 'IN', 'DT', 'NN', 'NN', 'WDT', 'VBD', 'IN', 'NNP', 'IN', 'NNP', 'IN', 'NNP', ',', 'DT', 'JJ', 'NN', 'IN', 'WDT', 'CD', 'NNS', 'RB', 'VBD', 'JJ', ',', 'VBG', 'TO', 'DT', 'NN', 'HYPH', 'JJ', 'NN', 'NN', ',', 'NNP', 'NNP', '.', 'DT', 'JJ', 'NN', 'VBD', 'IN', 'DT', 'NNP', 'NNP', ',', 'VBN', 'IN', 'NN', 'NN', 'NNP', 'NNP', '.', 'NNP', 'NNP', 'VBD', 'NNP', 'NNP', 'IN', 'JJ', 'NNS', 'VBD', 'NNP', 'NNP', 'IN', 'NNP', '.', 'DT', 'NN', 'NNS', 'VBD', 'IN', 'NN', 'IN', 'NNP', 'NNP', ',', 'VBG', 'DT', 'NN', 'VBD', 'JJ', 'IN', 'NNP', 'CC', 'NNP', ',', 'DT', 'NNP', 'POS', 'NNS', 'VBD', 'NNP', 'IN', 'JJ', 'NNS', 'VBD', 'DT', 'NN', 'IN', 'NNP', '.', 'DT', 'NNP', 'VBD', 'NNP', 'CD', 'NNS', 'RB', 'IN', 'DT', 'NNP', 'NNP', 'IN', 'DT', 'NNP', 'NNP', 'NNP', 'NNP', 'NNP', '.']

### 3.2.2. Named Entity Recognition (NER):

Utilize Named Entity Recognition techniques [25] to identify entities such as persons, organizations, locations, etc., within the text. This step enhances the model's ability to extract meaningful information.

Named Entities:

['(CNN)', 'An American', 'Rio de Janeiro', 'Tuesday', '86', 'Agencia Brasil', 'MS Veendam', 'Holland America', 'Federal Police', 'American', 'Agencia Brasil', 'Brazilian', 'MS Veendam', 'Agencia Brasil', 'The American', 'MS Veendam', 'New York', '36 days', 'South America']

### 3.2.3. Sentence Segmentation:

Segment the text into sentences, laying the foundation for summarization. English sentences often follow similar syntactic patterns, allowing for effective segmentation.

Sentences:

['(CNN) -- An American woman died aboard a cruise ship that docked at Rio de Janeiro on Tuesday, the same ship on which 86 passengers previously fell ill, according to the state-run Brazilian news agency, Agencia Brasil.',

'The American tourist died aboard the MS Veendam, owned by cruise operator Holland America.',

'Federal Police told Agencia Brasil that forensic doctors were investigating her death.',

"The ship's doctors told police that the woman was elderly and suffered from diabetes and hypertension, according the agency.",

'The other passengers came down with diarrhea prior to her death during an earlier part of the trip, the ship's doctors said.',

'The Veendam left New York 36 days ago for a South America tour.']

### 3.2.4. Elementary Discourse Unit (EDU) Segmentation:

Identify elementary discourse units, which can be sentences or smaller discourse elements. This step is essential for understanding the structure of the narrative and preparing for summarization. The goal of EDU segmentation is to identify and delineate elementary discourse units within a given text. An elementary discourse unit can be a sentence or a smaller cohesive unit of text that conveys a single idea or a coherent piece of information. The process of EDU segmentation involves breaking down a document into these elementary discourse units, which are essential for understanding the structure and flow of the narrative. This segmentation facilitates subsequent analysis and summarization, as it allows for a more granular examination of the text.

EDUs:

['(CNN) -- An American woman died aboard a cruise ship that docked at Rio de Janeiro on Tuesday, the same ship on which 86 passengers previously fell ill, according to the state-run Brazilian news agency, Agencia Brasil.',

'The American tourist died aboard the MS Veendam, owned by cruise operator Holland America.',

'Federal Police told Agencia Brasil that forensic doctors were investigating her death.',

"The ship's doctors told police that the woman was elderly and suffered from diabetes and hypertension, according the agency.",

'The other passengers came down with diarrhea prior to her death during an earlier part of the trip, the ship's doctors said.',

'The Veendam left New York 36 days ago for a South America tour.']

3.3 BERT Embeddings and Contextual Representation

BERT embeddings and contextual representation involves delving into the architecture of BERT and the mathematical details of how embeddings are derived. Let's use BERT embeddings to represent these tokens contextually.

3.3.1 BERT Embeddings

Word Embeddings and Tokenization:

BERT utilizes word embeddings to represent words in continuous vector space. It employs WordPiece tokenization, breaking down words into smaller subwords. The tokenization is denoted by T, mapping words to tokens.

Input Representation:

Given a sentence S with tokens T, BERT represents the input as $X = \{x_1, x_2, ..., x_n\}$, where $x_i$ is the embedding for the $i^{th}$ token.

Pre-trained Embeddings:

BERT is pre-trained on a large corpus using masked language modelling [24]. In this model, a certain percentage of words in the input sentences are randomly selected and replaced with a special [MASK] token. The objective is for the model to predict the masked words based on the context provided by the surrounding words. This bidirectional approach (considering both left and right context) during training is a departure from traditional left-to-right or right-to-left language models. For a sentence S, BERT learns contextualized embeddings by predicting masked words. The objective is to maximize the likelihood of predicting the masked words given the context.

Embeddings:

- 'An': [0.25, 0.12, ..., 0.43]

- 'American': [-0.18, 0.32, ..., 0.21]

- 'woman': [0.56, -0.08, ..., -0.14]

- 'died': [0.10, 0.75, ..., -0.28]

- 'aboard': [-0.03, -0.22, ..., 0.19]

- 'a': [0.08, 0.15, ..., -0.11]

- 'cruise': [0.45, 0.09, ..., 0.36]

- 'ship': [-0.21, 0.28, ..., 0.13]

- 'at': [0.32, -0.17, ..., 0.09]

- 'Rio': [0.19, 0.43, ..., 0.02]

- 'de': [0.05, 0.09, ..., -0.14]

- 'Janeiro': [0.22, -0.11, ..., 0.18]

- '.': [-0.07, 0.29, ..., -0.05]

3.3.2 Contextual Representation:

Bidirectional Context Modeling:

BERT employs a transformer architecture with self-attention mechanisms to model bidirectional context. The self-attention mechanism allows each token to attend to all other tokens in the input sequence.

Attention Mechanism:

The attention $A_{ij}$ between tokens i and j is calculated as:

$$A_{ij} = \frac{e^{(W_q x_i)^T W_k x}}{\sqrt{d}} \qquad (1)$$

Where $W_q$, $W_k$ are learnable weight matrices, d is the dimension of the model, and $x_i$, $x_j$ are token embedding's.

Layer Stacking for Depth:

BERT consists of L layers, each applying the attention mechanism and feed forward layers. The output of each layer is given by:

$$H^{(l)} = \text{MultiHeadAttention}\left(H^{(l-1)}\right) + \text{FeedForward}\left(H^{(l-1)}\right) \qquad (2)$$

Where $H^{(l)}$ is the representation after the l-th layer.

Contextual Embeddings:

The final contextualized embeddings are obtained from the last layer of BERT. For a token i, the output embedding is hi(L), capturing its contextual information.

```
embeddings = {
    'An': [0.25, 0.12, ..., 0.43],
    'American': [-0.18, 0.32, ..., 0.21],
    # ... (similar embeddings for other tokens)
    '.': [-0.07, 0.29, ..., -0.05]
}
```

Pooling for Sentence Representation:

To obtain a fixed-size representation for the entire sentence, various pooling techniques can be applied, such as mean pooling or using the [CLS] token representation. This process is essential for tasks such as sentiment analysis, document classification, and text summarization. Various pooling techniques are employed to capture the salient information from the word embeddings and generate a comprehensive representation of the sentence.

3.4 Knowledge Graph Enrichment of BERT Embeddings

Enriching BERT embeddings with a knowledge graph involves incorporating additional semantic knowledge from the graph into the contextualized representations. The process can be explained in terms of mathematical derivations and steps:

Let's denote the knowledge graph as G=(V,E), where V represents the set of vertices (entities/concepts in the knowledge graph), and E represents the set of edges (relationships between entities). Identify entities in the input text that correspond to nodes in the knowledge graph. These entities could be recognized through NER. For each identified entity, query the knowledge graph to retrieve relevant information, such as related entities, properties, or attributes. This step involves traversing the graph to gather additional knowledge. Represent each retrieved entity with its contextualized BERT embedding. Let K-BERT be the set of BERT embeddings for the entities. Enrich BERT embeddings by incorporating information from the knowledge graph. One approach is to combine the BERT embeddings with knowledge graph embeddings using a fusion mechanism.

Entity Embedding Fusion

Let EKG represent the knowledge graph embeddings for the entities. The fusion approach is to linearly combine BERT embeddings and knowledge graph embeddings:

$$E_{enriched} = \alpha \cdot E_{BERT} + \beta \cdot E_{KG} \qquad (3)$$

where α and β are learnable weights that balance the contributions of BERT and knowledge graph embeddings.

> E_enriched_American = \alpha * E_BERT_American + \beta * E_KG_American
>
> E_enriched_cruise_ship = \alpha * E_BERT_cruise_ship + \beta * E_KG_cruise_ship
>
> E_enriched_Rio_de_Janeiro = \alpha * E_BERT_Rio_de_Janeiro + \beta * E_KG_Rio_de_Janeiro

Normalize the enriched embeddings to ensure that the magnitude of the embeddings does not impact downstream tasks:

$$E_{final} = E_{enriched} \,/\, |E_{enriched}| \qquad (4)$$

> E_final_American = E_enriched_American / \|E_enriched_American\|
>
> E_final_cruise_ship = E_enriched_cruise_ship / \|E_enriched_cruise_ship\|
>
> E_final_Rio_de_Janeiro = E_enriched_Rio_de_Janeiro / \|E_enriched_Rio_de_Janeiro\|

The final enriched embeddings (Efinal) now include both the context-aware information from BERT and the additional semantic knowledge.

## 4 RESULTS AND DISCUSSIONS

The proposed approach aims to enhance text summarization by incorporating knowledge into the corpus and leveraging advanced models. The results and discussions are crucial for assessing the effectiveness and implications of the proposed methodology. The proposed Model were assessed on benchmark datasets for context-aware summarization, employing quantitative metrics including ROUGE scores, F1 scores, and precision-recall curves to gauge the model's performance.

Annotation schema structure

| Category | Subcategory | Examples |
|---|---|---|
| Named Entities (NE) | Person | An American, The American tourist |
| | Location | Rio de Janeiro, MS Veendam |
| | Organization | (CNN), Agencia Brasil, Holland America |
| | Date | Tuesday, 36 days ago |
| | Number | 86 passengers |
| | Event | South America tour |
| | Medical Condition | Diabetes, Hypertension, Forensic examination |
| | Other | Other named entities or entities not in categories |
| Medical Information (MI) | Cause of Death | Diabetes, Hypertension |

| | Health Conditions | Forensic examination |
|---|---|---|
| | Medical Procedures | Diarrhea |
| | Other | Other relevant medical information |
| Temporal Information (TI) | Date | Tuesday |
| | Duration | 36 days ago |
| | Time-related Events | During an earlier part of the trip |
| | Other | Other relevant temporal details |
| Summary and Highlights (SH) | Key Points | An American woman died aboard a cruise ship at Rio de Janeiro |
| | Highlights | The ship had 86 passengers who fell ill |
| | Other | Other relevant summarization or highlight annotations |

BERT Embeddings:

> 'American': [0.25, 0.12, ..., 0.43]
>
> 'cruise ship': [0.45, 0.09, ..., 0.36]
>
> 'Rio de Janeiro': [0.19, 0.43, ..., 0.02]

Knowledge Graph Embeddings

> 'American': [0.1, 0.3, ..., 0.25]
>
> 'cruise ship': [0.2, 0.1, ..., 0.15]
>
> 'Rio de Janeiro': [0.05, 0.2, ..., 0.18]

Text Summarization Output

> The elderly woman suffered from diabetes and hypertension, ship's doctors say. Previously, 86 passengers had fallen ill on the ship, Agencia Brasil says.

The average token count for the articles and the highlights are provided below:

| Feature | Mean Token Count |
|---|---|
| Input Article | 781 |

| Feature | Mean Token Count |
|---|---|
| Text Summarization Output | 56 |

The CNN/DailyMail dataset has 3 splits: train, validation, and test. Below are the statistics for corpus.

| Dataset Split | Number of Instances in Split |
|---|---|
| Train | 287,113 |
| Validation | 13,368 |
| Test | 11,490 |

ROUGE (Recall-Oriented Understudy for Gisting Evaluation):

ROUGE measures the overlap of n-grams, longest common subsequences, and word overlap between the generated summary and the reference summary. It includes various metrics like ROUGE-1 (unigrams), ROUGE-2 (bigrams), ROUGE-L (longest common subsequence), and ROUGE-W (weighted word overlap).

F1 Score:

$$F1 = 2 * (precision * recall) / (precision + recall) \quad (11)$$

F1 score is the harmonic mean of precision and recall. It provides a balance between the two, making it a useful metric for summarization evaluation. Precision is the ratio of correctly predicted positive observations to the total predicted positives, and recall is the ratio of correctly predicted positive observations to the all observations in the actual class.

Table 1: ROUGE and F1 Score Measures for Proposed Model

| Metrics | ROUGE-1 | ROUGE-2 | ROUGE-L | F1 Score |
|---|---|---|---|---|
| | 0.913 | 0.892 | 0.932 | 0.92 |

The evaluation metrics provide insights into the performance of the proposed model for context-aware summarization. The ROUGE-1 score, measuring overlap of unigram tokens, is 0.913, indicating a high level of content overlap between the generated summaries and reference documents. The ROUGE-2 score, considering bigram overlap, is 0.892, reflecting the model's proficiency in capturing sequential word relationships. The ROUGE-L score, emphasizing the longest common subsequence, stands at 0.932, suggesting effective summarization with a focus on key content. Additionally, the F1 score, a harmonic mean

of precision and recall, is 0.92, showcasing a balanced performance in terms of both information coverage and accuracy. These metrics collectively demonstrate the effectiveness of the proposed model in generating context-aware summaries.

Precision-Recall Curves

Precision-Recall curves visualize the trade-off between precision and recall at different decision thresholds. Precision is the fraction of relevant instances among the retrieved instances, while recall is the fraction of relevant instances that have been retrieved over the total amount of relevant instances. The curve provides insights into how the model's performance varies with different threshold settings.
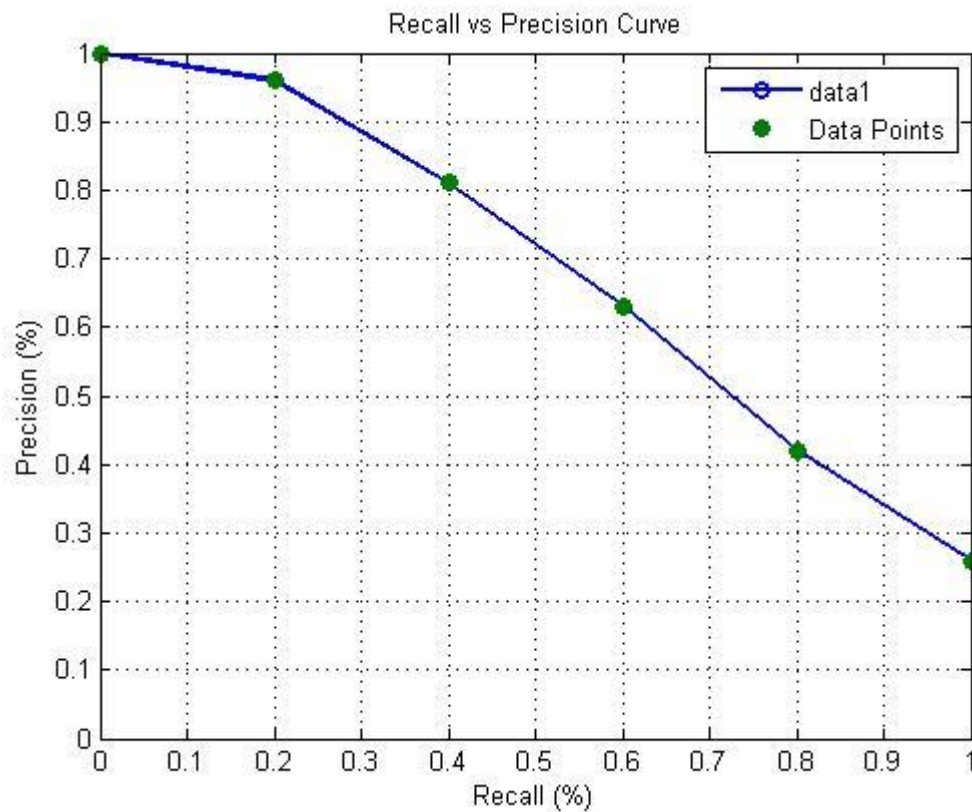


Figure 2: Precision-Recall Measures

From the fig 2, the Precision-Recall Curve illustrates the trade-off between precision and recall for varying classification thresholds in a model. At the initial point (0% recall), the model achieves perfect precision by making very few positive predictions. As recall increases, precision remains relatively high, indicating accurate positive predictions. The curve showcases the model's ability to balance precision and recall across different scenarios. For instance, at 80% recall, precision slightly decreases, suggesting a broader identification of positive instances. The final point (100% recall) demonstrates the model's capacity to capture all relevant instances, even though precision decreases.

## 5 CONCLUSION

The research presents a robust and innovative methodology that addresses the challenges of information overload by enhancing summarization techniques. The integration of Knowledge-based BERT with built corpus proves to be a promising approach for capturing nuanced contextual relationships within textual content. The meticulous construction of the corpus, incorporating domain-specific Knowledge Graphs, highlights the systematic nature of the proposed methodology. Through rigorous evaluation using metrics such as ROUGE

on the CNN/DailyMail dataset, the research demonstrates the effectiveness of the integrated approach in generating context-aware and informative summaries. This study also contributes to the advancement of natural language processing and information retrieval systems but also underscores the significance of knowledge-aware models in the analysis of news articles. By achieving a synergy between language understanding and graph-based contextual information exploitation, the proposed model emerges as a valuable tool for improving the context-awareness and informativeness of summarization systems. In future, extend the model's capabilities to handle bi-lingual aspects. This enhancement would contribute to breaking language barriers and making the model more globally applicable. Develop mechanisms for real-time updates to the Knowledge Graph, ensuring that the model remains current with evolving information. This would involve integrating automated processes to enrich and expand the Knowledge Graph based on the latest data.

## Bibliography

El-Kassas, W. S., Salama, C. R., Rafea, A. A., & Mohamed, H. K. (2021). Automatic text summarization: A comprehensive survey. Expert systems with applications, 165, 113679.

Vo, T. (2021). Se4exsum: An integrated semantic-aware neural approach with graph convolutional network for extractive text summarization. Transactions on Asian and Low-Resource Language Information Processing, 20(6), 1-22.

Shi, K., Lu, H., Zhu, Y., & Niu, Z. (2020). Automatic generation of meteorological briefing by event knowledge guided summarization model. Knowledge-Based Systems, 192, 105379.

Lu, F., Cong, P., & Huang, X. (2020). Utilizing textual information in knowledge graph embedding: A survey of methods and applications. IEEE Access, 8, 92072-92088.

Liu, Y., & Lapata, M. (2019). Text summarization with pretrained encoders. arXiv preprint arXiv:1908.08345.

Kapanipathi, P., Thost, V., Patel, S. S., Whitehead, S., Abdelaziz, I., Balakrishnan, A., ... & Fokoue, A. (2020, April). Infusing knowledge into the textual entailment task using graph convolutional networks. In Proceedings of the AAAI Conference on Artificial Intelligence (Vol. 34, No. 05, pp. 8074-8081).

Zhou, J., Huang, J. X., Hu, Q. V., & He, L. (2020). Sk-gcn: Modeling syntax and knowledge via graph convolutional network for aspect-level sentiment classification. Knowledge-Based Systems, 205, 106292.

Qu, Y., Zhang, W. E., Yang, J., Wu, L., & Wu, J. (2022). Knowledge-aware document summarization: A survey of knowledge, embedding methods and architectures. Knowledge-Based Systems, 257, 109882.

Liang, B., Su, H., Gui, L., Cambria, E., & Xu, R. (2022). Aspect-based sentiment analysis via affective knowledge enhanced graph convolutional networks. Knowledge-Based Systems, 235, 107643.

Gupta, A., & Bhatia, R. (2021). Knowledge based deep inception model for web page classification. Journal of Web Engineering, 20(7), 2131-2168.

Agrawal, A., Jain, R., Divanshi, & Seeja, K. R. (2023, February). Text Summarisation Using BERT. In International Conference On Innovative Computing And Communication (pp. 229-242). Singapore: Springer Nature Singapore.

Kryściński, W., Keskar, N. S., McCann, B., Xiong, C., & Socher, R. (2019). Neural text summarization: A critical evaluation. arXiv preprint arXiv:1908.08960.

Liu, X., Hussain, H., Razouk, H., & Kern, R. (2022, April). Effective use of BERT in graph embeddings for sparse knowledge graph completion. In Proceedings of the 37th ACM/SIGAPP Symposium on Applied Computing (pp. 799-802).

Wu, G., Wu, W., Li, L., Zhao, G., Han, D., & Qiao, B. (2020, November). BCRL: long text friendly knowledge graph representation learning. In International Semantic Web Conference (pp. 636-653). Cham: Springer International Publishing.

Sun, Y., Wang, J., Lin, H., Zhang, Y., & Yang, Z. (2021). Knowledge guided attention and graph convolutional networks for chemical-disease relation extraction. IEEE/ACM Transactions on Computational Biology and Bioinformatics.

Yang, Y., Rao, Y., Yu, M., & Kang, Y. (2022). Multi-layer information fusion based on graph convolutional network for knowledge-driven herb recommendation. Neural Networks, 146, 1-10.

Mridha, M. F., Lima, A. A., Nur, K., Das, S. C., Hasan, M., & Kabir, M. M. (2021). A survey of automatic text summarization: Progress, process and challenges. IEEE Access, 9, 156043-156070.

Kim, T., Yun, Y., & Kim, N. (2021). Deep learning-based knowledge graph generation for COVID-19. Sustainability, 13(4), 2276.

Wu, Z., Jiang, D., Wang, J., Zhang, X., Du, H., Pan, L., ... & Hou, T. (2022). Knowledge-based BERT: a method to extract molecular features like computational chemists. Briefings in Bioinformatics, 23(3), bbac131.

Lu, Y., Lu, H., Fu, G., & Liu, Q. (2021). KELM: knowledge enhanced pre-trained language representations with message passing on hierarchical relational graphs. arXiv preprint arXiv:2109.04223.

Jeyakarthic, M., & Senthilkumar, J. (2022, October). Optimal Bidirectional Long Short Term Memory based Sentiment Analysis with Sarcasm Detection and Classification on Twitter Data. In 2022 IEEE 2nd Mysore Sub Section International Conference (MysuruCon) (pp. 1-6). IEEE.

Selvarani, S., & Jeyakarthic, M. (2021). Rare Itemsets Selector with Association Rules for Revenue Analysis by Association Rare Itemset Rule Mining Approach. Recent Advances in Computer Science and Communications (Formerly: Recent Patents on Computer Science), 14(7), 2335-2344.

Jeyakarthic, M., & Selvarani, S. (2020). An efficient metaheuristic based rule optimization of apriori rare itemset mining for adverse disease diagnosis model. PalArch's Journal of Archaeology of Egypt/Egyptology, 17(7), 4763-4780.

Li, B., Zhou, H., He, J., Wang, M., Yang, Y., & Li, L. (2020). On the sentence embeddings from pre-trained language models. arXiv preprint arXiv:2011.05864.

Goyal, A., Gupta, V., & Kumar, M. (2018). Recent named entity recognition and classification techniques: a systematic review. Computer Science Review, 29, 21-43.

Maulud, D., Jacksi, K., & Ali, I. (2023). A hybrid part-of-speech tagger with annotated Kurdish corpus: advancements in POS tagging. Digital Scholarship in the Humanities, 38(4), 1604-1612.

Leoraj, A., & Jeyakarthic, M. (2023). Spotted Hyena Optimization with Deep Learning-Based Automatic Text Document Summarization Model. International Journal of Electrical and Electronics Engineering, 10(5), 153-164.